

# Retrospective estimation of latent COVID-19 infections before Omicron in the U.S.

Rachel Lobay<sup>a,1</sup>, Maria Jahja<sup>b</sup>, Ajitesh Srivastava<sup>c</sup>, Ryan J. Tibshirani<sup>d</sup>, and Daniel J. McDonald<sup>a</sup>

<sup>a</sup>Department of Statistics, The University of British Columbia

<sup>b</sup>Department of Statistics & Data Science, Carnegie Mellon University

<sup>c</sup>Department of Computer and Electrical Engineering, University of Southern California

<sup>d</sup>Department of Statistics, The University of California, Berkeley

Version: March 26, 2024

## Abstract

The true timing and magnitude of COVID-19 infections (rather than reported cases) are of interest to both the public and to public health, but these are challenging to pin down for a variety of data-driven and methodological reasons. Accurate estimates of latent COVID-19 infections can improve our understanding of the size and scope of the pandemic and provide more meaningful and timely quantification of disease patterns and burden. In this work, we estimate daily incident *infections* for each U.S. state. Rather than taking a model-based approach, our methods operate directly on data. We first deconvolve reported COVID-19 cases to their infection date using delay distributions estimated from the CDC linelist. We combine these deconvolved cases with serology data to scale up to unreported infections. Our results cover all states at the daily frequency, incorporate variant-specific incubation periods, and account for reinfections and waning antigenic immunity. This analysis also produces estimates for other important quantities such as the number of deconvolved cases specific to each variant and the infection-case-report ratio. We also discuss some implications of our results: a disease burden that appears earlier and more extensively than previously quantified; differential infection-hospitalization ratio estimates. Our findings help to better understand the impact of the pandemic in the U.S. prior to the onset of Omicron and its descendants.

## 1 Introduction

Reported COVID-19 cases are a staple in tracking the pandemic at varying geographic resolutions<sup>1–3</sup>. Yet, for every case that is eventually reported to public health, several infections are likely to have occurred, likely much earlier. To see why, it is important to understand *whose* cases are being reported and what differentiates them from the unreported cases as well as *when* these case reports happen. Figure 1 shows 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 omitted<sup>4</sup>. In addition, testing practices, availability, and uptake vary temporally and spatially<sup>5–7</sup>. 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 eventual submission to public health<sup>8;9</sup>. 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. Since there was no large-scale surveillance effort in the United States that reliably tracked symptom onset, let alone infection onset, ascertaining the onset of all *infections* is challenging.

---

<sup>1</sup>To whom correspondence should be addressed. E-mail: [rachel.lobay@stat.ubc.ca](mailto:rachel.lobay@stat.ubc.ca)

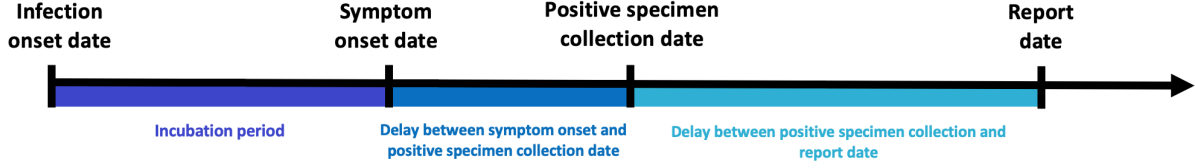


Figure 1: Idealized chain of events from infection onset to case report date for a symptomatic infection that is eventually reported to public health.

Explaining the course of the pandemic and investigating the effects of interventions, the burden facing various subgroups, and drawing insights for future pandemics is inhibited because the true spatial and temporal behaviour of infections is unknown. While reported cases provide a convenient proxy 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 explain its implications—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 understanding of the pandemic burden over space and time.

In this work, we provide a statistically rigorous, data-first reconstruction of daily incident infections for each U.S. state from June 1, 2020 to November 29, 2021. Using state-level line list data, we construct time-varying delay distributions for the time from symptom onset to positive specimen date and positive specimen to case report date. We combine these with variant-specific incubation period distributions to deconvolve daily reported COVID-19 cases back to their infection onset. Finally, the resulting deconvolved cases are adjusted to account for the unreported infections using seroprevalence and reinfection data to estimate adjust for the waning of antibody detectability over time. We examine some features of our infection estimates and the implications of using them rather than reported cases in assessing the impact of the pandemic. We produce simple 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 (and the resulting publicly-available estimates) will prove an important benchmark for others to undertake retrospective analyses.

## 2 Results

An important aspect of our methods is that deconvolution is not the same as a shift. **ATTN: We need a concise description of this somewhere, possibly with a graphic.**

### 2.1 Infection estimates reveal waves missed by reported cases

Outbreaks in infections precede those in reported cases and are reliably larger in magnitude. But simply shifting cases back in time and increasing them by some factor fails to capture the spatio-temporal dynamics of the pandemic. Hence, relative to reported cases, examining estimated infections reveals a rather different pattern. [Figure 2](#) shows estimates of the number of daily new infections per 100,000 inhabitants for each U.S. state from June 1, 2020 to November 29, 2021 compared with reported cases, and deconvolved cases (reported cases “pushed back” by the delays shown in [Figure 1](#)).

While the major Ancestral, Alpha, and Delta 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, a wave of infections is present in the spring of 2021 for North Dakota and South Dakota which is not visible in reported cases alone. **ATTN: we should search for more of these patterns.**

## 2.2 Spatial-temporal implications are ignored

ATTN: We need to expand the discussion of this figure along the lines we discussed in our meeting. Choose the dates carefully to emphasize that the spatial extent during peaks/troughs of waves is different if you look at infections instead of cases.

Figure 3 shows that for the earliest time of June 1, 2020, there is little discrepancy between case and infection rates, while for the later times there are immense differences in the rates, such that case rates tend to underrepresent infections to a great extent.

## 2.3 Infection/case ratios vary by state and VOC

ATTN: The rest of this section uses a lot of space to say, basically, that infections are bigger than cases. I think we should cut much of it (though see also below). I think a better use is to more explicitly emphasize how the case/infection ratio changes with time and VOC. Let's make that more focused, direct. Note that this isn't really "underreporting", they're still reporting all the tests, but those tests capture different numbers of infections.

With respect to variants of concern, consider the late 2020 Ancestral wave for the midwestern states of Illinois, Indiana, and Ohio. For the major Delta wave, some of the greatest discrepancies between cases and infections are visible in the western states of Idaho and Montana, the southern states of Louisiana and Georgia, and the midwestern states of Iowa and Nebraska (Figure 2). Earlier on in the pandemic, such discrepancies between cases and infections 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<sup>10,11</sup>.

Finally, while the main Delta wave is somewhat evident from the case counts for all states (Figure 2), our estimates suggest that case counts tend to severely underestimate infections during this time for many states. The lowest of all states was in New Jersey, where about 4.6% (95% confidence interval: [1.9, 67.7]) of the estimated infections were reported. This was followed by Maryland with 7.4% ([2.7, 83.8]), Connecticut with 8.0% ([3.1, 25.8]), and Florida with 8.7% ([4.8, 34.0]). This underreporting issue extends to most states as in 39 states less than 30% of infections were reported during this time. Only 4 states of Alaska, Maine, Vermont and Virginia reported at least 40% infections. No states were found to surpass 50% for reported infections for this time.

Similar patterns were observed during the earlier period of Alpha domination, where Louisiana had the lowest reported infections at 11.7% (95% confidence interval: [6.7, 31.5]) and was followed by California at 14.4% (95% confidence interval: [7.7, 68.2]). There were 23 states that reported at least 40% and 22 states that reported at least 50% of their infections.

Such patterns were comparatively less apparent during the earlier and larger period of Ancestral domination, where Ohio and Maryland held the lowest percentages of reported infections at 22.0% (95% confidence interval: [16.2, 34.0]) and 22.3% (95% confidence interval: [14.8, 40.5]), respectively. During this time, 28 states that reported at least 40% and 14 states that reported at least 50% of their infections.

## 2.4 Infections broken down by VOCs emphasize earlier outbreaks

Figure 4 examines the infection estimates for a selection of states more closely. This set has the largest infection/case ratios. ATTN: We need to say more about what these 2 figures show, uniquely. What's the point of looking at them? I think the section heading I wrote is the point. But we need to say this. The takeaways below are a bit too vague, I think. The top panel shows ATTN: .... The bottom panel divides estimated infections into buckets based on the circulating variant proportions at the time. From these plots, it is clear that few variant categories tends to dominate and drive infections at a time. The general progression in terms of variant starts with the Ancestral category from 2020 up to early 2021, to the Alpha variant in mid 2021, which eventually gets eclipsed by the Delta variant in mid to late 2021. This supports our division of our results by the three main variant-driven time periods.

## 2.5 The relationship between infections and hospitalizations is messy

ATTN: messy, but much more stable than CHR We systematically investigate the temporal relationship between infections and hospitalizations with Spearman's rank-correlation across different lags, shifting

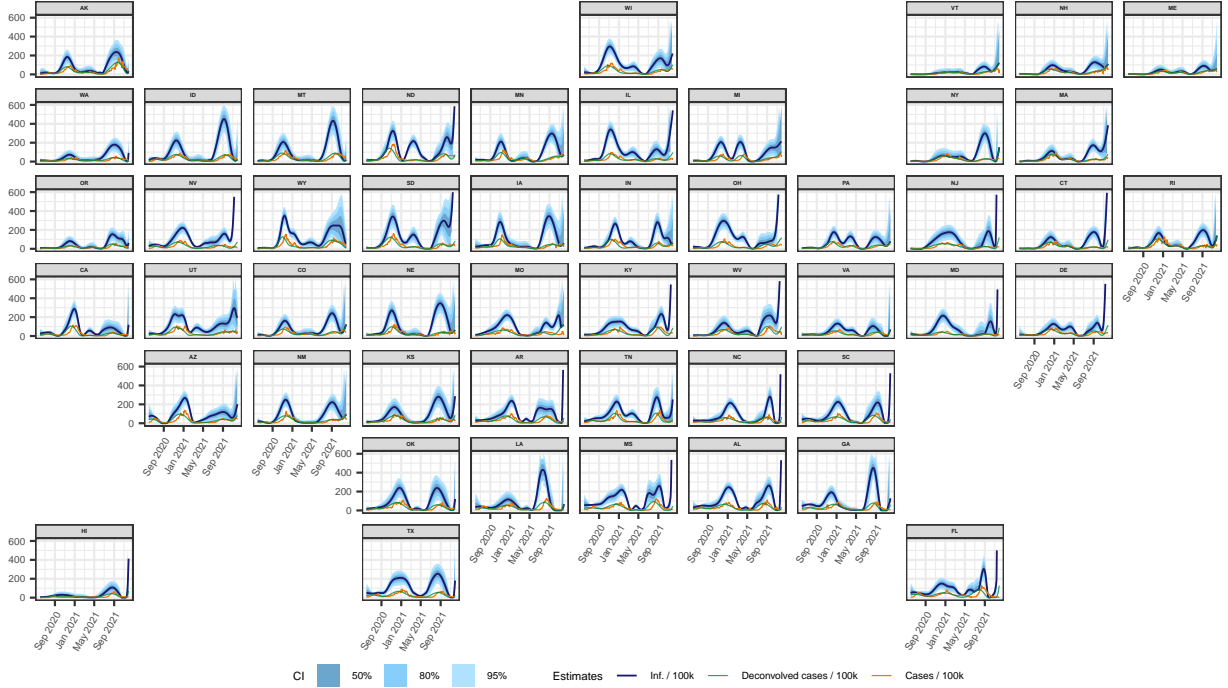


Figure 2: Estimates of the number of daily new infections per 100,000 population for each U.S. state from June 1, 2020 to November 29, 2021 (dark blue line). The blue shaded regions depict the 50, 80, and 95% confidence intervals for the estimates, while the teal line represents the number of new daily 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.

hospitalizations backward to align with infections. (Figure 5). The maximum average correlation across states is 0.513, occurring at a lag of 13 days. In contrast, we find that the greatest average Spearman correlation for cases is 0.691 and occurs at a lag of 1 day. That is, we find that case report rates are nearly contemporaneous to hospitalizations, while infection estimates clearly precede them.

The maximum correlation at a lag of 13 days is in similar to 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<sup>12</sup>.

Unsurprisingly, the deconvolved case and infection estimates achieve their maximum correlation at the same lag **ATTN: I don't think so....** And yet, the average correlation to hospitalizations tends to be greater for the deconvolved case estimates than for the infection estimates. This finding may stem from a difference in disease severity between the reported and unreported infections: unreported infections tend to be less severe and less likely to lead to hospitalization than those that are reported.

## 2.6 Estimating infection-hospitalization ratios

As a counterpart to the correlation analysis, we compute the time-varying infection-hospitalization ratios (IHRs) for each state using the correlation maximizing lag. We similarly compute the case-hospitalization ratios (CHRs) using their correlation maximizing lag for for comparison (Figure 6).

For each state, the CHRs tend to be larger and noiser relative to IHRs. 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 Ohio, Pennsylvania, and Virginia) tend to exhibit similar patterns in the IHRs and CHRs over time. In addition, there are similar spikes observed across many states during waves

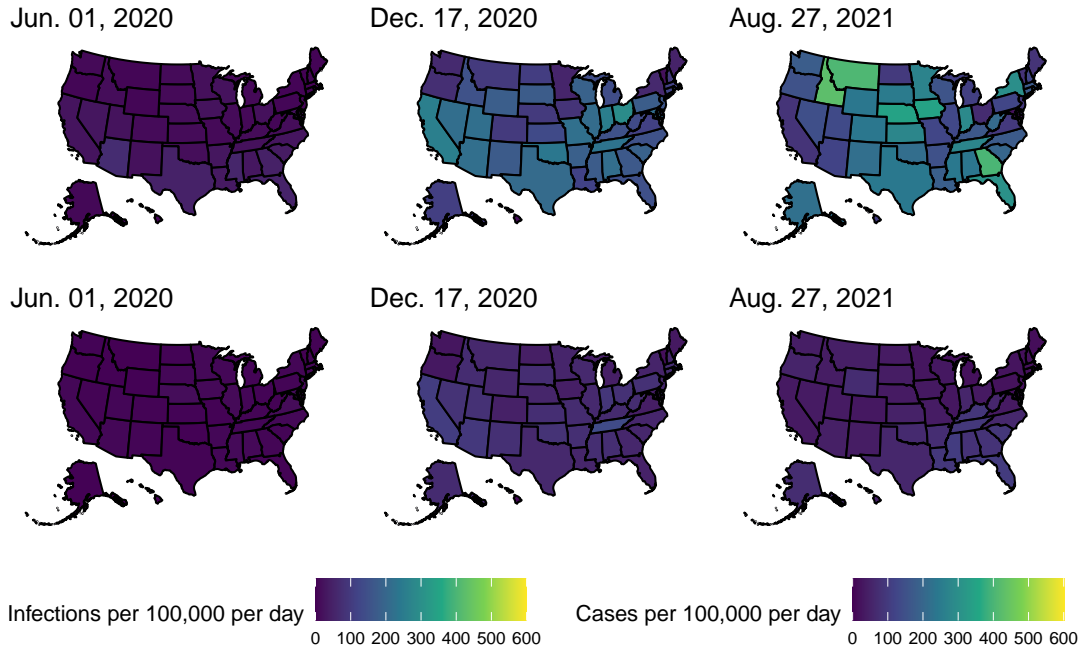


Figure 3: Choropleth maps of the state-level estimates of the number of daily new infections per 100,000 population (top row) and the daily new cases per 100,000 population (bottom row) for three times over the June 1, 2020 to November 29, 2021 period. The first date was chosen simply as a baseline, while the second and third dates were chosen based on the day that had the largest number of infections across the 50 states from each year.

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<sup>13</sup>. This finding aligns with previous studies that found an increased risk in hospitalizations with Delta in comparison to other variants<sup>14;15</sup>. Similarly, during the fall of 2020 there tends to be another spike in the IHRs that rivals or surpasses that observed during the time of Delta (which is the case for states like New York or Wyoming).

There does not tend to be a strict upward or downward trajectory or even a mild waning pattern in the IHRs, as one might expect with later variants that are more infectious but result in fewer hospitalization<sup>16;17</sup>. Overall, we observe intermittent spikes that punctuate longer periods where the IHRs tend to stabilize slightly below 0.2 hospitalizations per infection. These spikes tend to align with the emergence of new variants. **ATTN: 0.2 seems wrong. Is the IHR actually H/I? or are the units not quite right (e.g., H/1M / I/100K)?**

While we computed and compared CHR and IHRs for all states, it is important to note that both likely to vary within states and depend on confounding variables such as age and the presence of major comorbidities<sup>18</sup>. Therefore, it would be beneficial to account for such variables in their calculations by, for example, stratifying infections and hospitalizations by age to produce age-specific estimates of the IHRs for each state<sup>19</sup>.

## 2.7 Disease burden and viral transmission

**ATTN: I'm not yet sure what to do with this. Some of it should go into one of the sections above.**

From reconstructing the time series of COVID-19 infections per 100,000 population for each U.S. state from June 1, 2020 to November 29, 2021, we observe rates of infections that vary in intensity and disease burden across space and time (Figure 2, Figure 4). Most states present at least two major spikes in infections - the first starts in the fall of 2020 and extends into the winter season, while the second starts in the late

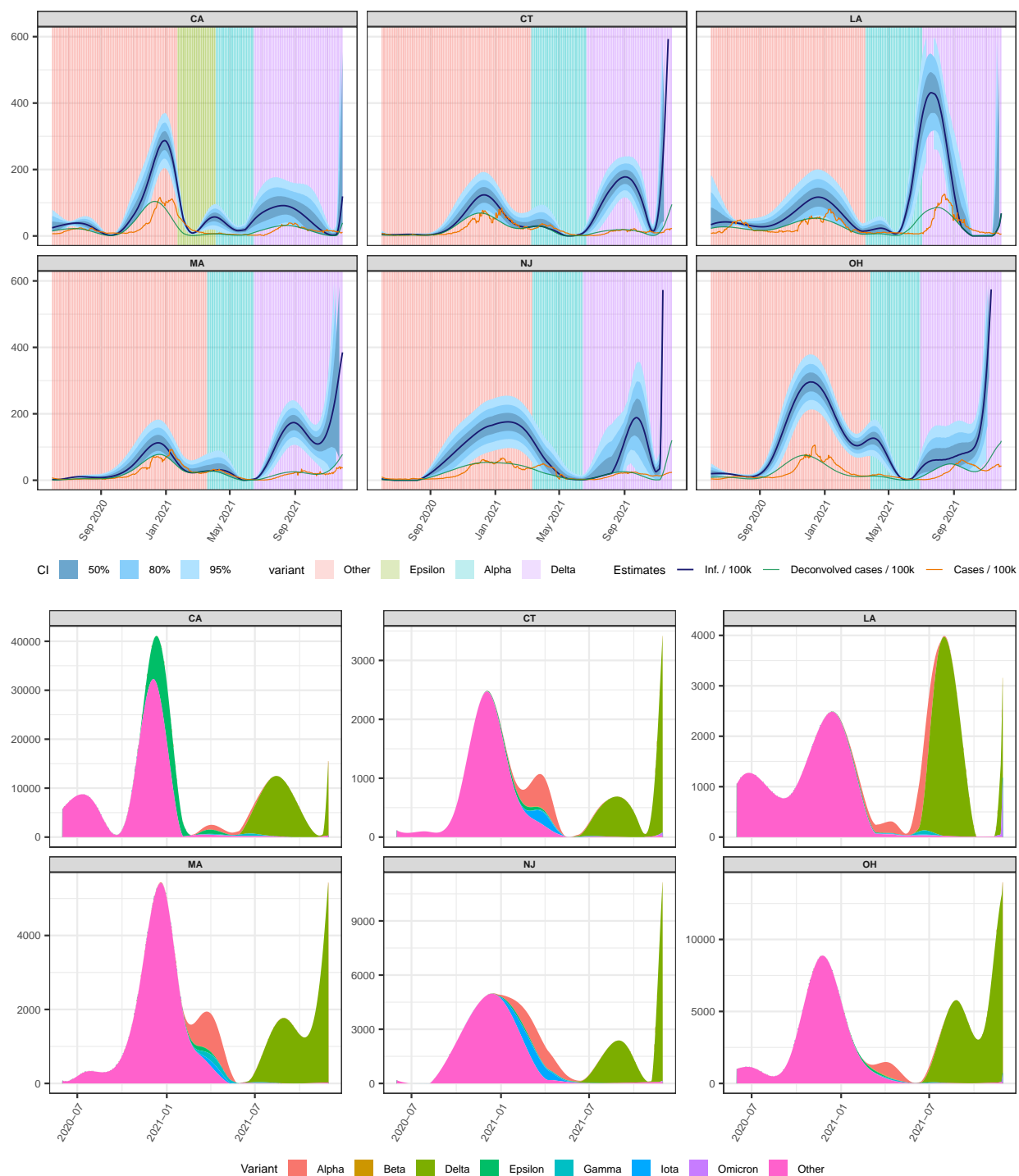


Figure 4: Top panel: Reported cases, deconvolved cases, and estimates of daily new infections (dark blue line) per 100K inhabitants. The blue shaded regions indicate the 50, 80, and 95% confidence bands, while the background is shaded to indicate the dominant variant in circulation at the time. **ATTN: I don't like the rotated x-axis.** Bottom panel: Deconvolved cases colored by variant per 100K inhabitants. **ATTN: Is this the correct units? or are these raw? They should be per 100K.**

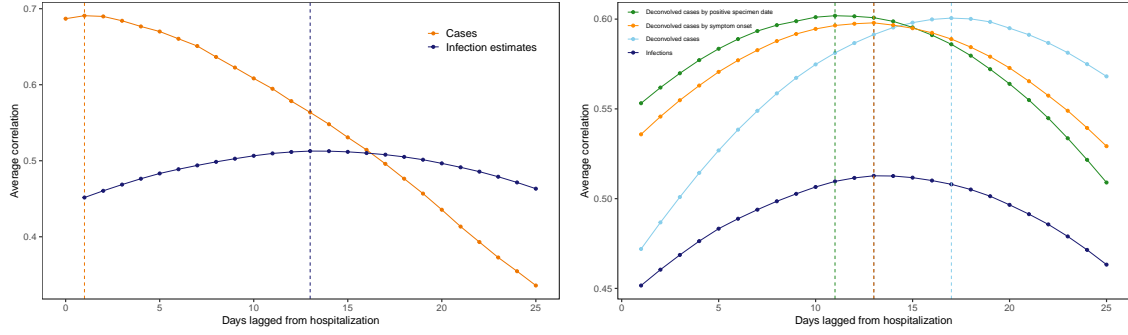


Figure 5: Spearman's correlation between the different case/infection rates and hospitalization rates per 100,000. These are calculated for each lag, state and rolling window of 61 days before averaging. The vertical dashed lines indicate the lags for which the highest average correlation is attained. **ATTN: Let's do one panel with cases, deconvolved cases, and infections.**

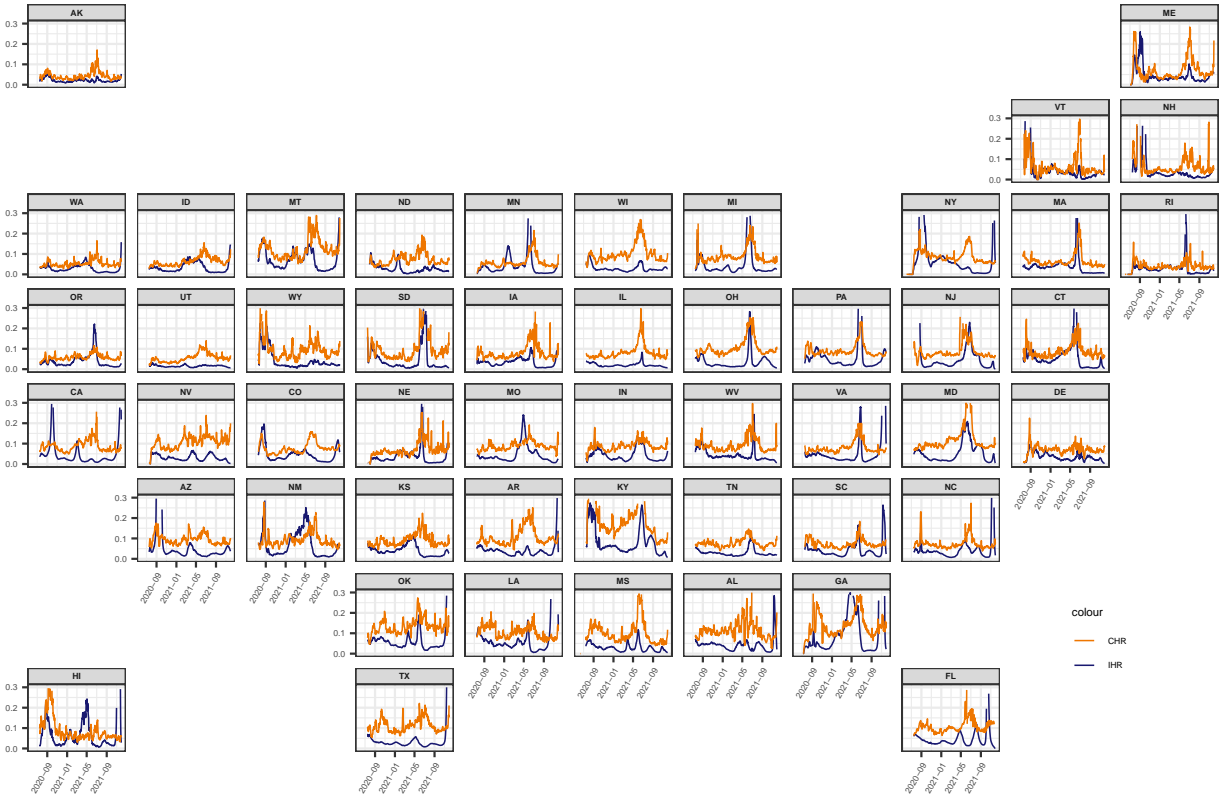


Figure 6: Time-varying IHR and CHR estimates for each state from June 1, 2020, to November 29, 2021, obtained using the corresponding optimal lag from the systematic lag analysis. Note that the infection, case, and hospitalization counts are subject to a center-aligned 7-day average to remove spurious day of the week effects. Also note that the different starting points across states are due to the availability of the hospitalization data.

summer of 2021 and proceeds into the mid-fall. These represent major waves driven by the Ancestral and Delta 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 found to emerge in the clusters of states that border each other.



To avoid encroaching upon possible boundary issues with ending the estimation during a time of volatility (the period of the Delta-Omicron transition), we focus on the infection estimates prior to November 1, 2021. The largest observed outbreaks prior to this time were observed in the late summer or early fall of 2021 in Georgia, Louisiana, Idaho, Montana, and Wyoming which suggests a similar spread of the virus in small clusters of states that are in close geographic proximity. During this time, the two states that have the highest rate of infections per 100,000 on single day are Georgia with about 451 infections per 100,000 on August 15, 2021 (95% confidence interval: [334, 567]) and Idaho with 451 on September 7, 2021 (95% confidence interval: [312, 590]). These are closely followed by Montana with 432 on September 8, 2021 (95% confidence interval: [282, 581]), Louisiana with 431 on July 20, 2021 (95% confidence interval: [252, 610]), and Wyoming with 350 on November 13, 2020 (95% confidence interval: [256, 444]).

Prior to the Delta wave, the state that has the highest rate of infections per 100,000 on single day is Louisiana with about 358 infections per 100,000 on July 3, 2021 (95% confidence interval: [177, 539]), followed by Wyoming with 349 on November 13, 2020 (95% confidence interval: [407, 546]), South Dakota with 342 infections per 100,000 on July 3, 2021 (95% confidence interval: [177, 539]), and Illinois with 340 infections per 100,000 on July 3, 2021 (95% confidence interval: [177, 539]). During this time, 74% of the top rates for each state were observed in the late fall or winter of 2020.

The period of lowest viral transmission is observed in the summer and fall of 2020. During this time, the state of New Hampshire achieves the lowest weekly rate of infections of 0.01 infections per 100,000 for the week of September 13, 2020. In the summer of 2020, Vermont maintains a rate under 10 infections per 100,000 from the week of June 1, 2020 to August 30, 2020, which is the longest continuous stretch observed for any state.

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 Idaho and Montana 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 Delta 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.

### 3 Discussion

We retrospectively estimated daily incident infections for each U.S. state over the period June 1, 2020 to November 29, 2021. Our estimates suggest both (a) that the pandemic impacted states earlier and at a larger scale than is indicated by cases and that (b) examining cases alone hides some spatio-temporal waves that become apparent by examining infections. We observe outbreaks in infections that are difficult to detect from cases alone such as the Delta wave in New Jersey, Connecticut, and Maryland. This suggests that cases paint an incomplete picture of the pandemic, especially when outbreaks are largely driven by unreported infections. Furthermore, since case reports generally follow symptom and infection onsets, cases 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<sup>9;20</sup>). Thus, while reported cases provide an indication of the trajectory of the pandemic, it is a delayed and incomplete version.

Our approach offers a number of advantages. For instance, we aim to incorporate as much state-specific information as possible when deriving our estimates. By using state-level case, line list, and variant circulation data, we are able to construct incubation and delay distributions that are specific to each state. Time-varying and state-specific seroprevalence data allows the reporting ratio estimates to similarly vary over space and time, a departure from existing work<sup>21;22</sup>. Existing approaches that use the delay distribution to generate infection estimates often only construct one delay distribution that is used for all states<sup>23;24</sup>. That is, our work avoids the assumption of geographic invariance, where it is assumed that all states have the same patterns of delay from symptom onset to case report. This assumption is unlikely to be true due to differences in reporting pipelines, pandemic response, and variants in circulation, among other issues.

Another limitation of previous approaches to estimate latent infections is that they do not account



for reinfections. While reinfections represent a small fraction of total infections until later in the pandemic, ignoring them means that the infection-reporting ratio will tend to be underestimated with seroprevalence data alone. By accounting for these as well as the waning of seropositivity (See [Section 4.5](#)), we more accurately estimate this ratio. 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 immunity is likely to be variant-dependent<sup>25</sup>, it follows that our model waning 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 that newer variants may escape detection<sup>26;27</sup>. 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.

Regarding reinfections, a major reason why we chose an end date of November 29, 2021 and ultimately decided to not tread into Omicron territory is because the Omicron variants come with substantial increase in the risk of reinfection in comparison to previous variants as Omicron has been shown to have an increased tendency towards immune escape<sup>28–30</sup>. So having quality reinfection data that is representative of each location under study is of the utmost importance for the Omicron era.

Using seroprevalence data to estimate the case-ascertainment ratio is subject to a number of issues, and precludes us from pushing the period of analysis past the Omicron wave in December 2021. While most state-level data suggests that reinfections still account for less than 20% of reported cases during Omicron [ATTN: cites](#), seropositivity rapidly reaches nearly 100% of the population, precluding its continued use. Due to these issues, alternative data sources for estimating the case-ascertainment ratio is necessary. For example, wastewater surveillance data is may be complementary to seroprevalence data, especially when testing is low<sup>31</sup>. However, viral detection is inconsistent across locations due to temperature, per-capita water use, and in-sewer travel time<sup>31–33</sup>. 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<sup>6;34</sup>.

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 is the development of a way of to model the waning of detectable antibody levels 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 quality and the quantity 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.

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.

## 4 Methods

In what follows, we provide details on how we estimate the daily incident infections for each state over the considered time period of June 1, 2020 to November 29, 2021 and the data we used to achieve this. [Figure 7](#) provides a visual summary of the data, analysis tasks, and the relationships between them. The major analysis tasks this figure aims to convey are as follows: First, we estimate variant-specific incubation periods and two types of delay distributions for each day over the considered time period. Next, each incubation period and symptom onset to positive specimen delay distribution are joined using convolution to obtain variant-specific infection onset to positive specimen distributions for each time. Then two types

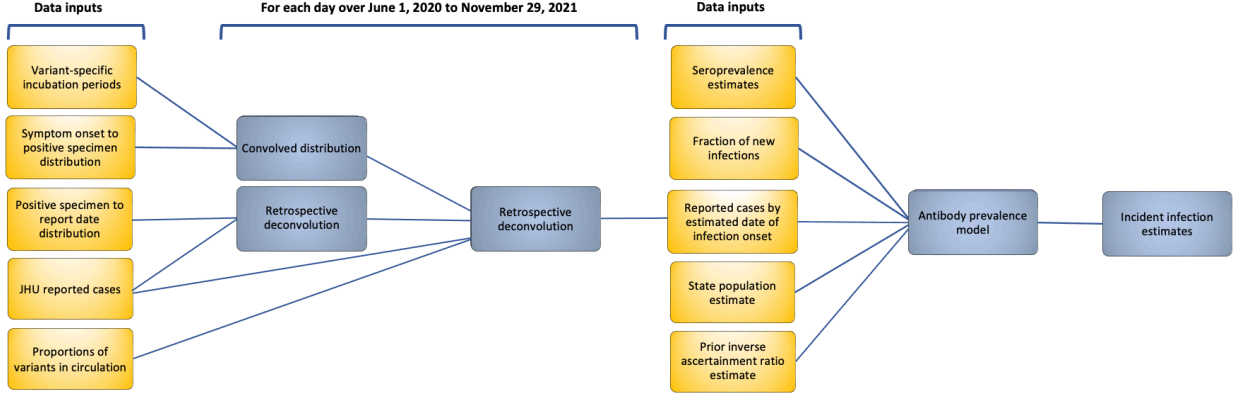


Figure 7: Flowchart of the inputted data and major analysis steps required to get from reported cases to incident infection estimates for each day over June 1, 2020 to November 29, 2021 for a state. Data sources are coloured in yellow, while data analysis steps are coloured in blue. The data sources that do not stem from an analysis step are literature estimates.

of deconvolution are performed. We first deconvolve from case report to positive specimen date. We then deconvolve from positive specimen to report date by variant. The resulting infection estimates are aggregated across the variant categories, and adjusted to account for the unreported infections by using state-specific, time-varying seroprevalence data in an antibody prevalence model. This lets us reach our ultimate goal of obtaining daily incident infection estimates.

#### 4.1 Estimating delay distributions from private line lists

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<sup>35;36</sup>, only the restricted dataset contains information on the state of residence. The three key dates of interest are those for symptom onset, positive specimen collection, and report to the CDC. Handling missingness and imputation in these dates is somewhat complicated, and additional details and justifications are deferred to [Section S1.6](#).

We use the line list to estimate the delay distribution for the pairs symptom onset to positive specimen and positive specimen to report. We provide the full procedure for the latter, before giving a brief description below for the former. First, define  $z_{\ell,t}$  to be a case report occurring at time  $t$  in location  $\ell$ , and let  $\pi_{\ell,t}(k)$  to be the probability that  $z_{\ell,t}$  has a positive specimen collected  $k$  days earlier. We assume that all positive specimens will be reported within 60 days and that no test will be reported on the same date as it was collected, that is,  $\pi_{\ell,t}(0) = 0$  and  $\pi_{\ell,t}(k) = 0$  whenever  $k > 60$ . Let  $N_{\ell,t}$  be the number of  $z_{\ell,s}$  with  $s \in [t - 75 + 1, t + 60] = \mathcal{S}_t$  and positive specimen date greater than  $s - 60$ . Then, we first compute

$$\tilde{p}_{\ell,t}(k) = \frac{1}{N_{\ell,t}} \sum_{s \in \mathcal{S}_t} (\# z_{\ell,s} \text{ with positive specimen at } s - k). \quad (1)$$

Next we compute a similar national quantity  $\tilde{p}_t(k) = \frac{1}{N_t} \sum_{s \in \mathcal{S}_t} (\# z_s \text{ with positive specimen at } s - k)$ , without restricting to location  $\ell$ . Next, let  $\alpha_{\ell,t}$  be the ratio of  $N_{\ell,t}$  to the number of cases reported by JHU CSSE<sup>1</sup> in the same window. Then, compute  $p_{\ell,t}(k) = \alpha_{\ell,t} \tilde{p}_{\ell,t}(k) + (1 - \alpha_{\ell,t}) \tilde{p}_t(k)$ . This construction allows for more reliance on the state estimate when there are more CDC cases relative to JHU (and vice versa). We calculate the mean  $m_{\ell,t}$  and variance  $v_{\ell,t}$  of  $\{p_{\ell,t}(k) : 0 < k \leq 60\}$  and estimate a gamma distribution by solving the moment equations  $m_{\ell,t} = \alpha_{\ell,t} \theta_{\ell,t}$  and  $v_{\ell,t} = \alpha_{\ell,t} \theta_{\ell,t}^2$  for the shape  $\alpha_{\ell,t}$  and scale  $\theta_{\ell,t}$ . Finally, we discretize the resulting gamma density to the support set of 1 to 60 days to produce an estimate  $\{\hat{\pi}_{\ell,t}(k) : 0 < k \leq 60\}$  of the delay distribution  $\pi_{\ell,t}$ .

Estimating the delay from symptom onset to positive specimen date follows the same procedure with a few minor adjustments. First, we allow  $k$  to range from  $-3$  to  $21$  (rather than  $1$  to  $60$ ). These upper

and lower bounds are based on the largest delay values for the state-wide 0.05 and 0.95 quantiles. This is reasonable because the median delay is very short at approximately 2 days, and an asymptomatic individual may test positive following a known exposure, before the onset of symptoms. Additional minor details are discussed in [Section S1.7](#).

## 4.2 Estimating the incubation period distributions

To model the incubation period, the time between infection and symptom onset, we use estimates from the existing literature, modified slightly to for coherence with each other: we model each incubation as a gamma distribution with different parameters. We focus on the following eight variants, which dominated at various points during our study period: Ancestral, Alpha, Beta, Epsilon, Iota, Gamma, Delta, and Omicron. Alpha, Beta, Delta, Gamma, and Omicron are all variants of concern<sup>37</sup>, while we include the Epsilon (California) and Iota (New York) variants because of large impact on those and neighbouring states<sup>38;39</sup>.

The Ancestral variant has been modelled as a gamma distribution<sup>40</sup>, so we simply use those reported parameters. For the Alpha, Beta, Gamma, Delta and Omicron variants, we use the reported mean and standard deviation of the number of days of incubation<sup>41–43</sup>. To match these moments to the gamma distribution, we solve the same moment equations described in [Section 4.1](#). Then, we discretize each resulting density to the support set, which is taken to be from 1 and 21 days. This range assumes that symptoms require at least 1 day to develop<sup>44</sup> and that an asymptomatic infection will resolve within 21 days<sup>45;46</sup>.

We were unable to locate incubation period estimates for the geo-specific Epsilon and Iota variants, so 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<sup>13</sup>. All other circulating variants are grouped together with the Ancestral variant. There was little available sequencing data prior to Alpha-emergence, but unfortunately, later in the pandemic, it is impossible to separate Ancestral from other rare variants.

**ATTN: Perhaps plot the incubation periods? Maybe sharing a panel with the circulation proportions below.**

## 4.3 Variant circulation proportions

To estimate the daily proportions of the variants circulating in each state, we obtain the GISAID genomic sequencing data from CoVariants.org<sup>13;47</sup>. These counts represent the total number of cases belonging to a particular variant using a sample of positive tests over a biweekly period. To estimate the population proportion of each variant, we apply multinomial logistic regression for the eight variant categories separately for each state.

We let  $V_{j\ell,t}$  to be the probability of a new cases at time  $t$  in location  $\ell$  corresponding to variant  $j$ . Let  $v_{j\ell,t}$  be the analogous observed proportion. Then the nonparametric multinomial logistic regression model is given as the system

$$\log\left(\frac{V_{j\ell,t}}{1 - V_{j\ell,t}}\right) = f_{j\ell}(t), \quad j = 1, \dots, J, \quad \text{subject to} \quad \sum_{j=1}^J \exp\{f_{j\ell}(t)\} = 1, \quad \forall t. \quad (2)$$

The constraint ensures that the estimated proportions will sum to 1 across all  $J$  variants. To encourage smoothness of the estimated proportions, we specify  $f_{j\ell}(t)$  as a third-order polynomial in time: that is  $f_{j\ell}(t) = \beta_{j\ell,0} + \beta_{j\ell,1}t + \beta_{j\ell,2}t^2 + \beta_{j\ell,3}t^3$ , computed such that the resulting matrix of covariates is orthogonal.

**ATTN: show the result for one state: the original proportions on the left, and the smoothed on the right. Then describe this to a small degree (2-3 sentences).**

## 4.4 Retrospective deconvolution: from cases to infections

Retrospective deconvolution estimates the daily number of new infections corresponding to each variant for each time and location, “pushing back” the dates that those cases were eventually reported to the time of infection. Because the circulating variant proportions in [Section 4.3](#) correspond to the positive specimen date, this requires two stages. The first is the deconvolution from report to positive specimen date, and the second is from positive specimen date to infection onset date.

We will start by describing the first type of deconvolution performed from report to positive specimen date in detail. For this problem, let  $t = 1, \dots, T'$  index the extended deconvolution period from March 1, 2020 to March 1, 2023, extended to minimize the effects of boundary issues. Define  $y_{\ell,t}$  to be the number of new cases reported in location  $\ell$  at time  $t$ , as reported by the John Hopkins Center for Systems Science and Engineering (JHU CSSE)<sup>1</sup> and retrieved with the The COVIDcast API<sup>34</sup>. Recall that  $\hat{\pi}_{\ell,t}(k)$  is the associated probability that these reported cases were collected  $k$  days earlier.

We estimate the deconvolved cases by positive specimen date by solving the following optimization problem:

$$\underset{\mathbf{x}}{\text{minimize}} \sum_{t=1}^{T'} \left( y_{\ell,t} - \sum_{k=1}^{60} \hat{\pi}_{\ell,t}(k) x_{t-k} \right)^2 + \lambda \sum_{t=4}^{T'} |x_t - 4x_{t-1} + 6x_{t-2} - 4x_{t-3} + x_{t-4}|. \quad (3)$$

**ATTN:** These sums aren't quite right: the outer  $t$  makes  $x_{t-k}$  run negative? Actually this wouldn't be a problem, except that we don't deconvolve before  $t = 1$ , we renormalize the delays. Perhaps too much detail, but the math annoys me. Ideas? The two parts of this optimization problem trade data fidelity (the sum of squared errors) with smoothness in the resulting estimates (the absolute error of the differences of  $\mathbf{x}$ ). The tuning parameter  $\lambda$  determines the relative importance of these competing goals. The solution to the problem is an adaptive piecewise cubic polynomial<sup>48;49</sup> and can be accurately computed easily<sup>50</sup>. We select  $\lambda$  with 3-fold cross validation<sup>24</sup> in which every third day is reserved for testing, and the value that results in the smallest out-of-sample mean squared error is chosen.

The result of this first deconvolution is  $\hat{x}_{\ell,t}$ , case estimates by positive specimen date for each state. To continue, pushing back to infection estimate, we first need the variant-specific delays from infection to positive specimen collection. These are calculated by convolving the location-time-specific symptom-to-test distributions from Section 4.1, denoted by  $\{q_{\ell,t}(k) : -3 \leq k \leq 21\}$ , with the variant-specific incubation periods from Section 4.2, denoted by  $\{i_j(k) : 0 < k \leq 21\}$ . The convolution of these yields a distribution  $\mathbf{q}_{\ell,t} * \mathbf{i}_j = \{\tau_{j\ell,t}(k) : -3 \leq k \leq 42\}$ . However, only a fraction of  $\hat{x}_{\ell,t}$  corresponds to each variant, so we must weight them by the variant proportions  $\hat{v}_{j\ell,t}$  estimated in Section 4.3. The analogous optimization problem is therefore:

$$\underset{\mathbf{u}}{\text{minimize}} \sum_{t=1}^{T'} \left( \hat{v}_{j\ell,t} \hat{x}_{\ell,t} - \sum_{k=-3}^{42} \tau_{j\ell,t}(k) u_{t-k} \right)^2 + \lambda \sum_{t=4}^{T'} |u_t - 4u_{t-1} + 6u_{t-2} - 4u_{t-3} + u_{t-4}|. \quad (4)$$

We call the solution  $\tilde{\mathbf{u}}_{j\ell}$  the *variant-specific deconvolved cases* and emphasize that these are those cases that will eventually be reported to public health. Because this deconvolution is done separately for each location and variant category, we ultimately obtain deconvolved case estimates by the date of infection onset that are separated by variant. Finally, we will denote the total deconvolved cases at location  $\ell$  as  $\hat{\mathbf{u}}_{\ell} = \sum_j \tilde{\mathbf{u}}_{j\ell}$ .

## 4.5 Inverse reporting ratio and the antibody prevalence model

To capture the unreported infections, it is necessary to adjust these deconvolved case estimates by the ratio of the true number of new infections to the new reported infections.

Because seroprevalence of anti-nucleocapsid antibodies represents the percentage of people who have at least one resolving or past infection<sup>51</sup>, we can use the change in subsequent seroprevalence measurements to estimate *all* new infections, rather than just those eventually appearing as cases. This intuition suggests modelling reported seroprevalence at time  $t+1$  as a fraction  $1 - \gamma$  of the previous seroprevalence measurement at  $t$  plus the reinfection-adjusted deconvolved cases multiplied by the inverse reporting ratio at time  $t$ :

$$s_{\ell,t+1} = (1 - \gamma_{\ell}) s_t + a_{\ell,t} z_{\ell,t} \hat{u}_{\ell,t} + \epsilon_{\ell,t}, \quad (5)$$

where  $\hat{u}_{\ell,t}$  is deconvolved cases,  $z_{\ell,t}$  is the percentage of reinfections,  $a_{\ell,t}$  is the inverse reporting ratio, and  $\epsilon_{\ell,t}$  represents noise. Note that  $\gamma_{\ell}$  is the percentage of people whose level of infection-induced antibodies falls below the detection threshold between time  $t$  and time  $t+1$ . Informally, we refer to  $\gamma$  as the waning parameter. Unfortunately, population seroprevalence is not directly observed, but is measured using various surveys.

For the proportion of the population in each state with evidence of previous infection across time, we use two major seroprevalence surveys: the 2020–2021 Blood Donor Seroprevalence Survey and the Nationwide Commercial Lab Seroprevalence Survey<sup>52;53</sup>. See [Section S1.8](#) for additional details. Each of these provides seroprevalence estimates along with confidence intervals. The daily fraction of new infections are based on surveillance work conducted by the Southern Nevada Health District<sup>54</sup>. **ATTN: These results are broadly similar to those in other locations with available data...**

In order to account for different surveys occurring on different dates with noisy estimates, we estimate the model on the weekly frequency, observed on Monday, and treat  $s_{\ell,t}$  as a latent variable. Therefore, we write,

$$r_{\ell,m}^1 = s_{\ell,m} + \tau_{\ell,m}, \quad \tau_{\ell,m} \sim \mathcal{N}(0, w_{\ell,m}^1 \sigma_{\ell,r}^2), \quad (6)$$

$$r_{\ell,m}^2 = s_{\ell,m} + \varphi_{\ell,m}, \quad \varphi_{\ell,m} \sim \mathcal{N}(0, w_{\ell,m}^2 \sigma_{\ell,r}^2), \quad (7)$$

$$s_{\ell,m+1} = (1 - \gamma_{\ell})s_{\ell,m} + a_{\ell,m}z_{\ell,m}\hat{u}_{\ell,m}^{\Sigma} + \epsilon_{\ell,m}, \quad \epsilon_{\ell,m} \sim \mathcal{N}(0, \sigma_{\ell,\epsilon}^2), \quad (8)$$

where  $r^1$  and  $r^2$  correspond to the two different seroprevalence surveys. These surveys each have measurement errors with variance  $\sigma_r^2$  that scale proportional to the observed confidence intervals for the estimates, respectively  $w_{\ell,m}^1$  and  $w_{\ell,m}^2$ . We denote  $\hat{u}_{\ell,m}^{\Sigma} = \sum_{t=m}^{m+1} \hat{u}_{\ell,t}$ . Finally, to ensure that  $\mathbf{a}_{\ell}$  is smooth over time, we complete the model with an additional equation that enforces smoothness,

$$a_{\ell,m+1} = 3a_{\ell,m} - 3a_{\ell,m-1} + a_{\ell,m-2} + \eta_{\ell,m}, \quad \eta_{\ell,m} \sim \mathcal{N}(0, \sigma_{\eta}^2). \quad (9)$$

This antibody prevalence model is a state-space model with latent variables  $\mathbf{s}_{\ell}$  and  $\mathbf{a}_{\ell}$  and unknown parameters  $\gamma_{\ell}$ ,  $\sigma_r^2$ ,  $\sigma_{\epsilon}^2$ , and  $\sigma_{\eta}^2$ . This model allows for convenient handling of missing data, extrapolation before and after the period of observed seroprevalence measurements, and maximum likelihood estimates of the errors. Details of this methodology and the computation of the associated uncertainty measurements are deferred to [Section S1.2](#).

## 4.6 Lagged correlation to hospitalizations and time-varying IHRs

From The COVIDcast API<sup>34</sup>, we 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). We use our infection estimates  $\hat{\mathbf{u}}_{\ell}$  to compute the lagged correlation with hospitalizations. The goal of this analysis is to find the lag between infection and hospitalization rates that gives the highest average rank-based correlation across U.S. 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. To remove day of the week effects, both the infection and hospitalization signals are averaged over a 7-day, center-aligned, moving window before their conversion to rates.

For each considered lag, we calculate Spearman’s correlation between the state infection and hospitalization rates for each observed between June 1, 2020 to November 29, 2021 with a center-aligned rolling window of 61 days. We then average these correlations 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. The IHR is computed by dividing the number of individuals who are hospitalized due to COVID-19 by the estimated total number who were infected on the lagged number of days before. **ATTN: Isn’t there a window here?**

## Data availability

## Code availability

## References

- [1] Ensheng Dong, Hongru Du, and Lauren Gardner. An interactive web-based dashboard to track COVID-19 in real time. *The Lancet Infectious Diseases*, 20(5):533–534, 2020.
- [2] The New York Times. Coronavirus in the U.S.: Latest map and case count. <https://www.nytimes.com/interactive/2021/us/covid-cases.html>, 2020.
- [3] The Washington Post. 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>, 2020.
- [4] Centers for Disease Control and Prevention. Estimated COVID-19 burden. <https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/burden.html>, 2022.
- [5] Virginia E. Pitzer, Melanie Chitwood, Joshua Havumaki, Nicolas A. Menzies, Stephanie Perniciaro, Joshua L. Warren, Daniel M. Weinberger, and Ted Cohen. The impact of changes in diagnostic testing practices on estimates of COVID-19 transmission in the United States. *American Journal of Epidemiology*, 190(9):1908–1917, 2021.
- [6] European Centre for Disease Prevention and Control. Strategies for the surveillance of COVID-19. Technical report, ECDC, Stockholm, Sweden, 2020.
- [7] Matt D.T. Hitchings, Natalie E Dean, Bernardo García-Carreras, Thomas J. Hladish, Angkana T. Huang, Bingyi Yang, and Derek A.T. Cummings. The usefulness of the test-positive proportion of severe acute respiratory syndrome coronavirus 2 as a surveillance tool. *American Journal of Epidemiology*, 190(7):1396–1405, 2021.
- [8] Lorenzo Pellis, Francesca Scarabel, Helena B. Stage, Christopher E. Overton, Lauren H.K. Chappell, Elizabeth Fearon, Emma Bennett, Katrina A. Lythgoe, Thomas A. House, Ian Hall, et al. Challenges in control of COVID-19: Short doubling time and long delay to effect of interventions. *Philosophical Transactions of the Royal Society B*, 376(1829):20200264, 2021.
- [9] Washington State Department of Health. COVID-19 data dashboard. <https://doh.wa.gov/emergencies/covid-19/data-dashboard>, 2020.
- [10] Ontario Agency for Health Protection and Promotion. 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](https://www.publichealthontario.ca/-/media/documents/ncov/voc/2022/01/covid-19-omicron-b11529-risk-assessment-jan-6.pdf?sc_lang=en), 2022.
- [11] Nigel Garrett, Asa Tapley, Jessica Andriesen, Ishen Seocharan, Leigh H. Fisher, Lisa Bunts, Nicole Espy, Carole L. Wallis, April Kaur Randhawa, Nzeera Ketter, et al. High rate of asymptomatic carriage associated with variant strain Omicron. *MedRxiv*, 2022.
- [12] Thomas Ward and Alexander Johnsen. Understanding an evolving pandemic: An analysis of the clinical time delay distributions of COVID-19 in the United Kingdom. *PLoS One*, 16(10):e0257978, 2021.
- [13] Emma Hodcroft. CoVariants: SARS-CoV-2 mutations and variants of interest. <https://covariants.org>, 2021.
- [14] Katherine A. Twohig, Tommy Nyberg, Asad Zaidi, Simon Thelwall, Mary A. Sinnathamby, Shirin Aliabadi, Shaun R. Seaman, Ross J. Harris, Russell Hope, Jamie Lopez-Bernal, et al. 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(1):35–42, 2022.
- [15] Tommy Nyberg, Neil M. Ferguson, Sophie G. Nash, Harriet H. Webster, Seth Flaxman, Nick Andrews, Wes Hinsley, Jamie Lopez Bernal, Meaghan Kall, Samir Bhatt, et al. 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(10332):1303–1312, 2022.



- [16] Ramon Lorenzo-Redondo, Egon A. Ozer, and Judd F. Hultquist. COVID-19: Is Omicron less lethal than Delta? *British Medical Journal*, 378, 2022.
- [17] Beth Blauer. 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>, 2022.
- [18] Clark D. Russell, Nazir I. Lone, and J. Kenneth Baillie. Comorbidities, multimorbidity and COVID-19. *Nature Medicine*, 29(2):334–343, 2023.
- [19] Spencer J. Fox, Emily Javan, Remy Pasco, Graham C. Gibson, Briana Betke, José L. Herrera-Diestra, Spencer Woody, Kelly Pierce, Kaitlyn E. Johnson, Maureen Johnson-León, et al. Disproportionate impacts of COVID-19 in a large US city. *PLOS Computational Biology*, 19(6):e1011149, 2023.
- [20] Stephanie Dunkel. COVID-19 case numbers: Why the delay in reporting? <https://www.tpchd.org/Home/Components/Blog/Blog/21448>, 2020.
- [21] H. Juliette T. Unwin, Swapnil Mishra, Valerie C. Bradley, Axel Gandy, Thomas A. Mellan, Helen Coupland, Jonathan Ish-Horowicz, Michaela A.C. Vollmer, Charles Whittaker, Sarah L. Filippi, et al. State-level tracking of COVID-19 in the United States. *Nature Communications*, 11(1):6189, 2020.
- [22] Center for the Ecology of Infection Diseases. COVID-19 portal. <https://www.covid19.uga.edu/nowcast.html>, 2020.
- [23] Melanie H. Chitwood, Marcus Russi, Kenneth Gunasekera, Joshua Havumaki, Fayette Klaassen, Virginia E Pitzer, Joshua A. Salomon, Nicole A. Swartwood, Joshua L. Warren, Daniel M. Weinberger, et al. 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(8):e1010465, 2022.
- [24] Maria Jahja, Andrew Chin, and Ryan J. Tibshirani. Real-time estimation of COVID-19 infections: Deconvolution and sensor fusion. *Statistical Science*, 37(2):207–228, 2022.
- [25] Nick Pooley, Salim S. Abdool Karim, Behazine Combadière, Eng Eong Ooi, Rebecca C. Harris, Clotilde El Guerche Seblain, Masoumeh Kisomi, and Nabila Shaikh. Durability of vaccine-induced and natural immunity against COVID-19: A narrative review. *Infectious Diseases and Therapy*, 12(2):367–387, 2023.
- [26] National Institutes of Health. 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>, 2022.
- [27] U.S. Food and Drug Administration. 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>, 2023.
- [28] Jia Wei, Nicole Stoesser, Philippa C Matthews, Tarnjit Khara, Owen Gethings, Ian Diamond, Ruth Studley, Nick Taylor, Tim EA Peto, A Sarah Walker, et al. Risk of sars-cov-2 reinfection during multiple omicron variant waves in the uk general population. *Nature Communications*, 15(1):1008, 2024.
- [29] Juliet RC Pulliam, Cari van Schalkwyk, Nevashan Govender, Anne von Gottberg, Cheryl Cohen, Michelle J Groome, Jonathan Dushoff, Koleka Mlisana, and Harry Moultrie. Increased risk of sars-cov-2 reinfection associated with emergence of omicron in south africa. *Science*, 376(6593):eabn4947, 2022.
- [30] Elias Eythorsson, Hrafnhildur Linnet Runolfsdottir, Ragnar Freyr Ingvarsson, Martin I Sigurdsson, and Runolfur Palsson. Rate of sars-cov-2 reinfection during an omicron wave in iceland. *JAMA Network Open*, 5(8):e2225320–e2225320, 2022.
- [31] Oliver McManus, Lasse Engbo Christiansen, Maarten Nauta, Lene Wulff Krogsgaard, Naja Stolberg Bahrenscheer, Lene von Kappelgaard, Tobias Christiansen, Mikkel Hansen, Nicco Claudio Hansen, Jonas Kähler, et al. Predicting COVID-19 incidence using wastewater surveillance data, Denmark, October 2021–June 2022. *Emerging Infectious Diseases*, 29(8):1589, 2023.

- [32] Olga E. Hart and Rolf U. Halden. 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, 2020.
- [33] Xuan Li, Shuxin Zhang, Samendra Sherchan, Gorka Orive, Unax Lertxundi, Eiji Haramoto, Ryo Honda, Manish Kumar, Sudipti Arora, Masaaki Kitajima, et al. 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, 2023.
- [34] Alex Reinhart, Logan Brooks, Maria Jahja, Aaron Rumack, Jingjing Tang, Sumit Agrawal, Wael Al Saeed, Taylor Arnold, Amartya Basu, Jacob Bien, et al. An open repository of real-time COVID-19 indicators. *Proceedings of the National Academy of Sciences*, 118(51):e2111452118, 2021.
- [35] Centers for Disease Control and Prevention. COVID-19 case surveillance public use data. <https://data.cdc.gov/Case-Surveillance/COVID-19-Case-Surveillance-Public-Use-Data/vbim-akqf>, 2020.
- [36] Centers for Disease Control and Prevention. COVID-19 case surveillance restricted access detailed data. <https://data.cdc.gov/Case-Surveillance/COVID-19-Case-Surveillance-Restricted-Access-Detail/mbd7-r32t>, 2020.
- [37] World Health Organization. Tracking SARS-CoV-2 variants. <https://www.who.int/activities/tracking-SARS-CoV-2-variants>, 2021.
- [38] Shangxin Yang, Peera Hemarajata, Evann E. Hilt, Travis K. Price, Omai B. Garner, and Nicole M. Green. 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(5):649–652, 2022.
- [39] Ralf Duerr, Dacia Dimartino, Christian Marier, Paul Zappile, Guiqing Wang, Jennifer Lighter, Brian Elbel, Andrea B Troxel, Adriana Heguy, et al. Dominance of Alpha and Iota variants in SARS-CoV-2 vaccine breakthrough infections in New York City. *The Journal of Clinical Investigation*, 131(18):e152702, 2021.
- [40] Lauren C. Tindale, Jessica E. Stockdale, Michelle Coombe, Emma S. Garlock, Wing Yin Venus Lau, Manu Saraswat, Louxin Zhang, Dongxuan Chen, Jacco Wallinga, and Caroline Colijn. Evidence for transmission of COVID-19 prior to symptom onset. *eLife*, 9:e57149, 2020.
- [41] Hideo Tanaka, Tsuyoshi Ogata, Toshiyuki Shibata, Hitomi Nagai, Yuki Takahashi, Masaru Kinoshita, Keisuke Matsubayashi, Sanae Hattori, and Chie Taniguchi. Shorter incubation period among COVID-19 cases with the BA. 1 Omicron variant. *International Journal of Environmental Research and Public Health*, 19(10):6330, 2022.
- [42] Rebecca Grant, Tiffany Charmet, Laura Schaeffer, Simon Galmiche, Yoann Madec, Cassandre Von Platen, Olivia Chény, Faïza Omar, Christophe David, Alexandra Rogoff, et al. 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, 2022.
- [43] Tsuyoshi Ogata, Hideo Tanaka, Fujiko Irie, Atsushi Hirayama, and Yuki Takahashi. Shorter incubation period among unvaccinated delta variant coronavirus disease 2019 patients in Japan. *International Journal of Environmental Research and Public Health*, 19(3):1127, 2022.
- [44] Public Health Agency of Canada. COVID-19 for health professionals: Transmission. <https://www.canada.ca/en/public-health/services/diseases/2019-novel-coronavirus-infection/health-professionals/transmission.html>, 2021.
- [45] Nazar Zaki and Elfadil A. Mohamed. The estimations of the COVID-19 incubation period: A scoping reviews of the literature. *Journal of Infection and Public Health*, 14(5):638–646, 2021.

- [46] Jordi Cortés Martínez, Daewoo Pak, Gabriela Abelenda-Alonso, Klaus Langohr, Jing Ning, Alexander Rombauts, Mireia Colom, Yu Shen, and Guadalupe Gómez Melis. 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):1–7, 2022.
- [47] Stefan Elbe and Gemma Buckland-Merrett. Data, disease and diplomacy: GISAID’s innovative contribution to global health. *Global Challenges*, 1(1):33–46, 2017.
- [48] Ryan J. Tibshirani. Adaptive piecewise polynomial estimation via trend filtering. *The Annals of Statistics*, 42(1):285–323, 2014.
- [49] Ryan J. Tibshirani. Divided differences, falling factorials, and discrete splines: Another look at trend filtering and related problems. *Foundations and Trends in Machine Learning*, 15(6):694–846, 2022.
- [50] Aaditya Ramdas and Ryan J. Tibshirani. Fast and flexible ADMM algorithms for trend filtering. *Journal of Computational and Graphical Statistics*, 25(3):839–858, 2016.
- [51] Centers for Disease Control and Prevention. COVID Data Tracker. <https://covid.cdc.gov/covid-data-tracker/#national-lab>, 2020.
- [52] Centers for Disease Control and Prevention. 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>, 2021.
- [53] Centers for Disease Control and Prevention. Nationwide commercial laboratory seroprevalence survey. <https://data.cdc.gov/Laboratory-Surveillance/Nationwide-Commercial-Laboratory-Seroprevalence-Su/d2tw-32xv>, 2021.
- [54] Jeanne Ruff, Ying Zhang, Matthew Kappel, Sfurti Rath, Kellie Watkins, Lei Zhang, and Cassius Lockett. Rapid increase in suspected SARS-CoV-2 reinfections, Clark County, Nevada, USA, December 2021. *Emerging Infectious Diseases*, 28(10):1977, 2022.
- [55] James Durbin and Siem Jan Koopman. *Time Series Analysis by State Space Methods*, volume 38. OUP Oxford, 2012.
- [56] Jouni Helske. KFAS: Exponential family state space models in R. *Journal of Statistical Software*, 78(10): 1–39, 2017.
- [57] U.S. Census Bureau, Population Division. 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>, 2022.
- [58] Jefferson M. Jones, Mars Stone, Hasan Sulaeman, Rebecca V. Fink, Honey Dave, Matthew E. Levy, Clara Di Germanio, Valerie Green, Edward Notari, Paula Saa, et al. Estimated US infection-and vaccine-induced SARS-CoV-2 seroprevalence based on blood donations, July 2020-May 2021. *JAMA*, 326(14):1400–1409, 2021.
- [59] Kristina L. Bajema, Ryan E. Wiegand, Kendra Cuffe, Sadhna V. Patel, Ronaldo Iachan, Travis Lim, Adam Lee, Davia Moyse, Fiona P. Havers, Lee Harding, et al. Estimated SARS-CoV-2 seroprevalence in the US as of September 2020. *JAMA Internal Medicine*, 181(4):450–460, 2021.

## 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<sup>47</sup>, on which this research is based.

ATTN: Grants, Delphi, etc.

## **Author contributions**

## **Competing interests**

The authors declare no competing interests.

## Online Supplement

### S1 Additional information about dataset used or estimation methodology

#### S1.1 Table on the percent pairwise occurrence of events in the CDC line list

Order of events	Percent pairwise occurrence	Handling
IO $\rightarrow$ SO $\rightarrow$ PS $\rightarrow$ RE	PS $\geq$ SO: 97.1 PS = SO: 33.6 PS $>$ RE: 1.74 PS = RE: 14.6	This is the idealized order of events and so we built the current support sets for SO $\rightarrow$ PS and PS $\rightarrow$ RE delay distribution constructions around this such that IO comes first by construction, SO typically precedes PS, but may be the same or come before, and RE comes after PS and SO
IO $\rightarrow$ PS $\rightarrow$ SO $\rightarrow$ RE	PS $<$ SO: 2.91 SO $\leq$ RE: 99.3 SO $<$ RE: 86.1	Allowed for negative delays up to the largest non-outlier value for the 0.05 quantile of delay from PS to SO by state
IO $\rightarrow$ PS $\rightarrow$ RE $\rightarrow$ SO	RE $<$ SO: 0.7 RE $<$ PS: 1.7	Nothing because current handling of the CDC of the line list ensures that the most concerning cases are handled where SO = PO = RE, SO = RE and PO = RE

Table S1: Percent pairwise occurrence for the different permutations of events considered in the restricted CDC line list. The abbreviation IO stands for infection onset, SO is symptom onset, PS is positive specimen, and RE is report date. We consider a restricted set of permutations because we assume that IO must come first and that PS must precede report date for a case to be legitimate. Finally, the underlying assumption for the percent pairwise occurrence calculations is that the cases must have both elements present (not missing).

#### S1.2 State space representation of the antibody prevalence model

The antibody prevalence model from Equation 6 is conceptualized as a Gaussian state space model (as in <sup>55;56</sup>).

In general, for  $t = 1, \dots, n$ , let  $\alpha_t$  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  $\eta$  is a  $k \times 1$  vector, the form of the linear Gaussian state space model is

$$y_t = Z\alpha_t + \epsilon_t, \quad \epsilon_t \sim N(0, H_t) \quad (10)$$

$$\alpha_{t+1} = T_t\alpha_t + R_t\eta_t, \quad \eta_t \sim N(0, Q_t) \quad (11)$$

where  $\alpha_1 \sim N(a_1, P_1)$  and there is independence amongst  $\alpha_1$ ,  $\epsilon_t$  and  $\eta_t$  <sup>55;56</sup>. 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  $\alpha_t$ , while the second equation is a first-order autoregressive model, which is Markovian in nature <sup>55</sup>.

The underlying idea behind the two equations is that we are assuming that the system evolves according to  $\alpha_t$  (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  $\alpha_t$  (as in the first equation) to drive the system forward <sup>55</sup>. So the objective of state space modeling is to obtain the latent states  $\alpha$  based on the observations  $y$  and this is achieved through Kalman filtering and smoothing.

Kalman filtering gives the following one-step-ahead predictions of the states

$$a_{t+1} = \mathbb{E}[\alpha_{t+1} \mid y_t, \dots, y_1]$$

with covariance,

$$P_{t+1} = \text{Var}(\alpha_{t+1} \mid y_t, \dots, y_1).$$

Then, the Kalman smoother works backwards to the first time to give

$$\hat{\alpha}_t = \mathbb{E}[\alpha_t \mid y_n, \dots, y_1] \quad (12)$$

$$V_t = \text{Var}(\alpha_t \mid y_n, \dots, y_1). \quad (13)$$

The filtering and smoothing steps are based on recursions that are described in Appendix A of Helske<sup>56</sup> as we use the R package KFAS to estimate our model.

To express the antibody prevalence model in state space form, we define the components in Equations 10 and 11 as follows:

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

where  $\sigma_o^2$  is the variance of observations,  $\sigma_s^2$  is the variance of the seroprevalence estimates, and  $\sigma_a^2$  is the trend variance. Since we expect the inverse ratios to be more variable than the seroprevalence estimates, we enforce that the estimate of  $\sigma_a^2$  is a multiple of  $\sigma_s^2$ . 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  $\alpha_1$  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  $\tilde{s}_1$ . 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.<sup>21</sup> 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  $\sigma_o^2$  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  $\sigma_s^2$  and  $\gamma$  values are fixed and from averaging the estimated values for all states on the real line (obtained under the starting conditions  $\sigma_s^2 = 3 \times 10^{-6}$ ,  $\gamma = 0.99$ , and  $\sigma_o^2$  as described).



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 12 and 13. Forwards and backwards extrapolation is then used to estimate the ratios and covariance outside of the observed seroprevalence range<sup>55</sup>, 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 antibody prevalence model (refer to S1.3 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.

### S1.3 Bayesian specification of the antibody prevalence model

In brief, the antibody prevalence 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  $\gamma$ ),  $D$  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

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

### S1.4 Ablation analysis of infection-hospitalization correlations

We undertake an ablation study for the lagged correlation of infections, the results of which are shown in Figure 5. From this, we can see that the deconvolved case or infection estimates from the intermediate steps 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 positive specimen date tend to precede hospitalizations by about 11 days, while those for the subsequent step indicate that the deconvolved cases by symptom onset tend to precede hospitalizations by a longer time of about 13 days. Finally, after adding the variant-specific incubation period data into the deconvolution and obtaining the deconvolved case estimates, we can observe that the reported infections precede hospitalizations by about 17 days.

### S1.5 Scaling by population

Annual estimates of the resident state populations as of July 1 of 2020 and 2021 are taken from the December 2022 press release on the U.S. Census Bureau website<sup>57</sup>. Unless otherwise specified, we use the July 1, 2020 estimates.

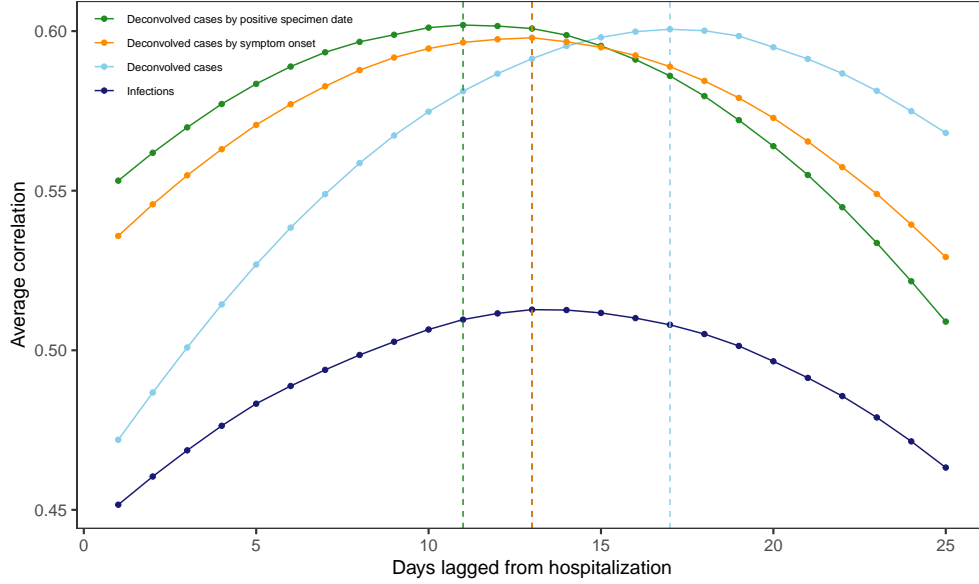


Figure S1: Lagged Spearman’s correlation between the infection and hospitalization rates per 100,000 averaged for each lag across U.S. states and days over June 1, 2020 to November 29, 2021, and taken over a rolling window of 61 days. The infection rates are based on the counts for the deconvolved case 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.

## S1.6 Additional details on the date fields in the CDC linelist

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 further significant revisions.

Table S1 presents the percent of pairwise occurrences for the different possible permutations of events in the line list. Essentially, most cases follow the idealized ordering shown by Figure 1 and so we adhere to this construction as much as possible.

We observe that the line list is prone to high percentages of missing data, notably with respect to our variables of interest. Approximately 62.3% of cases are missing the symptom onset date, 55.4% are missing positive specimen date, and 8.96% of cases are missing the report date. Relatedly, cases with missing report or positive specimen dates may be filled with their symptom onset date Jahja et al.<sup>24</sup>. So it is possible that all three variables may be imputed with the same date for a case. However, we only actually deal with select pairs of events; we do not use all three at once in our construction of the delay distributions or anywhere else in our analysis. Therefore, we restrict our investigation of missingness to the pairs of events. Figure S2 suggests that this issue impacts states differentially due to the inconsistent proportions of zero delay between positive specimen and report date across states.

Due to the contamination in the zero delay cases (the true extent of which is unknown to us), we omit all such cases where the positive specimen and report dates have zero delay from our analysis. We choose to allow for zero and negative delay for symptom onset to report because correspondence with the CDC confirms the distinct possibility that a person could test positive before symptom onset and it is a reasonable ordering to expect if, for example, the individual is aware that they have been exposed to an infected individual.

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 S2, from which it is clear the fraction of missing cases is substantial for many states, often surpassing 50%<sup>24</sup>. 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.

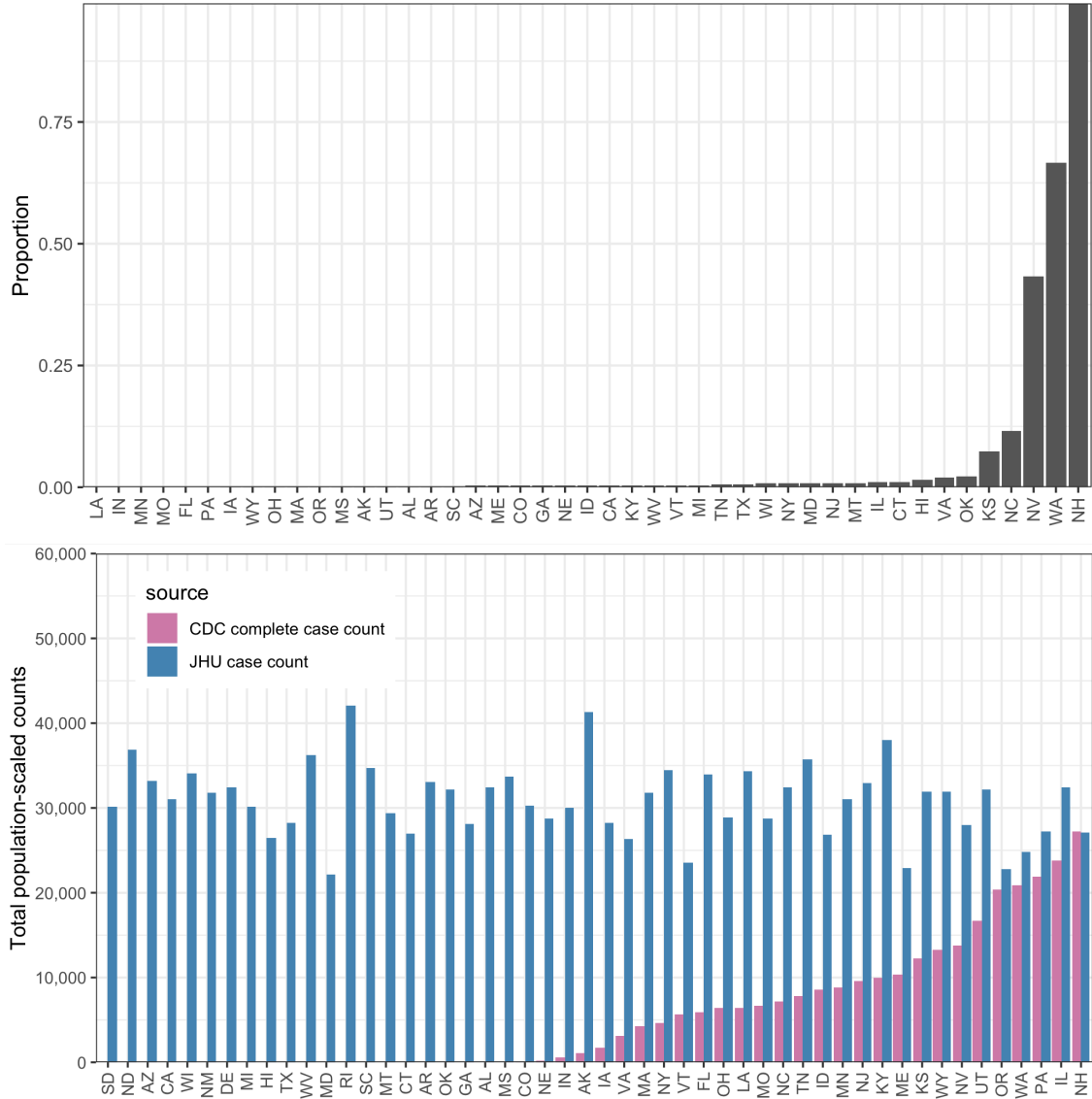


Figure S2: Top panel: Proportion of complete cases with zero delay between positive specimen and report date in the restricted CDC line list dataset. Bottom panel: Complete case counts by state in the CDC line list versus the cumulative complete case counts from JHU CSSE as of June 6, 2022. All counts have been scaled by the 2022 state populations as of July 1, 2022 from U.S. Census Bureau, Population Division <sup>57</sup>.

We consider such bias to be unavoidable in our analysis due to a lack of alternative line 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 positive specimen to case report. Therefore, we devise a simple, ad hoc approach to detect and prune these reporting backlogs.

First, we obtain the part of the line list intended for the positive specimen to case report delay estimation, where both such dates are present and where zero and negative delay cases have been omitted. Then, for each of the three dates of June 1, September 1, and December 1, 2020, we bin the reporting delays occurring from

50 days up to the maximum observed delay. 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).

## S1.7 Justifications for delay distribution calculations

Let  $y_t$  denote the count of new cases reported at time  $t$  and  $x_t$  denote the count of deconvolved cases with positive specimen at  $t$ . For all cases in the line list that had both a positive specimen and a report date, we can count the those that are reported at time  $t$  by enumerating them according to positive specimen date (similar to how symptom onset date was used in<sup>24</sup>):

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

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  $\pi_t(s) = \mathbb{P}(\text{case report at } t \mid \text{positive specimen date at } s)$  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 positive specimen date, except 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 these distributions. First, positive specimen tests that are reported to the CDC are always reported within  $d = 60$  days, which is true for the majority of the 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.<sup>24</sup>, 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  $k = 1, \dots, 60$ ,

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

Thirdly, there are times where the empirical probability was observed to be precisely 1 at zero delay and the proportion of CDC relative to JHU cases used for the weight was also 1. Since we believe that having zero delay for all cases is unrealistic and unlikely to be representative of all cases for the state, we inject a small amount of variance manually by setting the the CDC-to-JHU proportion to be the minimum shrinkage proportion observed for the affected state (such instances were isolated to the state of New Hampshire). Aside from these modifications, the construction of the delay distribution proceeds in precisely the same manner as for positive specimen to report date.

## S1.8 Details about seroprevalence data

In the former, the CDC collaborated with 17 blood collection organizations in the largest nationwide COVID-19 seroprevalence survey to date<sup>52</sup>. The blood donation samples were used to construct monthly seroprevalence estimates for nearly all states from July 2020 to December 2021<sup>58</sup>. 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,<sup>59</sup>). 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<sup>51;58;59</sup>. 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 S3). 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 S3, the final commercial seroprevalence measurement from 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. Note that North Dakota is the only state to which this exclusion does not apply as there are no commercial seroprevalence measurements beyond 2020.

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.

We convert our daily 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 seroprevalence source, then the average is used. We denote these changes by changing the time-based subscript from  $t$  to  $m$  where  $m$  indicates the Monday relative to our June 1, 2020 start date.

## S1.9 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 infection onset is the same as the positive specimen date (i.e., excluding the positive specimen to infection onset data and deconvolution);
2. those from the deconvolution procedure under the assumption that the infection onset is the same as the symptom onset date (excluding the incubation period data);
3. those from the deconvolution procedure when utilizing all incubation period and delay data (the deconvolved case estimates);
4. those from applying the antibody prevalence model to produce estimates for both the reported and the unreported cases (the infection estimates).

## S1.10 Possible investigation of reinfections

Possible change to the paper based on Ajitesh’s feedback - Main contribution is the model, shows without reinfection & here’s an extension that shows how to include reinfection data.

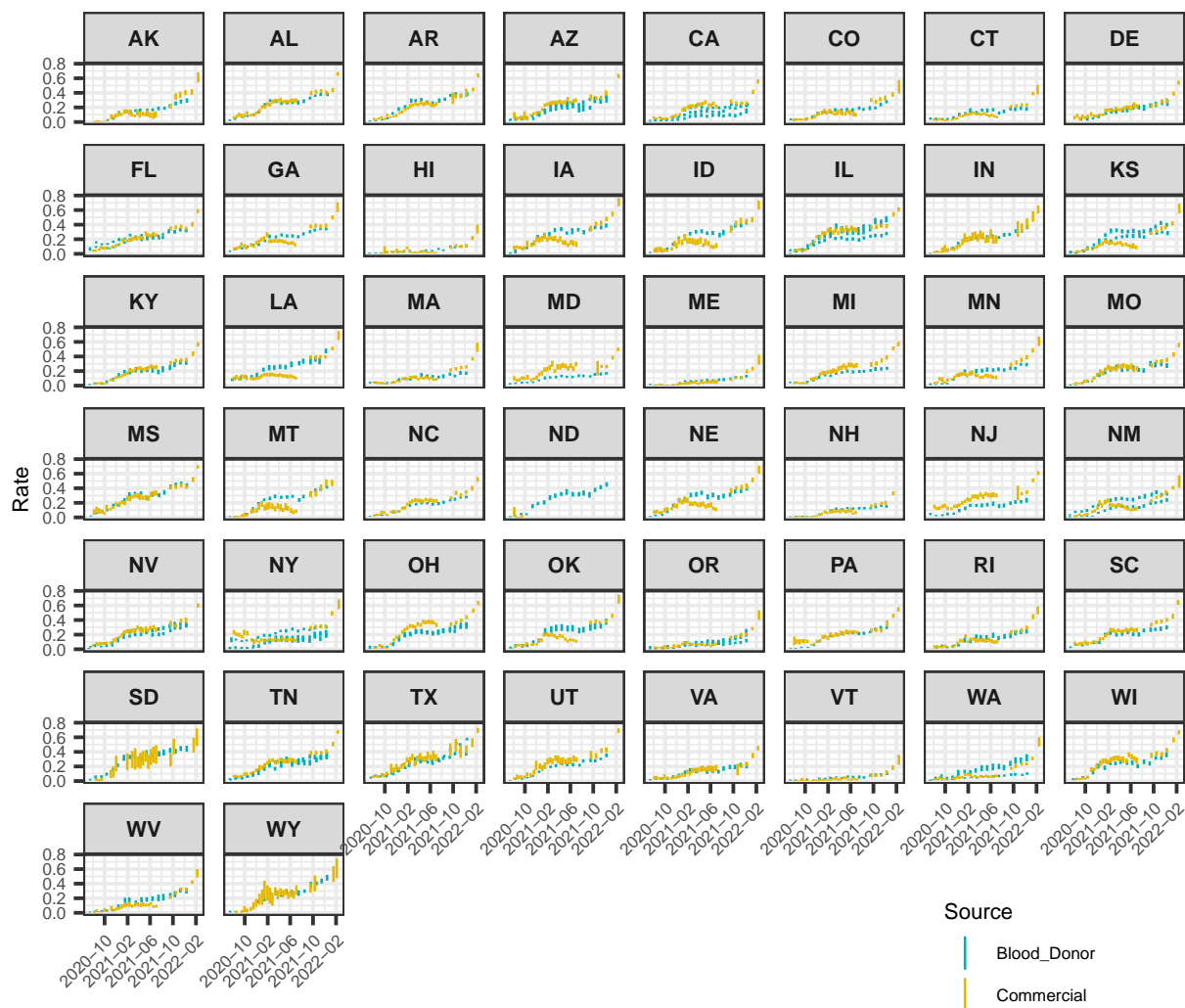


Figure S3: A comparison of the seroprevalence estimates from the Commercial Lab Seroprevalence Survey dataset (yellow) and the 2020–2021 Blood Donor Seroprevalence Survey dataset (blue). Note that the maximum and the minimum of the line ranges are the provided 95% confidence interval bounds to give a rough indication of uncertainty. **ATTN:** This figure doesn't use space very well. Let's remove the gap between panels and make the facet labels (the state names) normal sized. The x-axis could just show 2021, 2022 (rather than so many ticks). Alternatively, another map layout? And I don't think "Rate" is correct for the y-axis.