

Retrospective estimation of latent COVID-19 infections over the pandemic in US states

Rachel Lobay, Maria Jahja, Ajitesh Srivastava, Ryan J. Tibshirani, Daniel J. McDonald

Version: December 11, 2023

Abstract

The true timing and magnitude of infections from the COVID-19 pandemic are of interest to both the public and public health, but these are challenging to pin down for a variety of data-driven and methodological reasons. Nonetheless, accurate estimates of all latent infections can improve our understanding of the true size and scope of the pandemic and provide an indication of disease patterns and burden over time. In this work, we estimate the true daily incident infections for each U.S. state by deconvolving reported COVID-19 cases using estimated infection-onset-to-case-report distributions followed by a serology-based procedure to adjust for the unreported infections. We find clear variability in the timing and magnitude in the resulting estimates, indicating a differential impact of the pandemic across states and revealing a disease burden that appears earlier and more extensively than indicated by cases alone. Our findings help to better understand the impact of the pandemic in the U.S. at the state level.

1 Introduction

Reported COVID-19 cases are a staple in tracking the pandemic at varying geographic resolutions such as national, state and county levels ([Dong et al., 2020](#); [The New York Times, 2020](#); [The Washington Post, 2020](#)). Yet, for every case that is eventually reported to public health, several infections are likely to be missed. To see why, it is important to understand who’s cases are being reported and what differentiates them from the unreported cases. Refer to [Figure 1](#) for an illustration of the path of a symptomatic infection that is eventually reported to public health.

Using this figure, we can discern a number of sources of bias in the reporting pipeline. For instance, diagnostic testing mainly targets symptomatic individuals; thus, infected individuals exhibiting little to no symptoms are likely to be missed ([Centers for Disease Control and Prevention, 2022](#)). In addition, testing practices, availability, and uptake vary across space and time ([Pitzer et al., 2021](#); [European Centre for Disease Prevention and Control, 2020](#); [Hitchings et al., 2021](#)). Finally, cases provide a belated view of the pandemic’s progression because they are subject to delays due to the viral incubation period, the speed and severity of symptom onset, laboratory confirmation, test turnaround times, and submission to public health ([Pellis et al., 2021](#); [Washington State Department of Health, 2020](#)). For these reasons, reported cases are a lagging indicator of the course of the pandemic. Furthermore, they do not represent the actual number of new infections that occur on a given day, as indicated by exposure to the pathogen. Ascertaining infection onset is difficult because there was no large-scale surveillance effort in the United States that reliably tracked symptom onset, let alone infection onset.

Explaining the course of the pandemic and investigating the effects of interventions, the burden facing various subgroups, and drawing insights for future pandemics is challenging because the true spatial and temporal behaviour is unknown. While reported cases provide some understanding of the disease burden in a population, it is incomplete, delayed, and understates the true size of the pandemic. Regardless of these difficulties, it is important to the public and public health to perform a pandemic post-mortem and try to better estimate the true extent of its effect—to attempt to capture the true size and impact of the pandemic as much as we can. Estimates of daily incident infections are one such way to measure this and can guide public and professional understanding of the pandemic burden over space and time.

In this work, we provide a statistically justified reconstruction of daily incident infections for each U.S. state from March 9, 2020 to February 28, 2022. We achieve this by first breaking the task of estimating



Figure 1: Idealized chain of events from infection onset to case report date for a symptomatic infection that is eventually reported to public health.

infection onset from report date into the more manageable parts of estimating the time from infection to symptom onset and the time from symptom onset to report date (as depicted in Figure 1). Using state-level data, we construct state-time-specific incubation period and symptom-onset-to-case report delay distributions. We then use these estimated incubation period distributions in conjunction with their delay distributions to deconvolve daily reported COVID-19 cases to infection onset. The resulting infection estimates are adjusted to account for unreported infections by using seroprevalence data in a novel leaky immunity model that is defined by its ability to account for the waning of detectable immunity. We examine some features of the infection estimates and the implications of using them rather than reported cases in assessing the impact of the pandemic. We apply our infection estimates to get time-varying infection-hospitalization ratios (IHRs) for each state and compare those to similarly derived case-hospitalization ratios (CHRs). While these analyses provide a glimpse into the utility of our infection estimates, we believe that there is much more to be explored, and we hope that our work will prove an important benchmark for others to undertake retrospective analyses.

2 Methods

In what follows, we provide details on how we estimate the daily incident infections for each state over the considered time period of March 9, 2020 to February 28, 2022 and the data we used to achieve this. We start with a brief introduction to each data source used and follow this with a description of each major analysis task in the order they are performed. Figure 2 provides a visual summary of the data, analysis tasks, and the relationships between them. The five major analysis tasks this figure aims to convey are as follows: First, we estimate the incubation period and delay distribution for each day over the considered time period for a given state. Next, we join each of these two parts together using convolution to obtain a distribution from infection onset to case report for each time. We use the resulting probability estimates along with daily reported cases in retrospective deconvolution to estimate the infection onset dates for the reported cases. We adjust these infection estimates to account for the unreported infections by using state-specific, time-varying seroprevalence data in a leaky immunity model to reach our ultimate goal of obtaining daily incident infection estimates for each state.

2.1 Data

The variant-specific incubation periods are taken to be the same for all states. For the Ancestral variants, we build this using literature estimates of the gamma distribution parameters (Tindale et al., 2020). For the Alpha, Beta, Gamma, Delta and Omicron variants, we use the mean and standard deviation of the number of days of incubation as reported in Tanaka et al. (2022); Grant et al. (2022); Ogata et al. (2022).¹ Since the literature lacks reliable estimates for the incubation period of the Epsilon and Iota variants, we use the incubation period for Beta because Epsilon, Iota, and Beta are all children from the same parent in the phylogenetic tree of the Nextstrain Clades (as depicted in Hodcroft, 2021).

To estimate the daily proportions of the variants circulating in each state, we obtain the GISAID genomic

¹To clarify, we use the estimates for Alpha and Omicron from Tanaka et al. (2022), those for Beta and Gamma from Grant et al. (2022), and those for Omicron from Ogata et al. (2022).

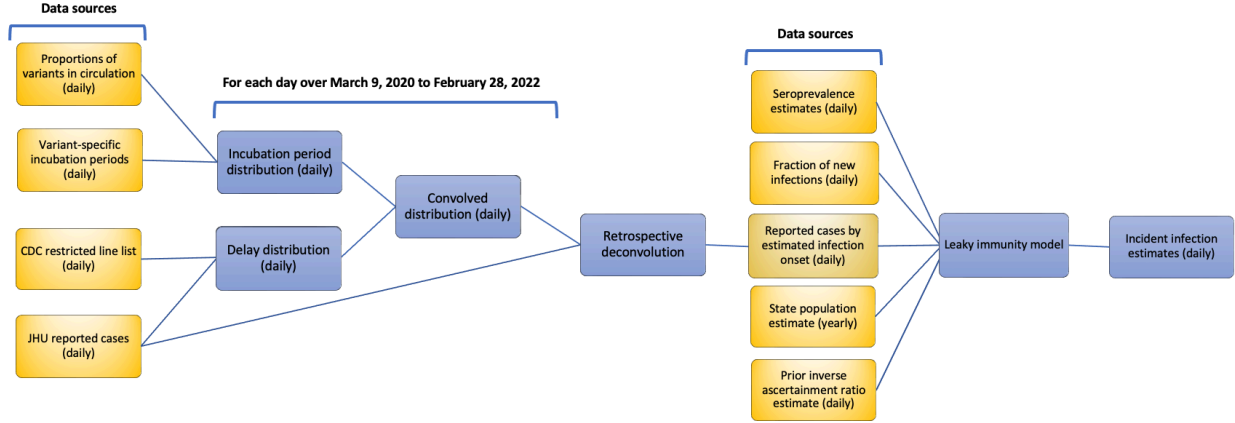


Figure 2: Flowchart of the inputted data and major analysis steps required to get from reported cases to incident infection estimates for each day over March 9, 2020 to February 28, 2022 for a state. Data sources are coloured in yellow, while data analysis steps are coloured in blue. The data sources that do not stem from an analysis step are literature estimates.

sequencing data counts from CoVariants.org (Hodcroft, 2021; Elbe and Buckland-Merrett, 2017).² Since these counts are biweekly totals, we use a simple convex optimization approach to interpolate daily numbers, where we enforce that the counts in each interval must sum to the right boundary (the biweekly total) and linear growth between the pairs of adjacent days.

The COVIDcast API (Reinhart et al., 2021) is used to retrieve the daily number of new confirmed COVID-19 cases for each state that are based on reports from the John Hopkins Center for Systems Science and Engineering (JHU CSSE, Dong et al., 2020). From the same API, we also retrieve the daily number of confirmed COVID-19 hospital admissions for each state that are collected by the U.S. Department of Health and Human Services (HHS). Both datasets are updated as of June 6, 2022.

We obtain de-identified patient-level line list data on COVID-19 cases from the CDC. Although there are both public and restricted versions of the dataset available containing the same patient records (Centers for Disease Control and Prevention, 2020a,b), the restricted dataset³ is selected because it contains information on the state of residence which is essential for constructing state-specific delay distributions. Since the restricted dataset is updated monthly and cases may undergo revision, we use a single version of it that was released on June 6, 2022. We consider this version to be finalized in that it well-beyond our study end date such that the dataset is unlikely to be subject to significant revisions.

In this dataset, the two key variables of interest are the dates of symptom onset and report to the CDC. However, we find that the line list is prone to high percentages of missing data, notably with respect to our variables of interest. Nearly 60% of cases are missing the symptom onset date, while about 9% of cases are missing the report date. In addition, we faced the fundamental issue that Jahja et al. (2022) described, in which cases with missing report dates may be filled with their symptom onset date. Figure 3 suggests that this impacts states differentially due to the inconsistent proportions of complete cases (those with both onset and report date) that have zero delay between onset and report across states. Due to this contamination in the zero delay cases (the true extent of which is unknown to us), we omit all such cases from our analysis.

For the same release date, the restricted line list contains 74,849,225 cases (rows) in total compared to 84,714,805 cases reported by the JHU CSSE; that is, line list is missing about 10 million cases. The extent that this issue impacts each state is shown in Figure 4, from which it is clear the fraction of missing cases is substantial for many states, often surpassing 50% (Jahja et al., 2022). In addition, the probability of being missing does not appear to be the same for states, so there is likely bias introduced from using the complete case line list data. We consider such bias to be unavoidable in our analysis due to a lack of alternative line

²The complete list of EPI_SET Identifiers that were used to produce the CoVariants data are provided in the Acknowledgements section of their website (Hodcroft, 2021).

³The CDC does not take responsibility for the scientific validity or accuracy of methodology, results, statistical analyses, or conclusions presented.

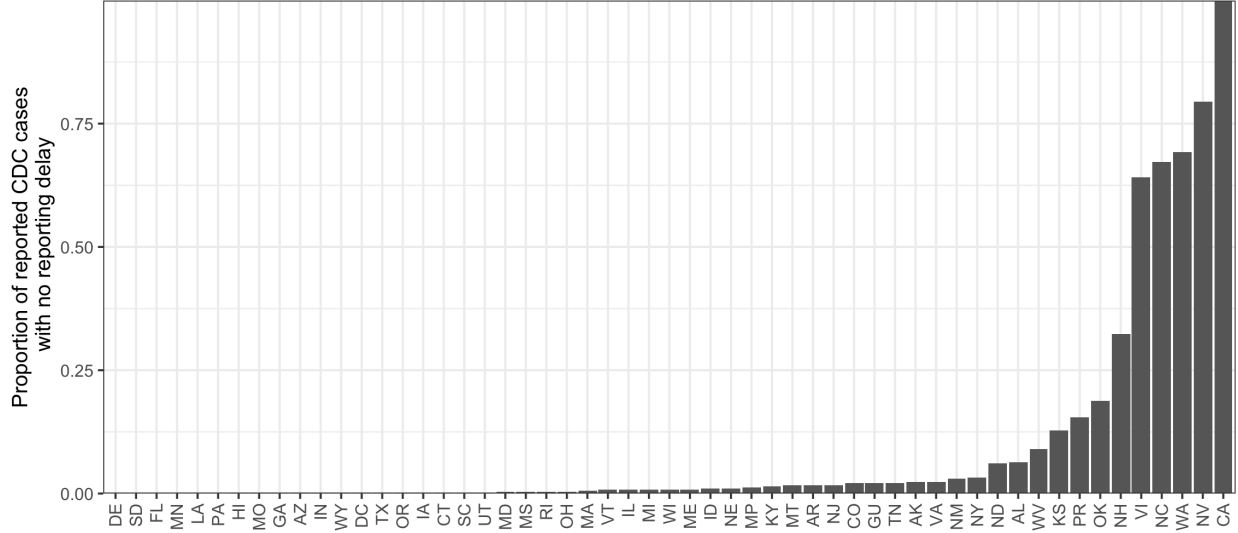


Figure 3: Proportion of complete cases with zero delay by state in the restricted CDC line list dataset.

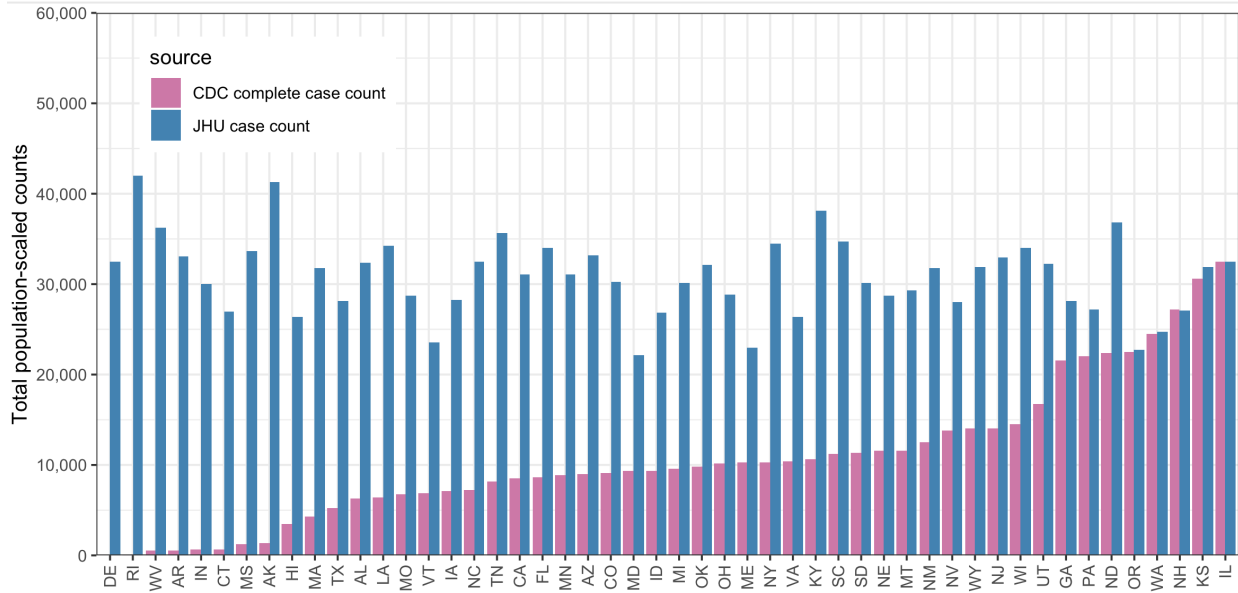


Figure 4: Complete case counts by state in the CDC line list versus the cumulative complete case counts from JHU CSSE as of June 6, 2022. All counts have been scaled by the 2022 state populations as of July 1, 2022 from [U.S. Census Bureau, Population Division \(2022\)](#).

list sources.

In the line list, we observe unusual jarring spikes in reporting in 2020 compared to 2021. Upon plotting by report date, we find that a few states are contributing unusually large case counts on isolated days very late in the reporting process (usually well beyond 50 days). We strongly suspect that these large accumulations of cases over time are due breakdowns of the reporting pipeline (which may be expected to occur more frequently in the year following its instantiation than later in time). Such anomalies are not likely to be reliable indicators of the delay from symptom onset to case report. Therefore, we devise a simple, ad hoc approach to detect and prune these reporting backlogs.

For each of the four dates of March 1, June 1, September 1, and December 1, 2020, we bin the reporting

delays occurring from 50 days up to the maximum observed delay. Then, for each bin, we obtain the total delay count for each state. We check whether each count on the log scale is at least the median (for the bin) plus 1.5 times the interquartile range and retain only those that exceed this criterion as potential candidates for pruning. Next, we compute the counts by report date for each candidate state. If there is a report date with a count greater than or equal to the pre-specified threshold, then we remove those cases from the line list. Based on inspection and intuition, we set the threshold to 2000 for the first two bins, and then lower it to 500 for the remaining bins. A similar trial and error approach is used to set the bin size (to 50 days).

To estimate the proportion of the population in each state with evidence of previous infection across time, we use two major seroprevalence surveys that were led by the CDC: the 2020–2021 Blood Donor Seroprevalence Survey and the Nationwide Commercial Lab Seroprevalence Survey ([Centers for Disease Control and Prevention, 2021a,b](#)). In the former, the CDC collaborated with 17 blood collection organizations in the largest nationwide COVID-19 seroprevalence survey to date ([Centers for Disease Control and Prevention, 2021a](#)). The blood donation samples were used to construct monthly seroprevalence estimates for nearly all states from July 2020 to December 2021 ([Jones et al., 2021](#)). In the latter survey, the CDC collaborated with two private commercial laboratories and used blood samples to test for the antibodies to the virus from people that were in for routine or clinical management (presumably unrelated to COVID-19, [Bajema et al., 2021](#)). The resulting dataset contains seroprevalence estimates for a number of multi-week collection periods starting in July 2020 to February 2022.

Both datasets are based on repeated, cross-sectional studies that aimed, at least in part, to estimate the percentage of people who were previously infected with COVID-19 using the percentage of people from a convenience sample who had antibodies against the virus ([Bajema et al., 2021](#); [Centers for Disease Control and Prevention, 2020c](#); [Jones et al., 2021](#)). Adjustments were made in both for age and sex to account for the demographic differences between the sampled and the target populations. However, both datasets are incomplete and they differ in the number and the timing of the data points for each state ([Figure 5](#)). Such limitations indicate that reliance upon only one seroprevalence survey is inadvisable. For example, in the commercial dataset, the last estimate for North Dakota is in September 2020. In the blood donor dataset, Arkansas does not have estimates available until October 2020. In addition, this blood donor dataset lacks measurements for any states in 2022 (as the corresponding survey ended in December 2021). Finally, as can be seen from [Figure 5](#), the final commercial seroprevalence measurement in 2022 shows a large increase relative to the immediately preceding measurement for each state. Since such an increase may signal unreliability or instability of the final estimates, we decided to remove them from our analysis.

The date variables that come with the two seroprevalence datasets are different and so the date variables that we are able to construct from them are not the same. For the commercial dataset, we use the midpoint of the provided specimen collection date variable. A major difference in the structure of the two datasets is that the commercial dataset always has the seroprevalence estimates at the level of the state, while the blood donor dataset can either have estimates for the state or for multiple separate regions within the state. For the blood donor dataset, we use the median donation date if the seroprevalence estimates are designated to be for entire state. If they are instead for regions in the state, since there is reliably one measurement per region per month, we aggregate the measurements into one per month per state by using a weighted average (to account for the given sample sizes of the regions). The median of the median dates is taken to be the date for the weighted average.

For adjusting our infection counts, annual estimates of the resident state populations as of July 1 of 2020, 2021, and 2022 are taken from the December 2022 press release on the U.S. Census Bureau website ([U.S. Census Bureau, Population Division, 2022](#)).

The daily fraction of new infections are estimated from the provided incidence of suspected reinfections over March 2020 to April 2022 in Clark County, which is based on surveillance work conducted by the Southern Nevada Health District (SNHD) and reported by [Ruff et al. \(2022\)](#). The proportion of new cases per week that are suspected reinfections are calculated by dividing the number of suspected reinfections by all new PCR-identified cases during the same week.

2.2 Estimating the incubation period distribution

For each state at each time over March 9, 2020 to February 28, 2022, we estimate the incubation period distribution from a finite and countable mixture of gamma distributions to account for the gradual decline in



Figure 5: A comparison of the seroprevalence estimates from the Commercial Lab Seroprevalence Survey dataset (yellow) and the 2020–2021 Blood Donor Seroprevalence Survey dataset (blue). Note that the maximum and the minimum of the line ranges are the provided 95% confidence interval bounds to give a rough indication of uncertainty.

the incubation period across variants. The variants we consider are Ancestral (for our purposes, all 2020 variants observed in the U.S.), as well as Alpha, Beta, Gamma, Delta and Omicron, which are designated as variants of concern by WHO based on their potential to cause new waves, dethrone the dominant variant, and lead to changes in public health policy ([World Health Organization, 2021](#)). In addition, we include the Epsilon (California) and Iota (New York) variants because of their impact on those and the surrounding states ([Yang et al., 2022](#); [Duerr et al., 2021](#)). We relegate all other variants to be in an Other category (so that the proportions circulating in a state at a time always sum to one). This decision is, in part, motivated by the lack of sequencing data for most states in 2020 as well as the presence of an others category in the sequencing data for that time. Then, for each variant in a state at a time, the proportion of the variant circulating (the mixture weight) is multiplied by the corresponding component gamma distribution for the incubation period. These distributions are the same for all states and based on literature estimates of the gamma parameters or the mean and standard deviation of the incubation period (in which case the method of moments is used to fit a gamma density). Finally, we discretized the resulting mixture density to the support

set, which is taken to be from 1 and 21 days. In other words, those are taken to be the lower and upper limits for the number of days that the virus could be incubating in someone. The implicit assumption for the lower bound is that there must be at least one day between infection and symptom onset (which follows the convention given in [Public Health Agency of Canada, 2021](#)). The assumption underlying the upper bound is that 21 days is the maximum number of days that the virus could be incubating in someone (which is reasonable based on [Zaki and Mohamed, 2021](#) and [Cortés Martínez et al., 2022](#)).

2.3 Estimating the delay from symptom onset to report date

We use the restricted CDC line list to estimate the distribution between symptom onset and report for each state at each time. We refer to this as the “delay distribution” (see [Figure 2](#) or [Jahja et al., 2022](#)). More formally, let y_t denote the count of new cases reported at time t and x_t denote the count of new infections with onset at t for a state. For all cases in the line list that had both an onset and a report date, we can count the those that are reported at time t by enumerating them according to onset (as in [Jahja et al., 2022](#)):

$$y_t = \sum_{s=1}^t \sum_{i=1}^{x_s} \mathbf{1}(\text{the } i^{\text{th}} \text{ infection at } s \text{ gets reported at } t).$$

Taking the conditional expectation of the above yields

$$\mathbb{E}(y_t \mid x_s, s \leq t) = \sum_{s=1}^t \pi_t(s) x_s,$$

where $\pi_t(s) = \mathbb{P}(\text{case report at } t \mid \text{infection onset at } s)$ for each $s \leq t$ are the delay probabilities and the $\{\pi_t(s) : s \leq t\}$ sequence comprises the delay distribution at time t . Notice that there are no time restrictions placed on the infection onset, save that it must have been between the start of the pandemic and the report date, inclusive. This is unlikely to be a realistic assumption to make as t moves farther away from s .

Thus, we make two key assumptions about the delay distributions. First, infections that are reported to the CDC are always reported within $d = 60$ days, which is true for the vast majority of reported cases. Second, the probability of zero delay is zero, which stems from the contamination of zero-delay in the line list. As in [Jahja et al. \(2022\)](#), we update the conditional expectation formula to reflect these two assumptions:

$$\mathbb{E}(y_t \mid x_s, s \leq t) = \sum_{k=1}^{60} p_t(k) x_{t-k}$$

where for $k = 1, \dots, 60$,

$$p_t(k) = \mathbb{P}(\text{case report at } t \mid \text{infection onset at } t - k).$$

For each state, we estimate the delay distribution at each t by using the empirical distribution of all non-zero lags between the complete cases whose onset dates fall in the center-aligned interval about t designated by $[t - 75 + 1, t + 60]$.

Now, the task of estimating the delay distribution for each state at each time requires four distinct steps. First, we obtain the empirical distribution of all lags (excluding zero) from all cases with onset dates falling in the center-aligned interval. Next, we weight the state-specific empirical distribution by the proportion of CDC cases to JHU cases. That is, we compare the number of CDC cases used to create the empirical distribution to the number of cases reported by JHU in the time window of $[t - 60 + 2, t + 75]$ (to correspond appropriately to the center aligned interval for the CDC cases). This proportion is used as the weight for the state’s empirical distribution, while the complement is used to weight the overall empirical distribution that is formed from the data for all states. This construction allows for more reliance on the state’s distribution when there are more CDC cases relative to JHU (and vice versa). After implementing the shrinkage method, we fit a gamma density to the resulting empirical distribution by the method of moments. Finally, we discretize the resulting density to the support set of 1 to 60 days.

2.4 Convolutional estimate of infection-to-report distributions

From the incubation period and delay distribution estimation, we acquire one delay and one incubation period distribution for each state at each time under consideration. We then convolve each pair of distributions to get the estimated infection-to-report distributions and, hence, the estimated probabilities for the delay from infection onset to case report.

2.5 Retrospective deconvolution for reported cases

The goal for retrospective deconvolution is to estimate the daily number of new infections that occurred at each time using the dates that those cases were eventually reported. For each state, we achieve this goal by solving the an optimization problem. Let \mathcal{T} represent the deconvolution period from June 1, 2020 to February 28, 2022. Let \hat{p}_t be probabilities from the estimated infection-to-report distribution for $t \in \mathcal{T}$, y_t the number of new cases reported, and $D^{(4)}x$ yields all 4th-order differences of the vector x (by using the discrete derivative matrix of order 4, $D^{(4)}$). From these, we estimate the latent infection counts for the reported cases across time by solving for the vector x in

$$\underset{x}{\text{minimize}} \sum_{t \in \mathcal{T}} \left(y_t - \sum_{k=1}^d \hat{p}_t(k) x_{t-k} \right)^2 + \lambda \|D^{(4)}x\|_1.$$

The above loss function decouples into two parts which trade data fidelity with desired smoothness (that encapsulate the classic bias-variance trade off). The first part represents minimizing the sum of squared errors between the JHU reported cases and the estimates, while the second part captures the smoothness of the estimates (smaller values being more smooth). The tuning parameter λ determines the relative importance of these competing goals.

We solve this trend-filtering-regularized least squares deconvolution problem by employing the ADMM algorithm from [Ramdas and Tibshirani \(2016\)](#) that is described in Appendix A of [Jahja et al. \(2022\)](#). The solution to the problem is an adaptive piecewise cubic polynomial ([Tibshirani, 2014, 2022](#)).

We select the tuning parameter, λ , by using 3-fold cross validation as in [Jahja et al. \(2022\)](#) in which every third infection count is reserved for testing and imputed with the average of the two surrounding counts. The tuning parameter that results in the smallest sum of squared errors is ultimately chosen.

2.6 Inverse reporting ratio and the leaky immunity model

The infection estimates from retrospective deconvolution are derived solely from the infection onset dates of the reported cases. To capture the unreported infections, it is necessary to adjust the estimates by a scaling factor that approximates the ratio of the true number of new infections to the new reported infections. We refer to this quantity as the inverse reporting ratio and denote it by a_t for time t . Our new goal is to estimate this quantity for every state at every time under consideration.

The number of new reported infections is obtained from our deconvolved case estimates. As for the true infections, since seroprevalence of anti-nucleocapsid antibodies is used to estimate the percentage of people who have at least one resolving or past infection ([Centers for Disease Control and Prevention, 2020c](#)), we can use the change in subsequent seroprevalence measurements to capture new infections, accounting for those who lose enough immunity to fall below the detection threshold. For each state, let s_t be a seroprevalence estimate at time t , w_t be the inverse variance weights corresponding to those estimates, and ΔR_t be the change in cumulative reported infections scaled by the state’s population. To account for reinfections, we multiply the change in reported infections at time t by the corresponding fraction of new infections, n_t . Using these components, we construct a model separately for each state

$$s_t = (1 - \gamma)s_{t-1} + a_t \Delta R_t n_t + \epsilon_t \tag{1}$$

where $\epsilon \sim N(0, w_t \sigma_\epsilon^2)$, and γ is the percentage of people who lose immunity between time t and time $t + 1$. Informally, we refer to γ as the leaky parameter and we call this model leaky immunity model. By “leaky” we mean the decrease in detectability of antibodies due to the natural degradation of infection-induced immunity over time. Since the true course of immunity over time is unknown ([Goldberg et al., 2022](#)), we take this

rather straightforward approach and simply model a singular γ to try avoid making gratuitous or overly restrictive assumptions.

To estimate this model, we will express it as a state-space model with weekly observations and enforce a_t to follow a second-order autoregressive model. This representation allows for convenient handling of missing data, extrapolation before and after the period of observed seroprevalence measurements, and maximum likelihood estimates of γ and σ_ϵ^2 . Details of this methodology and the computation of the associated uncertainty measurements are deferred to the Appendix.

2.7 Lagged correlation to hospitalizations and time-varying IHRs

We use our infection estimates in a lagged correlation analysis with confirmed COVID-19 hospitalizations. Our primary goal of this analysis is to find the lag between infection and hospitalization rates that gives the highest average rank-based correlation across US states. To that end, we consider a wide range of possible lag values ranging from 1 to 25 days. Zero and negative lags are not considered because COVID-19 infection onset must precede hospitalization due to the virus. To remove day of the week effects, both the infection and hospitalization signals are subject to a 7-day moving average (center-aligned) before their conversion to rates.

For each considered lag, we calculate the Spearman’s correlation between the state infection and hospitalization rates for each observed day over the March 9, 2020 to February 28, 2022 time period with a center-aligned rolling window of 61 days for each such computation. We then calculate the average correlation across all states and times for each lag. The lag that leads to the highest average correlation is used to estimate the time-varying IHRs for each state. To compute this for a given day, the number of individuals who are hospitalized due to COVID-19 on a day are divided by the estimated total number who were infected on the lagged number of days before.

2.8 Ablation study for the lagged correlation analysis

To better understand the contribution of the intermediate steps to the lagged correlation analysis, we carry out a brief ablation study in which we calculate the lagged correlation using the following infection estimates: 1. those from the deconvolution procedure under the assumption that the report date is the same as the symptom onset date (i.e., excluding the reporting delay data); 2. those from the deconvolution procedure when only the reporting delay data is used (excluding the incubation period data); 3. those from the deconvolution procedure when utilizing both the incubation period and delay data (the deconvolved case estimates); 4. those from applying the leaky model to produce estimates for both the reported and the unreported cases (the infection estimates).

3 Results

This work estimates incident infections for each U.S. state over March 9, 2020 to February 28, 2022 and to illustrate the disease burden and viral transmission dynamics at the level of the state across time. After converting the number of infections to rates (infections per 100,000 population), we perform a brief comparison between infection and case estimates within each state to see to what extent that surges in infections are evident in cases alone and point out instances where cases largely fail to capture surges in infections. Then, we look at patterns in infections across the states and postulate what may be contributing to these trends based on defining characteristics such as state healthcare performance and geographical contiguity.

3.1 Infection estimates compared to reported cases

Naturally, outbreaks in infections precipitate those in cases and are reliably larger in magnitude (Figure 6 and Figure 7). Hence, our infection estimates indicate that the pandemic had a differential impact across states earlier and at a larger scale than is suggested by cases. While the major Ancestral, Delta, and Omicron waves tend to be visible for most states, there are clear outbreaks in unreported infections that are not easily detectable from cases alone in the falls of 2020 and 2021. For example, take the Alpha wave from mid-2021 in the midwestern states of Michigan, North Dakota, South Dakota, and Illinois or the Delta wave from the fall of 2021 in the east coast states of Connecticut, Rhode Island, and Massachusetts. Earlier on in the

pandemic, such discrepancies may be more attributable to failures in the reporting pipeline, while later on in the pandemic, they more likely due to the rise in asymptomatic infections across variants ([Ontario Agency for Health Protection and Promotion, 2022](#); [Garrett et al., 2022](#)).

Finally, while the January 2022 Omicron wave is evident from the case counts for all states ([Figure 7](#)), our estimates suggest that case counts tend to severely underestimate infections during this time for many states. The lowest of all states was in Nevada, where about 9.1% (95% confidence interval: [4.0, 100.0]) of the infections were reported over January 2022. This was followed by Arizona with 12.2% ([4.8, 98.4]), and Minnesota with 13.9% ([8.0, 53.5]). More broadly, in 25 states less than 50% of infections were reported. Only 11 states of Florida, New York, Colorado, New Jersey, South Dakota, Delaware, Alabama, Arkansas, North Carolina, New Hampshire, and Maryland reported at least 70% infections over this time.

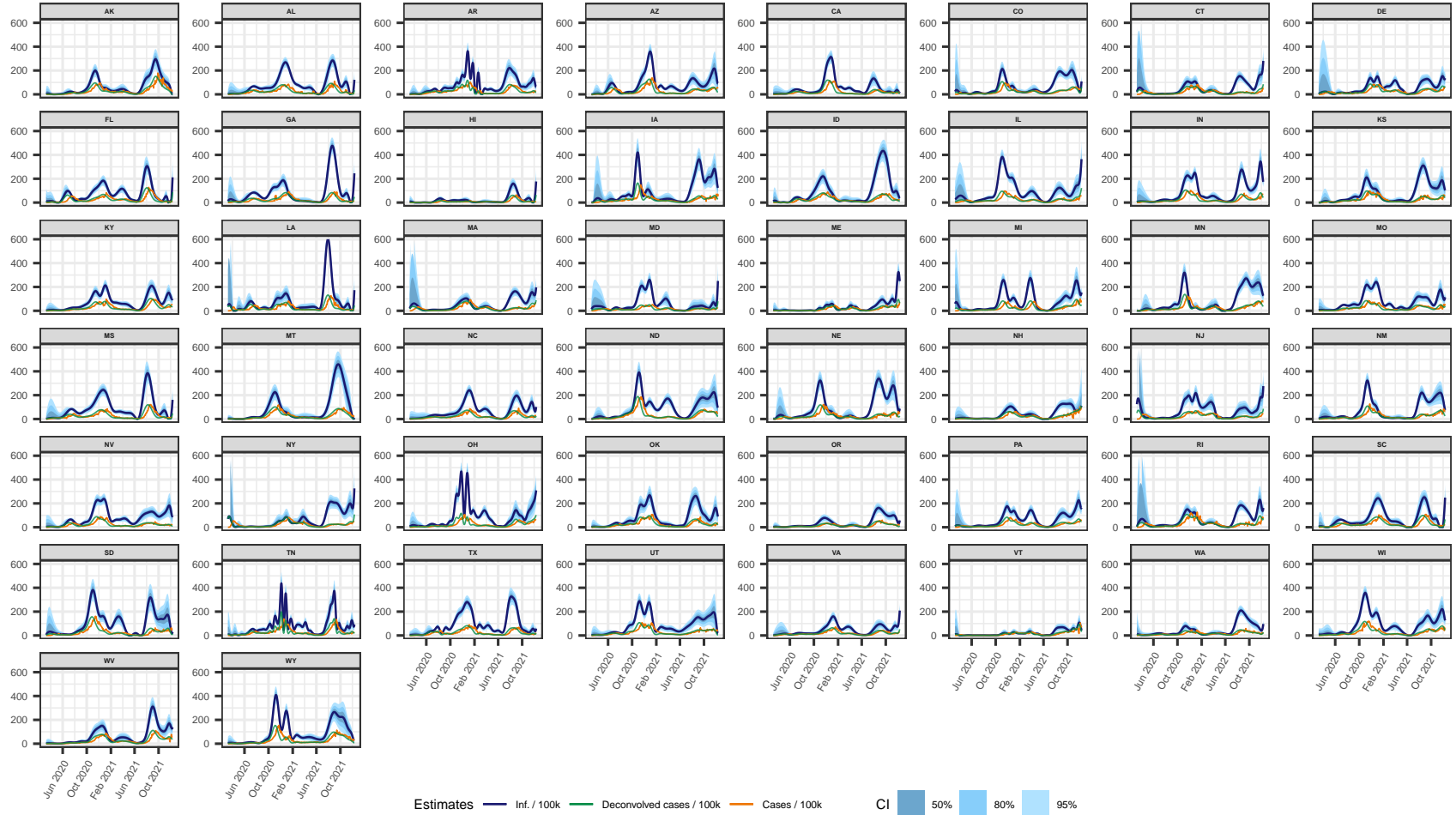


Figure 6: Estimates of the number of daily new infections per 100,000 population for each US state from March 9, 2020 to December 11, 2021 (dark blue line). The blue shaded regions depict the 50, 80, and 95% confidence intervals for the estimates, while the teal line represents the number of new daily new deconvolved cases per 100,000, and the dotted orange line represents the 7-day average of the new cases per 100,000 as of the same date.

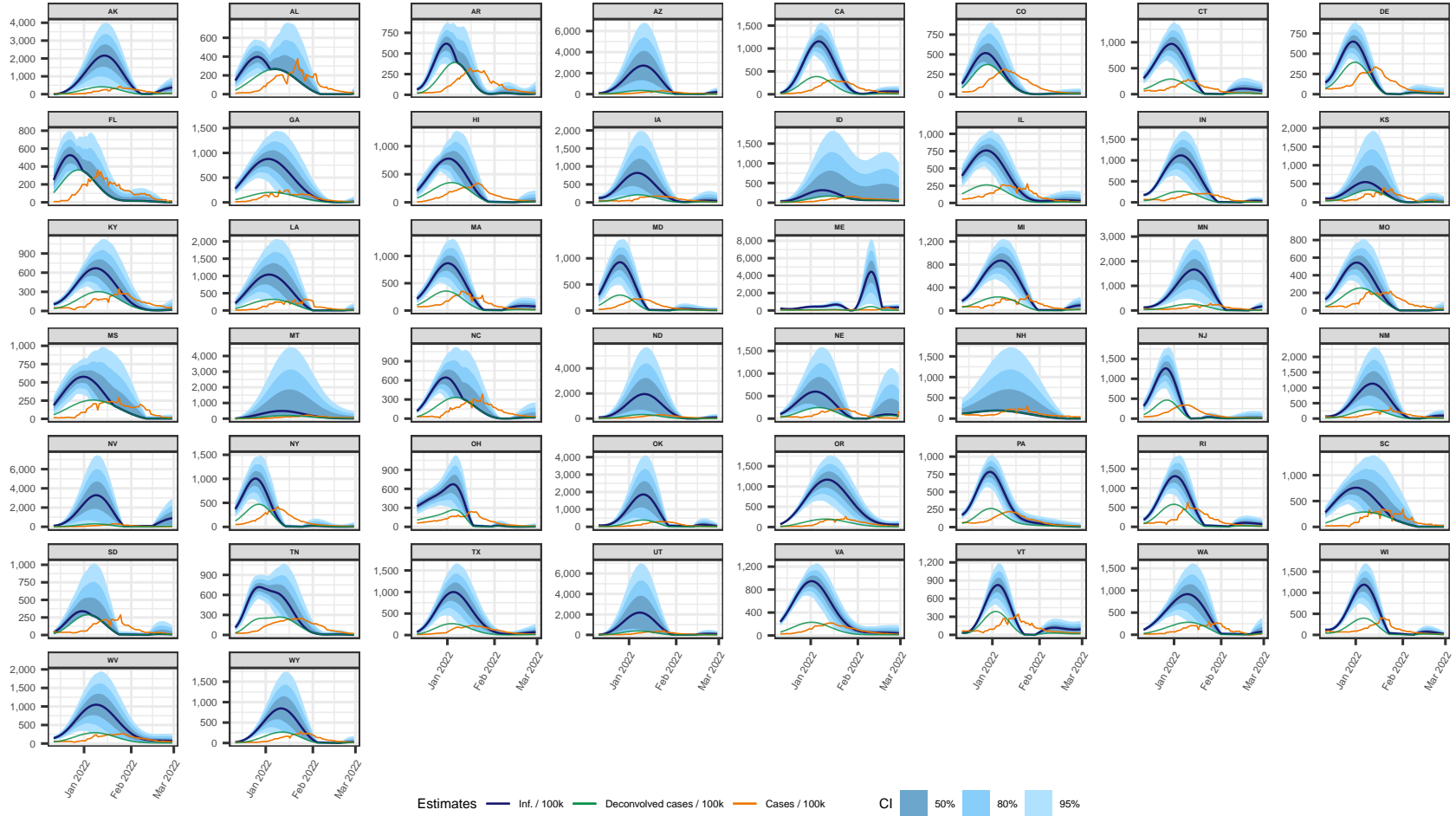


Figure 7: Estimates of the number of daily new infections per 100,000 population for each US state from December 11, 2021 to February 28, 2022 (dark blue line). The blue shaded regions depict the 50, 80, and 95% confidence intervals for the estimates, while the teal line represents the number of new daily new deconvolved cases per 100,000, and the dotted orange line represents the 7-day average of the new cases per 100,000 as of the same date.

We perform a lagged correlation analysis where we systematically investigate the rank-based (i.e., Spearman’s) correlation between our infection and confirmed hospitalization rates per 100,000 population over a broad range of lag values (Figure 8). Examining the correlation between infections and hospitalizations shows that the maximum average correlation across states of 0.651 is observed at a lag of 14 days. In contrast, we find that the greatest average rank-based correlations for cases with confirmed hospitalizations is achieved at a lag of 1. That is, we find that case report rates are nearly contemporaneous to hospitalizations, while infection estimates clearly precede them.

We undertake an ablation study for the lagged correlation of infections, the results of which are shown in Figure 9. From this, we can see that the infection estimates from each intermediate step are all leading indicators of hospitalizations. However, the degree that each such set of estimates lead hospitalizations depend on its location in the sequence of steps and how close the estimates are to infection onset. For example, the deconvolved cases by symptom onset tend to precede hospitalizations by about 11 days, while those for the subsequent step indicate that the deconvolved case estimates by infection onset precede hospitalizations by about 15 days. When we solely rely on incubation period data in the deconvolution (i.e., under the assumption that the report date is the same as the date of symptom onset), we can see that the reported infections tend to precede hospitalizations by about 6 days.

In terms of average correlation, the deconvolved case estimates by infection onset provide similar information indicative to hospitalizations as the deconvolved case estimates by symptom onset about 4 days before the latter tends to occur, which highlights a time-based benefit of opting for infection estimates by the date of infection onset over symptom onset.

Unsurprisingly, the deconvolved case and infection estimates achieve their maximum correlation at nearly the same lag. And yet, the average correlation to hospitalizations tends to be greater for the deconvolved case estimates than for the infection estimates (and the reported infections by symptom onset). If the goal is to find the part of this process that is most informative to hospitalizations, then it is clear that producing the deconvolved case estimates is the more informative step and these estimates are potentially the most meaningful signal for future hospitalizations. This finding could indicate that the unreported infections tend to be less severe and less likely to lead to hospitalization than those that are reported.

As a counterpart to our lagged correlation analysis, we compute the time-varying IHRs for each state using the optimal lag for infection and hospitalization rates. We also included the CHRs that are computed using the optimal lag for cases and hospitalizations for comparison. (Figure 10).

While the IHRs tend to be less than 0.1 hospitalizations per infection, the CHRs tend to present a more amplified version of the CHRs for each state. This supports our claim that the reported infections are more likely to require hospitalization than the unreported infections. Both the IHRs and CHRs exhibit similar geospatial and temporal trends as are noted for infections. Namely, states that are close in proximity (such as Pennsylvania and Virginia) tend to exhibit similar patterns in the IHRs over time. In addition, there are similar spikes observed across many states during waves of infections that are driven by prominent new variants. For example, many states exhibit a striking spike in hospitalizations in mid-2021, which coincides with the rapid takeover of the Delta variant during that time (Hodcroft, 2021). This finding aligns with previous studies that found an increased risk in hospitalizations with Delta in comparison to other variants (Twohig et al., 2022; Nyberg et al., 2022). Similarly, during the beginning of the wave driven by the new Omicron variant in early 2022, there tends to be another spike in the IHRs that rivals that observed during the time of Delta. However, this could be in part due to estimating near to the boundary. Aside from this, there appears to be a slight waning pattern over time for many states where the IHR is generally greater earlier than later in the pandemic over which time the virus mutates to variants that are generally more infectious, but that pose less of a risk to hospitalization (Lorenzo-Redondo et al., 2022; Blauer, 2022a).

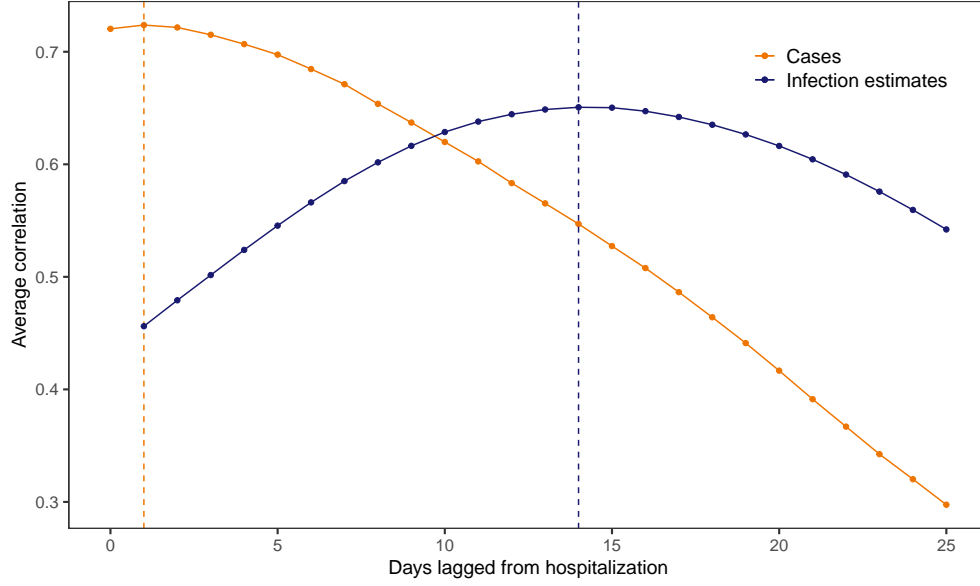


Figure 8: Lagged Spearman's correlation between infection and hospitalization rates per 100,000 as well as between case and hospitalization rates per 100,000. The averages shown are for each lag, across US states and days over March 9, 2020 to February 28, 2022, and taken over a rolling window of 61 days. Note that the infections, cases, and hospitalization counts are subject to a center-aligned 7-day averaging to remove spurious day of the week effects. The dashed lines indicate the lags for which the highest average correlation is attained.

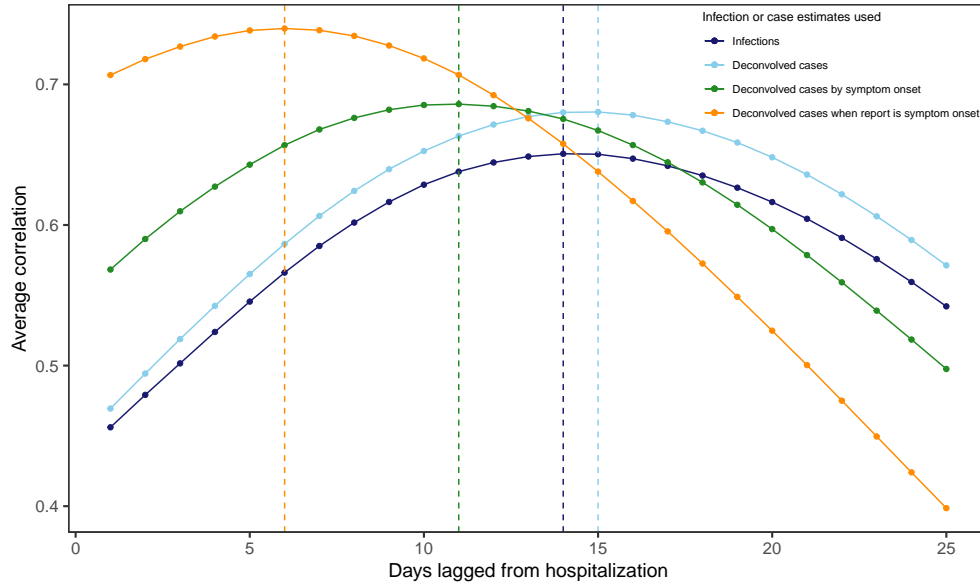


Figure 9: Lagged Spearman's correlation between the infection and hospitalization rates per 100,000 averaged for each lag across US states and days over March 9, 2020 to February 28, 2022, and taken over a rolling window of 61 days. The infection rates are based on the counts for the deconvolved case and infection estimates as well as the reported infections by symptom onset and when the report is symptom onset. Note that each such set of infection counts is subject to a center-aligned 7-day averaging to remove spurious day of the week effects. The dashed lines indicate the lags for which the highest average correlation is attained.

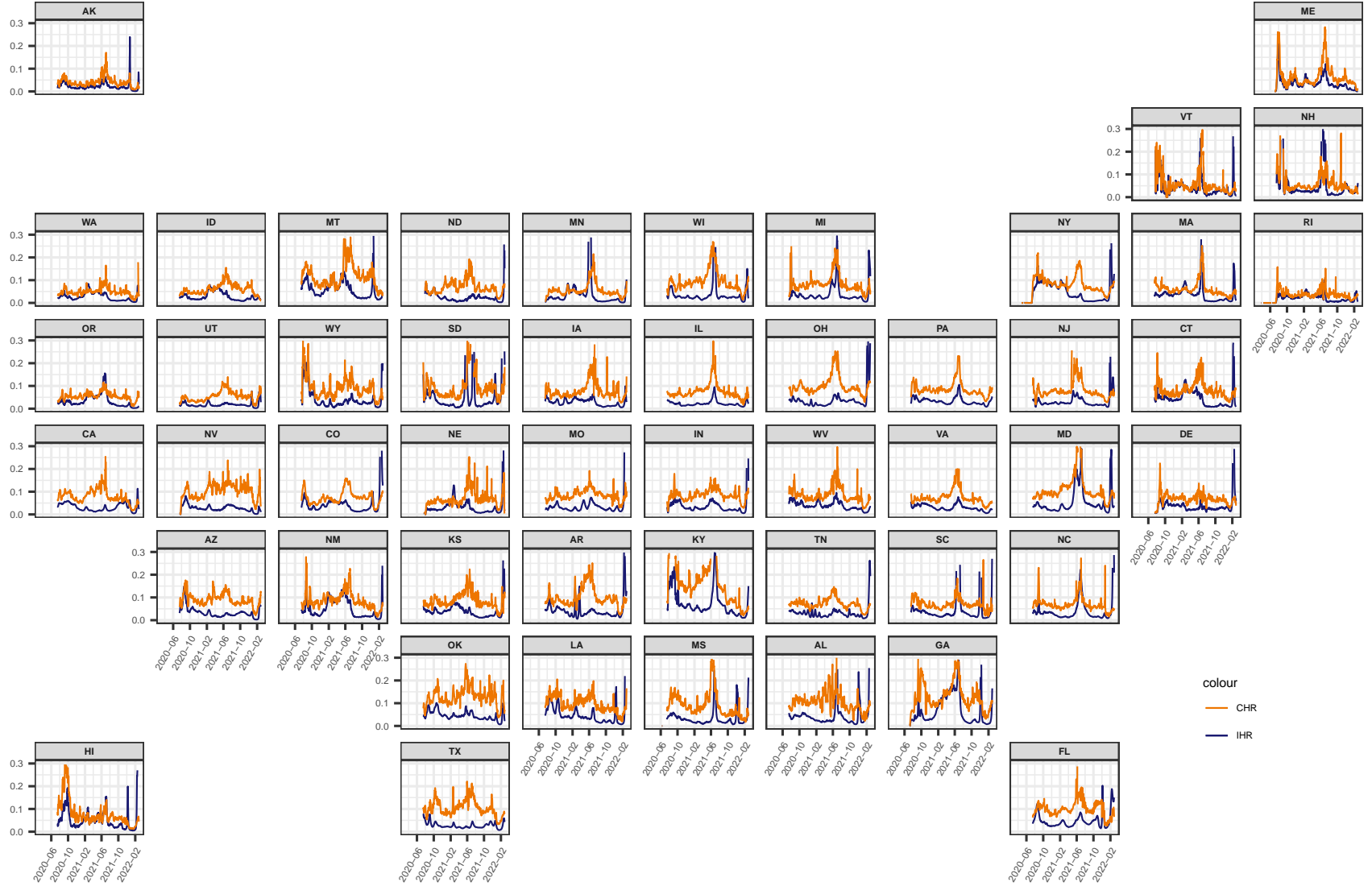


Figure 10: Time-varying IHR and CHR estimates for each state from March 9, 2020 to February 28, 2022, obtained using the corresponding optimal lag from the systematic lag analysis. Note that the infection, case, and hospitalization counts are subject to a center-aligned 7-day average to remove spurious day of the week effects. Also note that the different starting points across states are due to the availability of the hospitalization data.

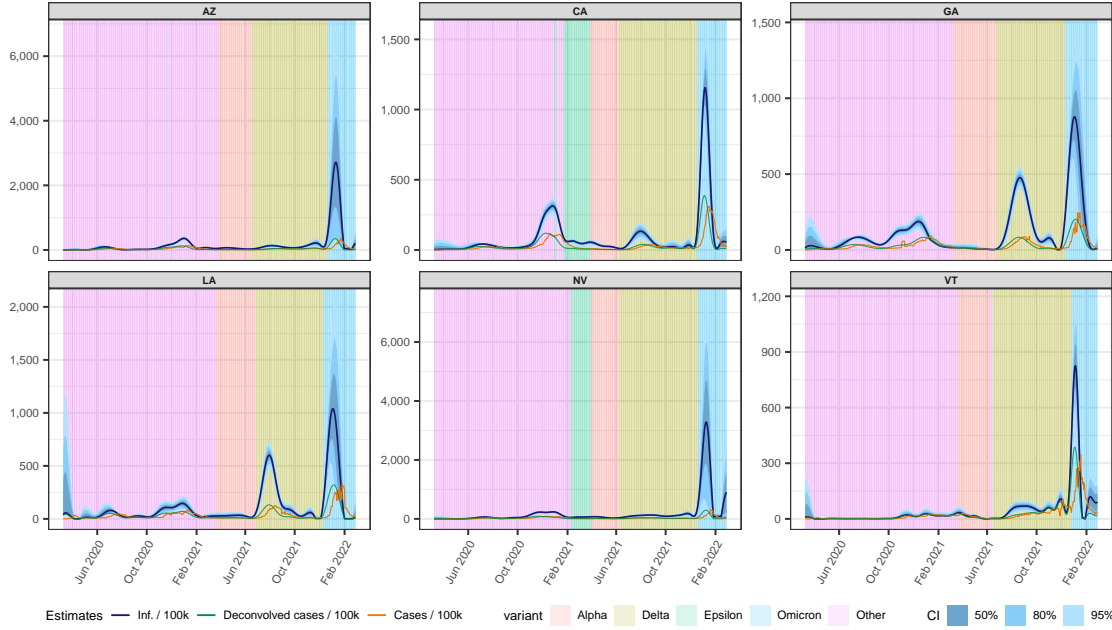


Figure 11: Estimates of the number of daily new infections per 100,000 for a sample of six US states from March 9, 2020 to February 28, 2022 (dark blue line). The blue shaded regions depict the 50, 80, and 95% confidence intervals. The background is shaded to indicate the top variant in circulation that is at least 60% of all variants in circulation at the time.

3.2 Disease burden and viral transmission

By reconstructing the time series of COVID-19 infections per 100,000 population for each US state from March 9, 2020 to February 28, 2022, we observe rates of infections that vary in intensity and disease burden across space and time (Figure 6, Figure 7, Figure 11, Figure 12). The largest observed outbreaks at the beginning of the Omicron era (January 2022) in Nevada, Arizona, and Utah which suggests a similar spread of the virus in states that are in close geographic proximity. During this time, the state that has the highest rate of infections per 100,000 on single day is Nevada with about 3289 infections per 100,000 on January 9, 2022 (95% confidence interval: [296, 7431]), followed by Arizona with 2718 on January 10, 2022 (95% confidence interval: [346, 6748]), and Utah with 2179 on January 8, 2022 (95% confidence interval: [383, 6955]). Interestingly, Maine has a delayed major spike in mid-February 2022 relative to the other states (which tend to present their major Omicron spikes in January 2022). On February 10, 2022, Maine presents the highest estimated infection rate on a day, having about 4464 infections per 100,000 (95% confidence interval: [718, 8211]).

Aside from Omicron, most states present smaller surges in infections from around November 2020 to January 2021, the late summer to fall of 2021 and leading up to the major surge in early 2022, which likely represent well-known waves driven by the Ancestral and Delta variants.

During the Delta wave in the summer of 2021, the state that has the highest rate of infections per 100,000 on single day is Louisiana with about 600 infections per 100,000 on July 30, 2021 (95% confidence interval: [474, 728]), followed by Georgia with 476 on August 8, 2021 (95% confidence interval: [407, 546]). These also represent the highest rates of infections per 100,000 for any state prior to the domination of the Omicron variants. Similar patterns in the major surges of infections are observed in nearly all states, though to varying degrees. In general, greater similarities in the strength and magnitude of outbreaks are observed in the clusters of states that border each other and in those that have similar state health system performance (in terms of indicators for the access, cost and the quality of the state's healthcare as well as the health-related outcomes as noted in Radley et al. (2020)).

The period of lowest viral transmission is observed in the spring of 2020. In April 2020, the states of

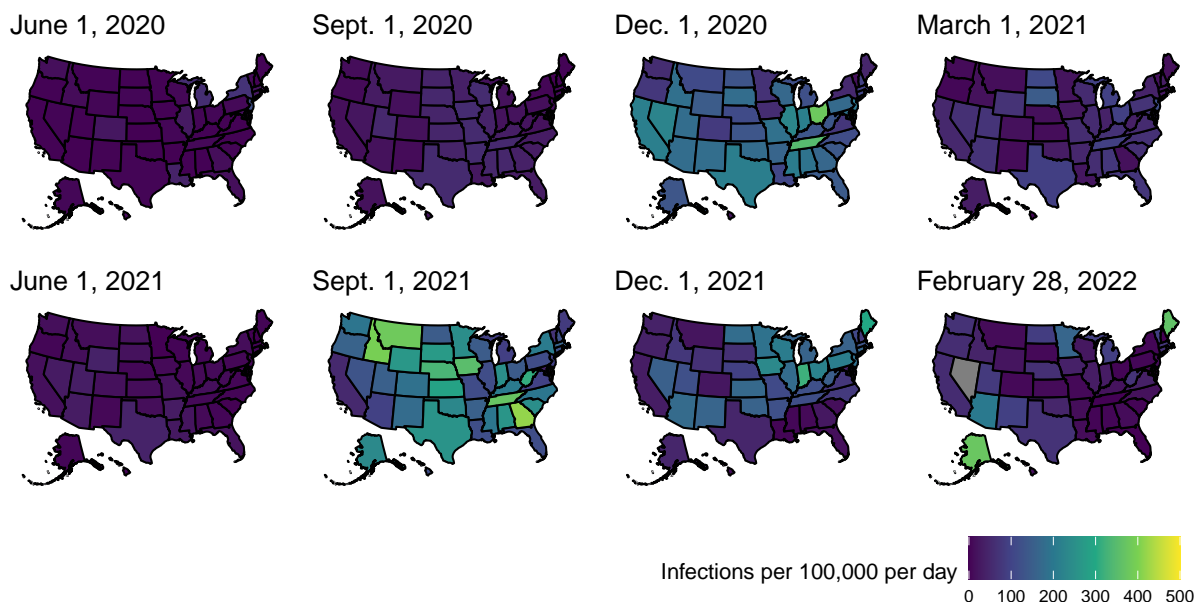


Figure 12: Choropleth maps of the state-level estimates of the number of daily new infections per 100,000 population for various times over March 9, 2020 to February 28, 2022. These maps are generated from the `usmap` package in R ([Di Lorenzo, 2023](#)).

Vermont and Hawaii achieve the lowest rate of infections over the month of 5.67 and 6.72 infections per 100,000, respectively. These are followed by Montana which achieves a rate of 10.6 infections per 100,000 in May 2020. In the spring of 2020, Montana maintains a rate under 10 infections per 100,000 from the week of April 12, 2020 to May 24, 2020. The states that consistently achieve the lowest rates of infections tend to be those that demonstrate better pre-pandemic healthcare performance such as Vermont, New Hampshire, and Hawaii ([Radley et al., 2020](#)).

From a brief inspection of the geo-contiguous states, we can observe similar patterns in surges and periods of waning over time, suggesting that states who share similarities in climate and topography performed similarly to each other. More precisely, we can observe neighboring states such as New Hampshire and Massachusetts or Washington and Oregon that present waves that mirror each other in amplitude and timing.

Interestingly, the two states that are geographically removed from the contiguous United States, Alaska and Hawaii, tend to perform quite differently from each other later in the pandemic. Alaska generally presents significantly greater rates of infections than Hawaii especially during the Omicron era. This suggests that it is not so much the non-contiguity aspect as it is other distinguishing factors that lead to lower infection rates.

4 Discussion

We obtained retrospective estimates of daily incident infections for each US state for March 9, 2020 to February 28, 2022. While all states present waves that are associated with the major emerging variant of the time, the clear variability in the magnitude of our estimates indicate that the intensity and disease burden are heterogenous across states. Yet, there are similar epidemic patterns in surges and periods of waning observed in clusters of neighbouring states. As well, states with lower reputed healthcare performance tend toward higher infection rates during major surges in infections. For instance, Nevada and Arizona, which exhibited the highest rates of infections during the major Omicron outbreak in January 2022, are both firmly located in the bottom quartile in the assessment for access and affordability from the 2020 Scorecard on State Health System Performance ([Radley et al., 2020](#)) and in the government estimates of per capita personal

health care spending by state for that year ([Centers for Medicare & Medicaid Services, 2020](#)). So perhaps the apparent difficulty in containing the spread during the Omicron wave is at least in part due to the lack of healthcare support. As noted in the report, both such states have a notably higher proportion of people who are uninsured and adults without a usual source of care. For example, the well-documented shortage of healthcare professionals impacting these states ([Do et al., 2023](#); [Gong et al., 2019](#)). Indeed, the proportion of adults without a usual source of care is estimated to be well above the national average in both states for the pre-pandemic year of 2018 ([Radley et al., 2020](#)). In addition, a lack of healthcare professionals means a lack of medical guidance at the level of the individual and a lack of opportunity for early testing/detection and treatment during the pandemic. This coupled with the typically milder symptoms of Omicron in comparison to previous variants may have made it more difficult for individuals to identify that they were infected with the virus and contributed to the increase in spread in those states. Furthermore, a lack of health insurance and the prospect of having to pay costly out-of-pocket medical expenses may have deterred people from seeking testing ([Embrett et al., 2022](#)) and denied opportunities to get tested when visiting a medical professional for reasons unrelated to COVID-19. In contrast, Hawaii, Vermont and New Hampshire, which are reported to be top performers in healthcare for the domains of access and affordability, prevention and treatment, avoidable use and cost, healthy lives, and income disparity in [Radley et al. \(2020\)](#), exhibit some of the lowest rates of infections during the pandemic and were routinely subject to surges that were lower in intensity.

Interestingly, the states with the highest rate of infections during the January 2022 Omicron wave, Utah, Arizona, and Nevada, are all in the top quartile for the percent change in population from 2020 to 2022 based on the annual estimates of the resident state populations from the US Census Bureau ([U.S. Census Bureau, Population Division, 2022](#)), suggesting that there may be a connection between states that are among the fastest growing and increased infection rates during that outbreak. Considering the larger time frame of 2010 to 2020 leading up to the pandemic, the US Census Bureau reports that Utah holds the greatest percent change in population at about 18.4%, Nevada is ranked fifth at 15%, and Arizona is ranked 9th at 11.9% ([United States Census Bureau, 2021](#)). It is reasonable that the faster growing states place more strain on the healthcare system than states that grow at a slower pace, limiting access to medical professionals and to resources for diagnosing and testing. To elucidate this possible connection, further exploration into the impact of state population growth on infection waves and healthcare system performance during the pandemic is warranted.

It is reasonable that Montana experienced some of the lowest infection rates in the spring of 2020 due to having a lower population and population density as well as from having a larger area relative to other states. This combination of favourable conditions likely contributed to lower infection rates earlier on in the pandemic. Similarly, Hawaii and Alaska were both states that maintained rates of under 10 infections per 100,000 for at least three weeks in the spring of 2020. It is plausible their geographic isolation in addition to the previous factors discussed for Montana contributed to such lower infections near to the beginning of the pandemic. Research by [Provenzano et al. \(2020\)](#) and [Carozzi et al. \(2022\)](#) on urban density and COVID-19 in US counties found that the geographic connectivity and social connectedness of denser areas likely impacted the timing of outbreaks (so that denser locations were more likely to have outbreaks earlier on), but by the end of 2020, density had little to no impact on time-adjusted COVID-19 cases. However, investigation beyond the first year of the pandemic is warranted. As well, this study only looks relation between density and cases, not infections. Since the relationship between density and infections has remained relatively unexplored, further research should be done to elucidate such connections and how they change over time.

Our infection estimates suggest that the pandemic has an impact in states earlier and at a larger scale than is indicated by cases. Since case reporting is not consistent across time and states, case counts underestimate the true number of infections and, hence, the impact of the pandemic ([Centers for Disease Control and Prevention, 2022](#); [Simon, 2021](#)). For example, some states report the number of individuals tested rather than the numbers of tests performed ([Schechtman, 2020](#); [Chitwood et al., 2022](#)).

We observe outbreaks in infections that are difficult to detect from cases alone such as the Delta wave that prevailed from July to November 2021 in Connecticut, Rhode Island, Massachusetts, and New Jersey. This suggests that cases paint an incomplete picture of the pandemic, especially when outbreaks are largely driven by unreported infections. Furthermore, since case report dates follow symptom and infection onset, cases are a fundamentally flawed indicator of disease burden because they have a built-in temporal bias. This is in addition to other biases from differences in reporting across states (such as temporary bottlenecks due to influxes of data or more persistent processing issues that increase the average time from case detection to

report ([Washington State Department of Health, 2020](#); [Dunkel, 2020](#)). So while reported cases provide an indication of the trajectory of the pandemic, it is a delayed and incomplete version. Estimating the new number of infections by symptom or infection onset date would more closely align with the definition of incidence as we know it ([Jahja et al., 2022](#)).

From the correlation analysis between daily infection estimates and hospitalizations, a lag of 14 days gives the maximum average correlation across states. This is in agreement with the early estimates of the average time from infection to hospitalization of 9.7 days (95% CI: [5.4, 17.0]) for cases reported in January, 2020 in Wuhan, China as well as with estimates from across the pandemic in the UK that ranged from an average of 8.0 to 9.7 days, more precisely, 8.0 days (95% interval: [2.7, 18.5]) for the first wave to 9.7 days (95% interval: [4.1, 19.6]) for the second wave, ([Ward and Johnsen, 2021](#)). However, we should note the first study is based on a small sample size for outbreak cases reported well before our study start date. As well, both sets of estimates depend upon the healthcare system and the population structure, amongst other things ([Ward and Johnsen, 2021](#)). Nevertheless, their relative agreement with our estimate of 14 days for the US states lends some credence to of our results.

While we computed IHRs for all states, it is important to note that the IHR is also likely to vary within states and depend on additional variables such as age and the presence of major comorbidities ([Russell et al., 2023](#)). Therefore, it would be beneficial to account for such variables in the IHR calculations by, for example, stratifying infections and hospitalizations by age to produce age-specific estimates of the IHRs for each state (similar to [Fox et al., 2023](#) though with the additional element of being time-varying). We strongly believe this would be a worthwhile direction to pursue in future work should the necessary information be available.

The remainder of our discussion consists of an in-depth look into the advantages and limitations of our approach and of other comparable approaches, followed by a high level summary of our work and its major contributions.

Our approach offers a number of advantages. To the best of our knowledge, no other modelling approach has been used to reconstruct the infection time series for every state over as much of the COVID-19 pandemic as in this study. Furthermore, we aim to incorporate as much state-specific information as possible when deriving our estimates. For instance, using variant circulation and line list data, we are able to construct incubation and delay distributions that are unique for each state. By using time-varying and state-specific seroprevalence data, we are able to allow the reporting ratio to vary over both time and state, which is an advantage over such ratios that are non-time varying but state-specific and those that are time-varying but the same for all states ([Unwin et al., 2020](#); [Center for the Ecology of Infection Diseases, 2020](#)). Existing approaches that use the delay distribution to generate infection estimates often only construct one delay distribution that is used for all states ([Chitwood et al., 2022](#); [Jahja et al., 2022](#)). That is, they operate under the assumption of geographic invariance, where it is assumed that all states have the same patterns of delay from infection onset to case report, which is unlikely to be true due to differences in reporting pipelines, pandemic response, and variants in circulation, amongst other things.

Another major limitation is that these models do not account for reinfections. Now, it may be contended that reinfections do not account for a substantial fraction of the infections until later in the pandemic, so they are not absolutely necessary to include in the earlier stages of the pandemic. Still, at no stage did infection with the COVID-19 virus confer lifelong immunity. Rather immunity is transient and wanes over time. And we believe it is important to account for such defining characteristics of the virus when tracking infections over time. Therefore, we account for reinfections and waning of detectable immunity in our custom leaky immunity model. However, we acknowledge that the extent to which each of these are accounted for could be improved upon in future work.

Since the waning of detectable immunity is likely to be variant-dependent ([Pooley et al., 2023](#)), it follows that the leaky parameter may be better posed as a mixture of parameters for different variants with weights determined by the proportion of the variants circulating at the time in the state. Related to this is the issue of how newer variants may escape detection ([National Institutes of Health, 2022](#); [U.S. Food and Drug Administration, 2023](#)). While in a retrospective analysis where finalized data is used this is less likely to be an issue, this could very well pose a problem for real-time estimates of infections.

As for reinfections, it would be ideal to have confirmed rates over time for each US state. However, we are unable to find such data available over the entire time period considered for even one state. So we have turned to suspected reinfection data over time for Clark County, USA, as that surveillance is amongst the most detailed that we have found for the United States. Nevertheless, using such localized data raises questions of

representativeness and the applicability of such estimates to Nevada and all other states. Furthermore, this data has no information available beyond suspected third infections, which imposes an irremediable bias. However, based on the third infection data available there, we expect that the probability of being reinfected more than three times is likely very low for time frame considered and so the omission of these would impact our infection estimates to a small extent.

The vast majority of issues we encountered when trying to reconstruct the infection time series for each state are due to an absence or a lack of data. Such is the primary issue we had with the restricted line list. In comparison to the number of JHU cases (which we are treating as a gold standard) for the same release date, we noted there are about 10 million cases that are unaccounted for in the CDC line list. Moreover, the missingness does not appear to be random and uniformly distributed across states. Rather it is unequally distributed, suggesting that the dataset is likely biased. However, more information on the cases that are missing versus present would be required to determine the extent the missing cases led to a nonrepresentative, and therefore, biased sample, and could be a topic of further study.

Seroprevalence data also runs the risk of being nonrepresentative of the intended population (Bajema et al., 2021). For example, in the blood donor dataset some states have region specific-estimates, which clearly do not stand for the entire state. Another source of systematic variation is in the characteristics of the individuals who opt for blood tests versus those who do not. For instance, there may be a healthy user bias, in which a number of those who opt for blood tests are generally more inclined to partake in proactive healthy behaviors (such as checking on basic health markers by taking an annual blood test) than those who do not (Parsley Health, 2018). Alternatively, a number of individuals may be recommended for blood tests by their doctors due to signs of ill-health (ex. mineral deficiencies or underlying medical conditions). The extent that each such bias persists depends on the purpose of the blood test and whether it was used as a proactive or reactive medical tool. Since such information is unavailable to us, all we can conclude is that participant-driven sources of bias impact the seroprevalence samples to an undetermined extent. There are additional concerns about the performance of antibody testing for individuals with mild or asymptomatic disease as well as the loss of immunity over time (Kaku et al., 2021; Seow et al., 2020; Ibarrondo et al., 2020).

In this work, we do not attempt to directly address infection underascertainment due to the increase in asymptomatic infections across variants (Public Health Ontario, 2023). We simply note that this would likely pose a greater problem later in the pandemic, particularly during the Omicron era (Fan et al., 2022). We hope that such infections would be largely represented by the seroprevalence and reinfection estimates, but there is undoubtedly increasing reliance on such estimates to be able to do this over time (owing to the simultaneous decline in the reporting cadence and the apparent rise in asymptomatic infections over time) (Ontario Agency for Health Protection and Promotion, 2022; Garrett et al., 2022; Blauer, 2022b; Ren et al., 2021). Consequently, there is an increasing uncertainty over time that is not expressed by the model or the estimates.

Due to such concerns with the seroprevalence data, one further area of research is on investigating the utility of various sources to estimate the incidence of infections. Intuitively, one might expect that leveraging data from multiple sources would likely lead to more accurate and stable estimates than those from using one source. Wastewater surveillance data is one promising source that may be complementary to seroprevalence data, especially when testing is low (McManus et al., 2023). However, there has been limited success in predicting incidence using such data and the extent that wastewater concentration data is a useful in estimating COVID-19 incidence is unclear owing to problems with viral occurrence and detectability in wastewater that render detection inconsistent across locations (ex. due to temperature, per-capita water use, and in-sewer travel time) (McManus et al., 2023; Hart and Halden, 2020; Li et al., 2023). Sentinel surveillance streams for influenza-like illness or acute respiratory infection may provide decent proxies for COVID-19 incidence, especially when testing for mild cases of COVID-19 is diminishing or has ceased completely. Finally, alternative surveillance streams (potentially outside of public health) such as those from surveys, helplines, or medical records could potentially be integrated if they provide at least a rough indication of the disease intensity over time (European Centre for Disease Prevention and Control, 2020).

Overall, we adopt a relatively simple deconvolution-based approach and devote much of our efforts to tailoring our approach to the available data. A major result of this was the development of a way of to model immunity and space-time-specific reporting ratios based on seroprevalence data. In a way, our approach is built for the data rather than trying to force the data to fit to an existing approach. However, our model is only as good as the quantity and the quality of the data provided to it. In our case, the lack of data is both a

barrier to entry and a continual roadblock. The assumptions we are required to make as a consequence of this clearly limit the generalizability and call into question the reliability of the results. So while we highlight some interesting trends and numerical findings, these results are not definitive, but rather exploratory and intended to stimulate discussion on the challenging task of estimating infections. Despite these limitations, we are encouraged by the ability to use routine data to produce sensible estimates of infections in the United States and the plausibility of the apparent geospatial and temporal trends.

Our approach is predicated upon having case, line list, viral circulation, and seroprevalence data for each state, all of which are readily available (or available upon request in the case of restricted line list data). As a result of this, we are able to demonstrate the feasibility of estimating COVID-19 infections at the state level by using standard sources of data.

Our framework is quite versatile as it lends itself to more localized, county or community level estimates, or globalized, country-specific estimates. Fundamentally, to produce estimates of infections for different geographic regions, one would simply need to input the required data and re-run the pipeline. In this way, one could readily adapt our approach to generate estimates for the provinces in Canada or the regions in England.

Well-informed, localized estimates of COVID-19 infections over time can help us to have a more clear and comprehensive understanding of the course of the pandemic. Such estimates contribute important information on the timing and magnitude of disease burden for each location and they highlight trends that may not be visible from case data alone. Therefore, our infection estimates provide key information for the ongoing debate on the true size and impact of the pandemic.

Acknowledgements

We gratefully acknowledge all data contributors, i.e., the Authors and their Originating laboratories responsible for obtaining the specimens, and their Submitting laboratories for generating the genetic sequence and metadata and sharing via the GISAID Initiative ([Elbe and Buckland-Merrett, 2017](#)), on which this research is based.

References

- Bajema, K. L., Wiegand, R. E., Cuffe, K., Patel, S. V., Iachan, R., Lim, T., Lee, A., Moyse, D., Havers, F. P., Harding, L. et al. (2021) Estimated SARS-CoV-2 seroprevalence in the US as of September 2020. *JAMA Internal Medicine*, **181**, 450–460.
- Blauer, B. (2022a) Comparing cases, deaths, and hospitalizations indicates Omicron is less deadly. <https://coronavirus.jhu.edu/pandemic-data-initiative/data-outlook/comparing-cases-deaths-and-hospitalizations-indicates-omicron-less-deadly>.
- (2022b) Reduce data reporting cadence for an endemic disease? Not quite yet. <https://coronavirus.jhu.edu/pandemic-data-initiative/data-outlook/reduce-data-reporting-cadence-for-an-endemic-disease-not-quite-yet>.
- Carozzi, F., Provenzano, S. and Roth, S. (2022) Urban density and COVID-19: understanding the US experience. *The Annals of regional science*, 1–32.
- Center for the Ecology of Infection Diseases (2020) COVID-19 portal. <https://www.covid19.uga.edu/nowcast.html>.
- Centers for Disease Control and Prevention (2020a) COVID-19 case surveillance public use data. <https://data.cdc.gov/Case-Surveillance/COVID-19-Case-Surveillance-Public-Use-Data/vbim-akqf>.
- (2020b) COVID-19 case surveillance restricted access detailed data. <https://data.cdc.gov/Case-Surveillance/COVID-19-Case-Surveillance-Restricted-Access-Detail/mbd7-r32t>.
- (2020c) COVID Data Tracker. <https://covid.cdc.gov/covid-data-tracker/#national-lab>.
- (2021a) 2020-2021 nationwide blood donor seroprevalence survey infection-induced seroprevalence estimates. <https://data.cdc.gov/Laboratory-Surveillance/2020-2021-Nationwide-Blood-Donor-Seroprevalence-Su/mtc3-kq6r>.
- (2021b) Nationwide commercial laboratory seroprevalence survey. <https://data.cdc.gov/Laboratory-Surveillance/Nationwide-Commercial-Laboratory-Seroprevalence-Su/d2tw-32xv>.
- (2022) Estimated COVID-19 burden. <https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/burden.html>.
- Centers for Medicare & Medicaid Services (2020) Health expenditures by state of residence, 1991–2014.
- Chitwood, M. H., Russi, M., Gunasekera, K., Havumaki, J., Klaassen, F., Pitzer, V. E., Salomon, J. A., Swartwood, N. A., Warren, J. L., Weinberger, D. M. et al. (2022) Reconstructing the course of the COVID-19 epidemic over 2020 for US states and counties: Results of a Bayesian evidence synthesis model. *PLOS Computational Biology*, **18**, e1010465.
- Cortés Martínez, J., Pak, D., Abelenda-Alonso, G., Langohr, K., Ning, J., Rombauts, A., Colom, M., Shen, Y. and Gómez Melis, G. (2022) SARS-CoV-2 incubation period according to vaccination status during the fifth COVID-19 wave in a tertiary-care center in Spain: A cohort study. *BMC Infectious Diseases*, **22**, 1–7.
- Di Lorenzo, P. (2023) *usmap*. URL: <https://usmap.dev>. R package version 0.6.2.
- Do, K., Do, J., Kawana, E., Zhang, R., Do, K. H. and Zhang, R. Y. (2023) Nevada’s healthcare crisis: A severe shortage of physicians and residency positions. *Cureus*, **15**, e41700.
- Dong, E., Du, H. and Gardner, L. (2020) An interactive web-based dashboard to track COVID-19 in real time. *The Lancet Infectious Diseases*, **20**, 533–534.
- Duerr, R., Dimartino, D., Marier, C., Zappile, P., Wang, G., Lighter, J., Elbel, B., Troxel, A. B., Heguy, A. et al. (2021) Dominance of Alpha and Iota variants in SARS-CoV-2 vaccine breakthrough infections in New York City. *The Journal of Clinical Investigation*, **131**, e152702.

- Dunkel, S. (2020) COVID-19 case numbers: Why the delay in reporting? <https://www.tpchd.org/Home/Components/Blog/Blog/21448>.
- Durbin, J. and Koopman, S. J. (2012) *Time Series Analysis by State Space Methods*, vol. 38. OUP Oxford.
- Elbe, S. and Buckland-Merrett, G. (2017) Data, disease and diplomacy: GISAID’s innovative contribution to global health. *Global Challenges*, **1**, 33–46.
- Embrett, M., Sim, S. M., Caldwell, H. A., Boulos, L., Yu, Z., Agarwal, G., Cooper, R., Aj, A. J. G., Bielska, I. A., Chishtie, J. et al. (2022) Barriers to and strategies to address COVID-19 testing hesitancy: A rapid scoping review. *BMC Public Health*, **22**, 1–10.
- European Centre for Disease Prevention and Control (2020) Strategies for the surveillance of COVID-19. *Technical report*, ECDC, Stockholm, Sweden.
- Fan, Y., Li, X., Zhang, L., Wan, S., Zhang, L. and Zhou, F. (2022) SARS-CoV-2 Omicron variant: Recent progress and future perspectives. *Signal Transduction and Targeted Therapy*, **7**, 141.
- Fox, S. J., Javan, E., Pasco, R., Gibson, G. C., Betke, B., Herrera-Diestra, J. L., Woody, S., Pierce, K., Johnson, K. E., Johnson-León, M. et al. (2023) Disproportionate impacts of COVID-19 in a large US city. *PLOS Computational Biology*, **19**, e1011149.
- Garrett, N., Tapley, A., Andriesen, J., Seocharan, I., Fisher, L. H., Bunts, L., Espy, N., Wallis, C. L., Randhawa, A. K., Ketter, N. et al. (2022) High rate of asymptomatic carriage associated with variant strain Omicron. *MedRxiv*.
- Goldberg, Y., Mandel, M., Bar-On, Y. M., Bodenheimer, O., Freedman, L. S., Ash, N., Alroy-Preis, S., Huppert, A. and Milo, R. (2022) Protection and waning of natural and hybrid immunity to SARS-CoV-2. *New England Journal of Medicine*, **386**, 2201–2212.
- Gong, G., Phillips, S. G., Hudson, C., Curti, D. and Philips, B. U. (2019) Higher US rural mortality rates linked to socioeconomic status, physician shortages, and lack of health insurance. *Health Affairs*, **38**, 2003–2010.
- Grant, R., Charmet, T., Schaeffer, L., Galmiche, S., Madec, Y., Von Platen, C., Chény, O., Omar, F., David, C., Rogoff, A. et al. (2022) Impact of SARS-CoV-2 Delta variant on incubation, transmission settings and vaccine effectiveness: Results from a nationwide case-control study in France. *The Lancet Regional Health–Europe*, **13**, 100278.
- Hart, O. E. and Halden, R. U. (2020) Computational analysis of SARS-CoV-2/COVID-19 surveillance by wastewater-based epidemiology locally and globally: Feasibility, economy, opportunities and challenges. *Science of the Total Environment*, **730**, 138875.
- Helske, J. (2017) KFAS: Exponential family state space models in R. *Journal of Statistical Software*, **78**, 1–39.
- Hitchings, M. D., Dean, N. E., García-Carreras, B., Hladish, T. J., Huang, A. T., Yang, B. and Cummings, D. A. (2021) The usefulness of the test-positive proportion of severe acute respiratory syndrome coronavirus 2 as a surveillance tool. *American Journal of Epidemiology*, **190**, 1396–1405.
- Hodcroft, E. (2021) CoVariants: SARS-CoV-2 mutations and variants of interest. <https://covariants.org>.
- Ibarrondo, F. J., Fulcher, J. A., Goodman-Meza, D., Elliott, J., Hofmann, C., Hausner, M. A., Ferbas, K. G., Tobin, N. H., Aldrovandi, G. M. and Yang, O. O. (2020) Rapid decay of anti-SARS-CoV-2 antibodies in persons with mild COVID-19. *New England Journal of Medicine*, **383**, 1085–1087.
- Jahja, M., Chin, A. and Tibshirani, R. J. (2022) Real-time estimation of COVID-19 infections: Deconvolution and sensor fusion. *Statistical Science*, **37**, 207–228.

- Jones, J. M., Stone, M., Sulaeman, H., Fink, R. V., Dave, H., Levy, M. E., Di Germanio, C., Green, V., Notari, E., Saa, P. et al. (2021) Estimated US infection-and vaccine-induced SARS-CoV-2 seroprevalence based on blood donations, July 2020-May 2021. *JAMA*, **326**, 1400–1409.
- Kaku, N., Nishimura, F., Shigeishi, Y., Tachiki, R., Sakai, H., Sasaki, D., Ota, K., Sakamoto, K., Kosai, K., Hasegawa, H. et al. (2021) Performance of anti-SARS-CoV-2 antibody testing in asymptomatic or mild COVID-19 patients: A retrospective study in outbreak on a cruise ship. *PLoS One*, **16**, e0257452.
- Li, X., Zhang, S., Sherchan, S., Orive, G., Lertxundi, U., Haramoto, E., Honda, R., Kumar, M., Arora, S., Kitajima, M. et al. (2023) Correlation between SARS-CoV-2 RNA concentration in wastewater and COVID-19 cases in community: A systematic review and meta-analysis. *Journal of Hazardous Materials*, **441**, 129848.
- Lorenzo-Redondo, R., Ozer, E. A. and Hultquist, J. F. (2022) COVID-19: Is Omicron less lethal than Delta? *British Medical Journal*, **378**.
- McManus, O., Christiansen, L. E., Nauta, M., Krogsgaard, L. W., Bahrenscheer, N. S., von Kappelgaard, L., Christiansen, T., Hansen, M., Hansen, N. C., Kähler, J. et al. (2023) Predicting COVID-19 incidence using wastewater surveillance data, Denmark, October 2021–June 2022. *Emerging Infectious Diseases*, **29**, 1589.
- National Institutes of Health (2022) Assessing how SARS-CoV-2 mutations might affect rapid tests. <https://www.nih.gov/news-events/nih-research-matters/assessing-how-sars-cov-2-mutations-might-affect-rapid-tests>.
- Nyberg, T., Ferguson, N. M., Nash, S. G., Webster, H. H., Flaxman, S., Andrews, N., Hinsley, W., Bernal, J. L., Kall, M., Bhatt, S. et al. (2022) Comparative analysis of the risks of hospitalisation and death associated with SARS-CoV-2 Omicron (B. 1.1. 529) and Delta (B. 1.617. 2) variants in England: A cohort study. *The Lancet*, **399**, 1303–1312.
- Ogata, T., Tanaka, H., Irie, F., Hirayama, A. and Takahashi, Y. (2022) Shorter incubation period among unvaccinated delta variant coronavirus disease 2019 patients in Japan. *International Journal of Environmental Research and Public Health*, **19**, 1127.
- Ontario Agency for Health Protection and Promotion (2022) COVID-19 variant of concern Omicron (B.1.1.529): Risk assessment. https://www.publichealthontario.ca/-/media/documents/ncov/voc/2022/01/covid-19-omicron-b11529-risk-assessment-jan-6.pdf?sc_lang=en.
- Parsley Health (2018) 5 essential blood tests you need every year. <https://www.parsleyhealth.com/blog/5-essential-blood-tests-need-every-year/>.
- Pellis, L., Scarabel, F., Stage, H. B., Overton, C. E., Chappell, L. H., Fearon, E., Bennett, E., Lythgoe, K. A., House, T. A., Hall, I. et al. (2021) Challenges in control of COVID-19: Short doubling time and long delay to effect of interventions. *Philosophical Transactions of the Royal Society B*, **376**, 20200264.
- Pitzer, V. E., Chitwood, M., Havumaki, J., Menzies, N. A., Perniciaro, S., Warren, J. L., Weinberger, D. M. and Cohen, T. (2021) The impact of changes in diagnostic testing practices on estimates of COVID-19 transmission in the United States. *American Journal of Epidemiology*, **190**, 1908–1917.
- Pooley, N., Abdool Karim, S. S., Combadière, B., Ooi, E. E., Harris, R. C., El Guerche Seblain, C., Kisomi, M. and Shaikh, N. (2023) Durability of vaccine-induced and natural immunity against COVID-19: A narrative review. *Infectious Diseases and Therapy*, **12**, 367–387.
- Provenzano, S., Roth, S. and Carozzi, F. (2020) Urban density and COVID-19. *CEP Discussion Paper*.
- Public Health Agency of Canada (2021) COVID-19 for health professionals: Transmission. <https://www.canada.ca/en/public-health/services/diseases/2019-novel-coronavirus-infection/health-professionals/transmission.html>.

- Public Health Ontario (2023) COVID-19 Omicron variant of concern and communicability – What we know so far. <https://www.publichealthontario.ca/-/media/documents/ncov/covid-wksf/2022/01/wksf-omicron-communicability.pdf>.
- Radley, D. C., Collins, S. R. and Baumgartner, J. C. (2020) 2020 scorecard on state health system performance.
- Ramdas, A. and Tibshirani, R. J. (2016) Fast and flexible ADMM algorithms for trend filtering. *Journal of Computational and Graphical Statistics*, **25**, 839–858.
- Reinhart, A., Brooks, L., Jahja, M., Rumack, A., Tang, J., Agrawal, S., Al Saeed, W., Arnold, T., Basu, A., Bien, J. et al. (2021) An open repository of real-time COVID-19 indicators. *Proceedings of the National Academy of Sciences*, **118**, e2111452118.
- Ren, R., Zhang, Y., Li, Q., McGoogan, J. M., Feng, Z., Gao, G. F. and Wu, Z. (2021) Asymptomatic SARS-CoV-2 infections among persons entering China from april 16 to october 12, 2020. *Jama*, **325**, 489–492.
- Ruff, J., Zhang, Y., Kappel, M., Rath, S., Watkins, K., Zhang, L. and Lockett, C. (2022) Rapid increase in suspected SARS-CoV-2 reinfections, Clark County, Nevada, USA, December 2021. *Emerging Infectious Diseases*, **28**, 1977.
- Russell, C. D., Lone, N. I. and Baillie, J. K. (2023) Comorbidities, multimorbidity and COVID-19. *Nature Medicine*, **29**, 334–343.
- Schechtman, K. (2020) Counting COVID-19 tests: How states do it, how we do it, and what’s changing. <https://covidtracking.com/analysis-updates/counting-covid-19-tests>.
- Seow, J., Graham, C., Merrick, B., Acors, S., Pickering, S., Steel, K. J., Hemmings, O., O’Byrne, A., Kouphou, N., Galao, R. P. et al. (2020) Longitudinal observation and decline of neutralizing antibody responses in the three months following SARS-CoV-2 infection in humans. *Nature Microbiology*, **5**, 1598–1607.
- Simon, S. (2021) Inconsistent reporting practices hampered our ability to analyze COVID-19 data. Here are three common problems we identified. <https://covidtracking.com/analysis-updates/three-covid-19-data-problems>.
- Tanaka, H., Ogata, T., Shibata, T., Nagai, H., Takahashi, Y., Kinoshita, M., Matsubayashi, K., Hattori, S. and Taniguchi, C. (2022) Shorter incubation period among COVID-19 cases with the BA. 1 Omicron variant. *International Journal of Environmental Research and Public Health*, **19**, 6330.
- The New York Times (2020) Coronavirus in the U.S.: Latest map and case count. <https://www.nytimes.com/interactive/2021/us/covid-cases.html>.
- The Washington Post (2020) Tracking U.S. COVID-19 cases, deaths and other metrics by state. <https://www.washingtonpost.com/graphics/2020/national/coronavirus-us-cases-deaths/?state=US>.
- Tibshirani, R. J. (2014) Adaptive piecewise polynomial estimation via trend filtering. *The Annals of Statistics*, **42**, 285–323.
- (2022) Divided differences, falling factorials, and discrete splines: Another look at trend filtering and related problems. *Foundations and Trends in Machine Learning*, **15**, 694–846.
- Tindale, L. C., Stockdale, J. E., Coombe, M., Garlock, E. S., Lau, W. Y. V., Saraswat, M., Zhang, L., Chen, D., Wallinga, J. and Colijn, C. (2020) Evidence for transmission of COVID-19 prior to symptom onset. *eLife*, **9**, e57149.
- Twohig, K. A., Nyberg, T., Zaidi, A., Thelwall, S., Sinnathamby, M. A., Aliabadi, S., Seaman, S. R., Harris, R. J., Hope, R., Lopez-Bernal, J. et al. (2022) Hospital admission and emergency care attendance risk for SARS-CoV-2 Delta (B. 1.617. 2) compared with Alpha (B. 1.1. 7) variants of concern: A cohort study. *The Lancet Infectious Diseases*, **22**, 35–42.

- United States Census Bureau (2021) 2020 Census: Percent change in resident population for the 50 states, the District of Columbia, and Puerto Rico: 2010 to 2020. <https://www.census.gov/library/visualizations/2021/dec/2020-percent-change-map.html>.
- Unwin, H. J. T., Mishra, S., Bradley, V. C., Gandy, A., Mellan, T. A., Coupland, H., Ish-Horowicz, J., Vollmer, M. A., Whittaker, C., Filippi, S. L. et al. (2020) State-level tracking of COVID-19 in the United States. *Nature Communications*, **11**, 6189.
- U.S. Census Bureau, Population Division (2022) Annual estimates of the resident population for the United States, regions, states, District of Columbia, and Puerto Rico: April 1, 2020 to July 1, 2022. <https://www.census.gov/data/tables/time-series/demo/popest/2020s-state-total.html>.
- U.S. Food and Drug Administration (2023) SARS-CoV-2 viral mutations: Impact on COVID-19 tests. <https://www.fda.gov/medical-devices/coronavirus-covid-19-and-medical-devices/sars-cov-2-viral-mutations-impact-covid-19-tests>.
- Ward, T. and Johnsen, A. (2021) Understanding an evolving pandemic: An analysis of the clinical time delay distributions of COVID-19 in the United Kingdom. *PLoS One*, **16**, e0257978.
- Washington State Department of Health (2020) COVID-19 data dashboard. <https://doh.wa.gov/emergencies/covid-19/data-dashboard>.
- World Health Organization (2021) Tracking SARS-CoV-2 variants. <https://www.who.int/activities/tracking-SARS-CoV-2-variants>.
- Yang, S., Hemarajata, P., Hilt, E. E., Price, T. K., Garner, O. B. and Green, N. M. (2022) Investigation of SARS-CoV-2 Epsilon variant and hospitalization status by genomic surveillance in a single large health system during the 2020-2021 winter surge in Southern California. *American Journal of Clinical Pathology*, **157**, 649–652.
- Zaki, N. and Mohamed, E. A. (2021) The estimations of the COVID-19 incubation period: A scoping reviews of the literature. *Journal of Infection and Public Health*, **14**, 638–646.

Online Supplement

S1 Additional information about estimation methodology

S1.1 State space representation of the leaky immunity model

To estimate the leaky immunity model, [Equation 1](#), we express it as a Gaussian state space model (as in [Durbin and Koopman, 2012](#); [Helske, 2017](#)).

In general, for $t = 1, \dots, n$, we let α_t be the $m \times 1$ vector of latent state processes at time t and y_t be the $p \times 1$ vector of observations at time t . Under the assumption that η is a $k \times 1$ vector, the form of the linear Gaussian state space model is

$$y_t = Z\alpha_t + \epsilon_t \text{ (observation equation)} \quad (2)$$

$$\alpha_{t+1} = T_t\alpha_t + R_t\eta_t \text{ (state equation)} \quad (3)$$

where $\epsilon_t \sim N(0, H_t)$, $\eta_t \sim N(0, Q_t)$, and $\alpha_1 \sim N(a_1, P_1)$ independently of each other ([Helske, 2017](#)). For notational compactness, we let $\alpha = (\alpha_1^\top, \dots, \alpha_n^\top)$ and $y = (y_1^\top, \dots, y_n^\top)$.

The observation equation can be viewed as a linear regression model with the time-varying coefficient α_t , while the second equation is a first-order autoregressive model, which is Markovian in nature ([Durbin and Koopman, 2012](#)).

The main idea of the two equations is that the system evolves over time according to α_t (as in the second equation), but since those states are not directly observed, we turn to the observations y_t and use their relationship with α_t (as in the first equation) to drive the system forward ([Durbin and Koopman, 2012](#)). So the objective of state space modeling is to obtain the latent states α based on the observations y and this is achieved through Kalman filtering and smoothing.

Kalman filtering gives the one-step-ahead predictions and prediction errors:

$$\begin{aligned} a_{t+1} &= \mathbb{E}[\alpha_{t+1} \mid y_t, \dots, y_1] \\ v_t &= y_t - Za_t \end{aligned}$$

with covariance,

$$\begin{aligned} P_{t+1} &= \text{Var}(\alpha_{t+1} \mid y_t, \dots, y_1) \\ \text{Var}(v_t) &= ZP_tZ^\top + H_t. \end{aligned}$$

Then, the state smoothing equations are run back in time to give

$$\hat{a}_t = \mathbb{E}[\alpha_t \mid y_n, \dots, y_1] \quad (4)$$

$$V_t = \text{Var}(\alpha_t \mid y_n, \dots, y_1). \quad (5)$$

The filtering and smoothing steps are based on recursions that are described in Appendix A of [Helske \(2017\)](#) as we use the R package KFAS to estimate our model.

For our situation, the Kalman filter and smoothing approach offers a number of advantages over the penalized regression approach. Perhaps most notably, the parameters are estimated all at once (so cross validating for model parameter tuning is not necessary). Another major benefit is that it can handle unevenly spaced time series (refer to [Durbin and Koopman, 2012](#) for further details).

To help manage the sparseness in the seroprevalence data, we convert our data to weekly by summing the reported infections and shifting the observed seroprevalence measurements to the nearest Monday. If there are multiple measurements in a week from a source, then the average of those measurements is used (and similarly for the weights). We denote these changes by changing the time-based subscript from t to m where m indicates the Monday relative to our March 9, 2020 start date and ΔR_m is the change in weekly reported infections.

The leaky immunity model on weekly data can be written in state space form by defining the corresponding components in [Equations 2 and 3](#) as follows:

$$\begin{aligned}
R &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} & Z &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} & H_m &= \begin{bmatrix} w_{m,c}\sigma_o^2 & 0 \\ 0 & w_{m,b}\sigma_o^2 \end{bmatrix} \\
\alpha_m &= \begin{bmatrix} s_m \\ a_m \\ a_{m-1} \\ a_{m-2} \end{bmatrix} & T_m &= \begin{bmatrix} \gamma & \Delta R_m n_m & 0 & 0 \\ 0 & 3 & -3 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} & Q &= \begin{bmatrix} \sigma_s^2 & 0 \\ 0 & \sigma_a^2 \end{bmatrix} \\
a_1 &= \begin{bmatrix} \tilde{s}_1 \\ \tilde{a}_1 \\ \tilde{a}_1 \\ \tilde{a}_1 \end{bmatrix} & P_1 &= \begin{bmatrix} \sigma_{\tilde{s}_1}^2 & 0 & 0 & 0 \\ 0 & \sigma_{\tilde{a}_1}^2 & 0 & 0 \\ 0 & 0 & \sigma_{\tilde{a}_1}^2 & 0 \\ 0 & 0 & 0 & \sigma_{\tilde{a}_1}^2 \end{bmatrix}
\end{aligned}$$

where σ_o^2 is the variance of observations, σ_s^2 is the variance of the seroprevalence estimates (in terms of smoothness), and σ_a^2 is the trend variance. Since we expect the inverse ratios to be more variable than the seroprevalence estimates, we enforce that the estimate of σ_a^2 is a multiple of σ_s^2 . Letting the subscripts b and c denote the blood donor and commercial datasets, $w_{m,c}$ and $w_{m,b}$ are the time-varying inverse variance weights computed from the commercial and blood donor datasets, respectively.

For each source, we compute the weights for the observed seroprevalence estimates using the standard formula for the standard error of a proportion. These weights are then re-scaled so they sum to the number of observed seroprevalence measurements for the source. All days that are unobserved (i.e., lack seroprevalence measurements) are given weights of one. Finally, the ratio of the average observed weights for the sources is used as a multiplier to scale all of the weights for one source. For example, if the average weight of the commercial source is double the average weight of the blood donor source (for an arbitrary state), then we scale all of the weights in the commercial source (including the ones) by two. The main purpose of this step is to ensure that the source with a greater sample size contributes more weight in the model on average.

The prior distribution for α_1 is estimated using both data-driven constraints and externally sourced information. To obtain the initial value of the seroprevalence component, \tilde{s}_1 , we extract the first observed seroprevalence measurement from each source, round down to two decimal places, and take the average to be the estimated initial value \tilde{s}_1 . The corresponding initial variance estimate, $\sigma_{\tilde{s}_1}^2$, is taken to be the mean of the standard errors of the two seroprevalence estimates. For all of the initial values of the trend components, we use the inverse of the ascertainment ratio estimate as of June 1, 2020 for each state from Table 1 in [Unwin et al. \(2020\)](#) and denote this by \tilde{a}_1 . The initial variance estimate of $\sigma_{\tilde{a}_1}^2$ is based on the variance implied by the given inverse ascertainment ratio distribution.

The initial σ_o^2 is taken to be the average of the estimated variances from the linear models for the sources where the observed seroprevalence measurements are regressed on the enumerated dates. The initial value of the multiplier is set to be 100 for all states. The σ_s^2 and γ values are fixed and obtained by averaging the final values for all states on the real line.

Following the maximum likelihood estimation of the two non-fixed parameters we use the Kalman filtering and smoothing to obtain the smoothed estimates of the weekly inverse reporting ratios and their covariance matrices as shown in Equations 4 and 5. Forwards and backwards extrapolation is then used to estimate the ratios and covariance outside of the observed seroprevalence range ([Durbin and Koopman, 2012](#)), followed by linear interpolation to fill-in estimates for each day in our considered time period. After we obtain one vector of inverse reporting ratios for each state in this way, we take each inverse reporting ratio and multiply it by the corresponding deconvolved case estimate (that has undergone linear interpolation to correct instances of 0 reported infections) to obtain an estimate of new infections. We are able to convert these numbers of infections to infections per 100,000 population by simple re-scaling (enabled by the fact that normality is preserved under linear transformations).

The 50, 80, and 95% confidence intervals are constructed by taking a Bayesian view of the leaky immunity model (refer to the Online Supplement [S1.2](#) for the Bayesian specification of the model). That is, for each time, t , we obtain an estimate of the posterior variance of a_t , apply the deconvolved case estimate as a constant multiplier, and then use resulting variance to build a normal confidence interval about the infection

estimate. We additionally enforce that the lower bound must be at least the deconvolved case estimate for the time under consideration.

S1.2 Bayesian specification of the leaky immunity model

In brief, the leaky immunity model where we let $\beta = \{\gamma, a_1, \dots, a_t\}$ and X be the design matrix, corresponds to a Bayesian model with prior

$$\beta \sim N\left(0, \frac{\sigma^2}{\lambda} (A^T D^T D A)^{-1}\right)$$

and likelihood

$$s|X, \beta \sim N(X\beta, \sigma^2 W^{-1}),$$

where A is indicator matrix save for the first column of 0s (corresponding to γ), D represents the discrete derivative matrix of order 3, and W is the inverse variance weights matrix. Then, the posterior on a_t is normally distributed with mean

$$(X^T W X + \lambda A^T D^T D A)^{-1} X^T W s$$

and variance

$$\sigma^2 (X^T W X + \lambda A^T D^T D A)^{-1}.$$