

Accurately Estimating Total COVID-19 Infections using Information Theory

Jiaming Cui¹, Arash Haddadan², A S M Ahsan-Ul Haque³, Jilles Vreeken⁴, Bijaya Adhikari⁵, Anil Vullikanti^{2,3}, and B. Aditya Prakash^{1,*}

¹College of Computing, Georgia Institute of Technology, Atlanta, GA 30332, US

²Biocomplexity Institute, University of Virginia, Charlottesville, VA 22904, US

³Department of Computer Science, University of Virginia, Charlottesville, VA 22904, US

⁴CISPA Helmholtz Center for Information Security, Saarbrücken 66123, Germany

⁵Department of Computer Science, The University of Iowa, Iowa City, IA 52242, US

Abstract

One of the most significant challenges in the early combat against COVID-19 was the difficulty in estimating the true magnitude of infections. Unreported infections drove up disease spread in numerous regions, made it very hard to accurately estimate the infectivity of the pathogen, therewith hampering our ability to react effectively. Despite the use of surveillance-based methods such as serological studies, identifying the true magnitude is still challenging today. This paper proposes an information theoretic approach for accurately estimating the number of total infections. Our approach is built on top of Ordinary Differential Equations (ODE) based models, which are commonly used in epidemiology and for estimating such infections. We show how we can help such models to better compute the number of total infections and identify the parameterization by which we need the fewest bits to describe the observed dynamics of reported infections. Our experiments show that our approach leads to not only substantially better estimates of the number of total infections but also better forecasts of infections than standard model calibration based methods. We additionally show how our learned parameterization helps in modeling more accurate what-if scenarios with non-pharmaceutical interventions. Our results support earlier findings that most COVID-19 infections were unreported and non-pharmaceutical interventions indeed helped to mitigate the spread of the outbreak. Our approach provides a general method for improving epidemic modeling which is applicable broadly.

*To whom correspondence should be addressed. E-mail: badityap@cc.gatech.edu

Introduction

The COVID-19 pandemic has emerged as one of the most formidable public health challenges in recent history. By Nov 1, 2022, there were already more than 98 million reported infections and 1.07 million deaths in the United States alone. Worldwide, the reported infections summed to 636 million with at least 6.61 million deaths [19]. The devastating effects of COVID-19 extends to the economy as well. For example, in the US, the unemployment rate peaked at 15.8 percent in April 2020 [6], and US GDP contracted at a 3.5% annualized rate for 2020 [1]. Similar economic impacts have been observed worldwide.

One of the most significant challenges in the early combat against COVID-19 was estimating the number of total infections. A significant number of COVID-19 infections were unreported, due to various factors such as the lack of testing and asymptomatic infections [13, 11, 57, 55, 39]. The inability in estimating these unreported infections allowed them to drive up disease transmission in many regions. For example, phylogenetic studies revealed that COVID-19 had locally spread in Washington state before early 2020, when active community surveillance was implemented [14]. There were only 23 reported infections in five major U.S. cities by March 1, 2020, but it has been estimated that there were in fact more than 28,000 total infections by then [5]. Similar trends were observed in other countries, such as in Italy, Germany, and the UK [60]. Despite having more advanced surveillance techniques such as serological studies, estimating the total number of infections continues to be a challenge for COVID-19 response even today [8, 30].

An accurate estimation of the number of total infections is a fundamental epidemiological question and critical for pandemic planning and response. Notwithstanding its importance, there is not even a commonly agreed upon metric. One proposal is the *case ascertainment rate*, which is defined as the ratio of reported symptomatic infections to the actual number of symptomatic infections [52]. Another popular proposal is the *reported rate* α_{reported} , which is defined as the ratio of reported infections to total infections [46]. This definition includes asymptomatic infections, which are known to contribute substantially to community transmission [58, 41]. In this paper, we focus on this particular measure.

However, estimating the reported rate is challenging, and as a result all current methods have their limitations. One of the most effective current methods to identify the reported rate in a region is through large-scale serological studies [56, 26, 64]. These surveys use blood tests to identify the prevalence of antibodies against SARS-CoV-2 in a large population. The CDC COVID Data Tracker portal [2, 26] summarizes the results of serological studies conducted by commercial laboratories at a national level as well as at 10 specific sites. For example, the estimated reported rate was at most 0.1 in Minneapolis and South Florida as of April 2020. This means that there were at least 10 times more total infections than reported infections. While serological studies can give an accurate estimation, they are expensive and are not sustainable in the long run [4]. Furthermore, it is also challenging to obtain real-time data using such studies since there are unavoidable delays between sample collection and laboratory tests [2, 26]. Additional difficulties include sampling biases that make it necessary to use carefully designed heuristics to account for them [9]. Other methods include exploiting existing surveillance systems of related diseases like influenza, and using them to estimate symptomatic infections [40]. However, this can also be unreliable and requires ad-hoc corrections to account for the similarities between COVID-19 and influenza symptoms.

In the face of these challenges, data scientists and epidemiologists have devoted much time and effort to estimate the reported rate α_{reported} through epidemiological models O_M . By now, there exist carefully constructed models that capture the transmission dynamics of COVID-19 well [39, 52, 12, 50, 36, 38, 61, 25, 33, 62, 63, 17, 43]. In general, an epidemiological model O_M has a set of parameters Θ that we estimate from *observed data* using a so-called calibration procedure,

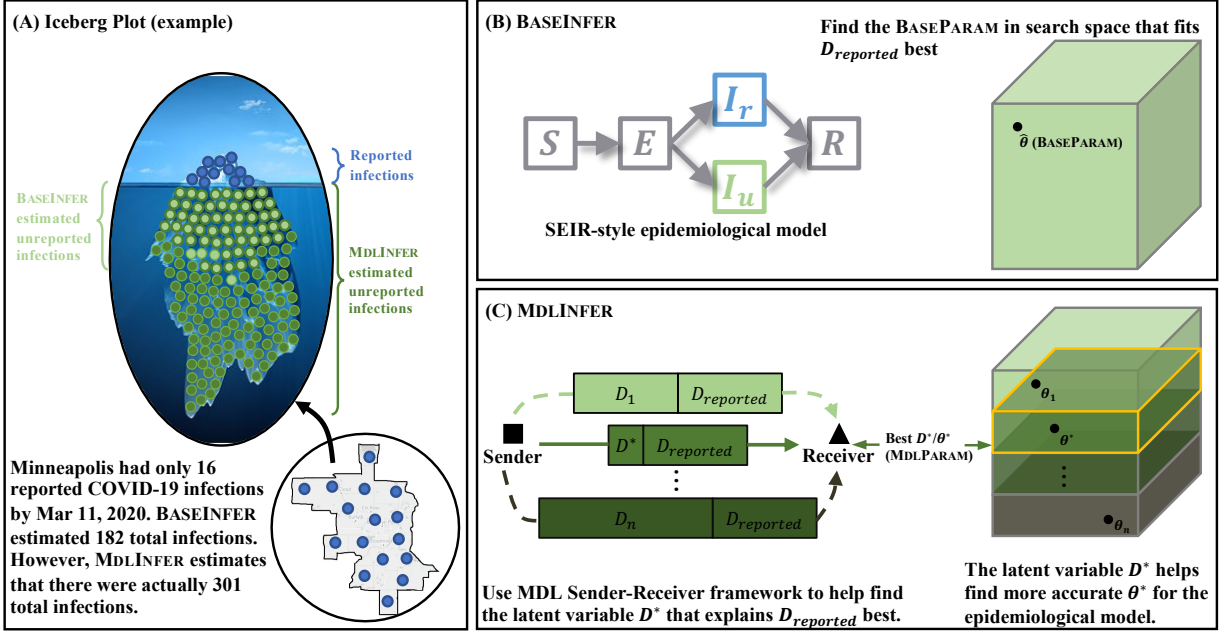


Figure 1: **Overview of our problem and methodology.** (A) We visualize the idea of reported rates using the iceberg. The visible portion above water are the reported infections, which is only a fraction of the whole iceberg representing total infections. Light green corresponds to the unreported infections estimated by typical current practice used by researchers (182 in this example). We call it as the basic approach, or BASEINFER. In contrast, dark green corresponds to the more accurate and much larger 301 unreported infections found by our approach MDLINFER. (B) The usual practice is to calibrate an epidemiological model to reported data and compute the reported rate from the resultant parameterization of the model. Here, an SEIR-style model with explicit compartments for reported-vs-unreported infection is shown in the figure as an example. (C) Our new approach MDLINFER instead aims to compute a more accurate reported rate by finding a ‘best’ parameterization *for the same epidemiological model* (i.e., SEIR-style model in this example) using a principled information theoretic formulation - two-part ‘sender-receiver’ framework. Assume that a hypothetical Sender S wants to transmit the reported infections as the DATA to a Receiver R in the cheapest way possible. Hence S will find/solve for the best D^* , intuitively, the MODEL that takes the fewest number of bits to encode the DATA. Using D^* , we can find the best θ^* by exploring a smaller search space.

CALIBRATE. In practice, the data we use for calibration can be the time series of the number of reported infections, which we call $D_{reported}$. To estimate the number of total infections, these models often explicitly include reported rate $\alpha_{reported}$ as one of their parameters, or include multiple parameters that jointly account for it. There are many calibration procedures commonly used in literature, such as RMSE-based [23] or Bayesian approaches [33, 25].

We call the above general methodology the basic approach to estimate the reported rate, or BASEINFER for short. It takes the epidemiological model O_M , a calibration procedure CALIBRATE, and observed data $D_{reported}$ as input. The output of BASEINFER is then a baseline parameterization $\hat{\Theta}$ and, by extension, an estimated reported rate $\hat{\alpha}_{reported}$. Calibrating a parameterization is generally a complex, high-dimensional problem, since $\hat{\Theta}$ consists of multiple interacting parameters. To make matters worse, there exist many possible parameterizations that show similar performance

(e.g. in RMSE, likelihood) yet correspond to vastly different reported rates. BASEINFER cannot select between these competing parameterizations in a principled way: the parameterization $\hat{\Theta}$ it results in may or may not overfit the reported infections and may or may not predict future infections well. One method for selecting them is to take a Bayesian approach. That is, we choose a prior distribution, and then select the best parameterization that maximizes the posterior probability. Choosing such a prior, however, is ad-hoc and does not generalize well across different models O_M . As we will see in the experimental evaluation, minor differences in estimates of reported rates can indeed lead to very different forecasts of future trends and therewith intervention policy recommendations.

Instead, we propose a new information theory-based approach named MDLINFER. It takes the same input as BASEINFER, but uses a principled approach to determine the best parameterization Θ^* . It is based on the following central intuition: Suppose an oracle also gives us the time series of the number of *total infections* D in addition to the already known reported number of infections D_{reported} , and we are asked to describe D_{reported} as succinctly as possible. As we know both D and D_{reported} , it is trivial to estimate $\alpha'_{\text{reported}}$. If we know D and $\alpha'_{\text{reported}}$, it is trivial to describe D_{reported} , as it is simply $D \times \alpha'_{\text{reported}}$ plus a little bit of noise. Now to most succinctly describe D , we have to calibrate O_M to obtain Θ' . The only things we now have to describe are Θ' , $\alpha'_{\text{reported}}$, the (small) errors that O_M makes in predicting D , and the (small) errors that we make predicting D_{reported} using D and $\alpha'_{\text{reported}}$. In practice, we are of course not given D , but the key idea of this paper is to estimate D as a latent variable such that we can most succinctly describe (most accurately reconstruct) the dynamics of D_{reported} .

In practice, we need both a way of measuring how well a latent MODEL (i.e., D and its corresponding $\alpha'_{\text{reported}}$) describes the DATA (i.e., reported infections D_{reported}), as well as a way to find the best such MODEL. To do so, the Minimum Description Length (MDL) principle provides a statistically sound approach. MDL has been widely used for numerous optimization problems ranging from network summarization [34], causality inference [16], and failure detection in critical infrastructures [10]. MDL has also previously been used for some epidemiological problems, mainly in inferring patient-zero and associated infections in cascades over contact networks [49]. However, we are the first to propose an MDL-based approach on top of ODE-based epidemiological models, which are harder to formulate and optimize.

Specifically, we use two-part MDL (aka sender-receiver framework) consisting of hypothetical actors S and R : Sender S has the DATA and wants to transmit it to receiver R using as few bits as possible [24]. Hence, sender S searches for the best possible MODEL, which minimizes the overall cost of encoding and transmitting both the MODEL and the DATA given the MODEL. Following the convention in information theory, we use $L(\text{MODEL})$ to denote the number of bits required to encode the MODEL; and $L(\text{DATA}|\text{MODEL})$ to denote the number of bits required to encode the DATA, D_{reported} , given the MODEL. The overall objective of our optimization problem is to infer an optimal MODEL*, which minimizes $L(\text{MODEL}) + L(\text{DATA}|\text{MODEL})$. To put MDL to practice for our problem, we carefully design our MDL cost to minimize the discrepancy in fitting D_{reported} . This cost ensures the generalizability of our learned D^* and $\alpha^*_{\text{reported}}$ - it can avoid overfitting on D_{reported} and predict the future reported infections well. Our later experiments exactly show this. Our approach, MDLINFER, can be applied to any ODE model since two-part MDL does not assume about the nature of the DATA or the MODEL.

We compare MDLINFER and BASEINFER using two different ODE-based epidemiological models: SAPHIRE [25] and SEIR + HD [33] as O_M . Following their literature [25, 33], we use Markov Chain Monte Carlo (MCMC) as the calibration procedure CALIBRATE for SAPHIRE and iterated filtering (IF) for SEIR + HD, both of which are Bayesian approaches [29]. Both these epidemiological models have previously been shown to perform well in fitting reported infections and provided insight

that was beneficial for the COVID-19 response. SAPHIRE focuses on two key features of the outbreak: high covertness and high transmissibility that drove the outbreak of COVID-19 in Wuhan. SEIR + HD investigates how non-pharmaceutical interventions like social distancing will be needed to maintain epidemic control. These models are broadly representative to show that MDLINFER gives consistent performance across multiple epidemiological models with different dynamics. The experiments clearly show that our proposed MDL-based approach MDLINFER performs superior to the state of the art. To illustrate, we give an example in Fig. 1. By March 11, 2020, the Minneapolis Metro Area had only 16 COVID-19 reported infections. BASEINFER estimated 182 total infections, which are colored as light green in the iceberg. On the other hand, our MDLINFER gives an estimate of 301 total infections shown below the sea level, which is closer to the total infections estimated from serological studies [26, 2]. Additionally, MDLINFER also leads to better fits and *future* projections on reported infections. We also demonstrate that MDLINFER can aid policy making by analyzing counter-factual non-pharmaceutical interventions, while inaccurate BASEINFER estimates lead to wrong non-pharmaceutical intervention conclusions.

Results

Next, we present our empirical findings on a large set of experiments in different geographical regions and time periods. We choose 8 regions and periods based on the severity of the outbreak and the availability of serological studies and symptomatic surveillance data. In each region, we divide the timeline into two time periods: (i) observed period, when only the number of reported infections are available, and both BASEINFER and MDLINFER are used to learn the baseline parameterization (BASEPARAM) $\hat{\Theta}$ and MDL parameterization (MDLPARAM) Θ^* , and (ii) forecast period, where we evaluate the forecasts generated by the parameterizations learned in the observed period. To handle the time-varying reported rates, we divide the observed period into multiple sub-periods and learn different reported rates for each sub-period separately.

(A) Estimating total infections: MDLINFER estimates total infections more accurately than BASEINFER

Here, we use the point estimates of the total infections calculated from serological studies as the ground truth (black dots shown in Fig. 2). We call it SEROSTUDY_{Tinf}. We also plot MDLINFER’s estimation of total infections, MDLPARAM_{Tinf}, in the same figure (red curve). To compare the performance of MDLINFER and BASEINFER with SEROSTUDY_{Tinf}, we use the cumulative value of estimated total infections. Note that values from the serological studies are not directly comparable with the total infections because of the lag between antibodies becoming detectable and infections being reported [2, 26]. In Fig. 2, we have already accounted for this lag following CDC study guidelines [2, 26] (See Methods section for details). The vertical black lines shows a 95% confidence interval for SEROSTUDY_{Tinf}. The blue curve represents total infections estimated by BASEINFER, BASEPARAM_{Tinf}. As seen in the figure, MDLPARAM_{Tinf} falls within the confidence interval of the estimates given by serological studies. Significantly, in Fig. 2B and Fig. 2F for South Florida, BASEINFER for SAPHIRE model [25] overestimates the total infections, while for SEIR + HD model underestimates the total infections. However, MDLINFER consistently estimates the total infections correctly. This observation shows that as needed, MDLPARAM_{Tinf} can improve upon the BASEPARAM_{Tinf} in either direction (i.e., by increasing or decreasing the total infections). Note that the MDLPARAM_{Tinf} curves from both models are closer to the SEROSTUDY_{Tinf} even when the BASEPARAM_{Tinf} curves are different. The results of better accuracy in spite of various

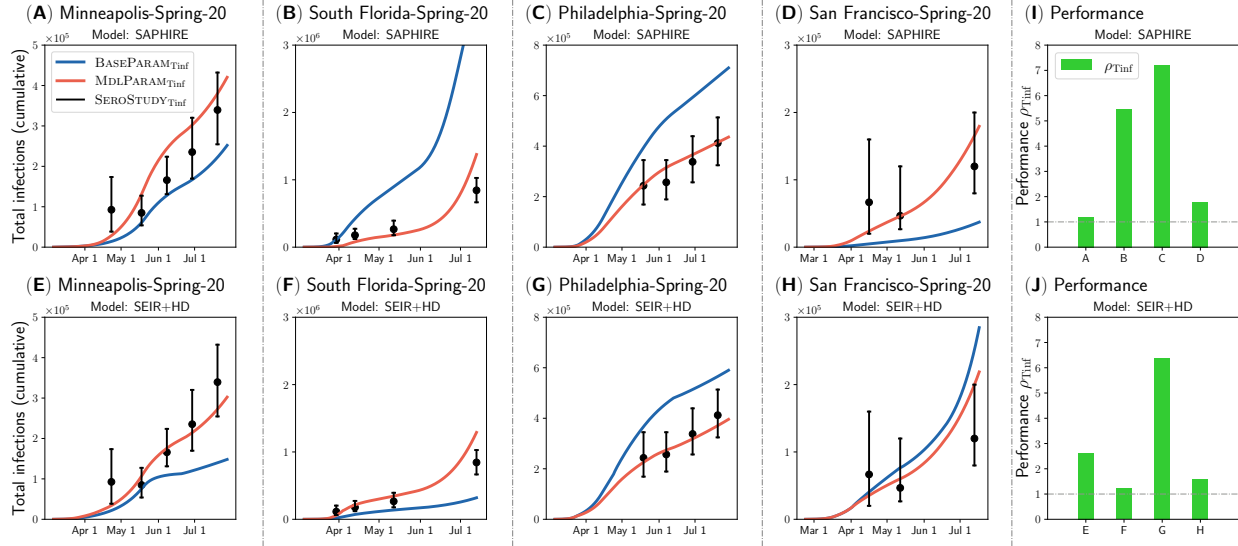


Figure 2: **MDLINFER (red) gives a closer estimation of total infections to serological studies (black) than BASEINFER (blue) on various geographical regions and time periods.** Note that both approaches try to fit the serological studies without being informed with them. (A)-(H) The red and blue curves represent MDLINFER’s estimation of total infections, MDLPARAM_{Tinf}, and BASEINFER’s estimation of total infections, BASEPARAM_{Tinf}, respectively. The black point estimates and confidence intervals represent the total infections estimated by serological studies [2, 26], SEROSTUDY_{Tinf}. (A)-(D) use SAPHIRE model and (E)-(H) use SEIR + HD model. (I)-(J) The performance metric, ρ_{Tinf} , comparing MDLPARAM_{Tinf} against BASEPARAM_{Tinf} in fitting serological studies is shown for each region. (I) is for SAPHIRE model in (A)-(D), and (J) is for SEIR + HD model in (E)-(H). Here, the values of ρ_{Tinf} are 1.20, 5.47, 7.21, and 1.79 in (I), and 2.62, 1.22, 6.39, and 1.58 in (J). Note that ρ_{Tinf} larger than 1 means that MDLPARAM_{Tinf} is closer to SEROSTUDY_{Tinf} than BASEPARAM_{Tinf}. We show more experiments in the Supplementary Information.

geographical regions and time periods show that MDLINFER is consistently able to estimate total infections more accurately.

To quantify the performance gap between the two approaches, we first compute the root mean squared error (RMSE) between SEROSTUDY_{Tinf} and BASEPARAM_{Tinf}. We also compute the same between SEROSTUDY_{Tinf} and MDLPARAM_{Tinf}. We then compute the ratio, ρ_{Tinf} , of the two RMSE errors as $\frac{RMSE(BASEPARAM_{Tinf}, SEROSTUDY_{Tinf})}{RMSE(MDLPARAM_{Tinf}, SEROSTUDY_{Tinf})}$. Note that the values of ρ_{Tinf} being greater than 1 implies that the MDLPARAM_{Tinf} is closer to SEROSTUDY_{Tinf} estimates than BASEPARAM_{Tinf}. In Fig. 2I and Fig. 2J, we plot ρ_{Tinf} . Overall, the ρ_{Tinf} values are greater than 1 in Fig. 2I and Fig. 2J, which indicates that MDLINFER performs better than BASEINFER. Note that even when the value of ρ_{Tinf} is 1.20 for Fig. 2A, the improvement made by MDLPARAM_{Tinf} over BASEPARAM_{Tinf} in terms of RMSE is about 12091. Hence, one can conclude that MDLINFER is indeed superior to BASEINFER, when it comes to estimating total infections. We show more experiments in the Supplementary Information.

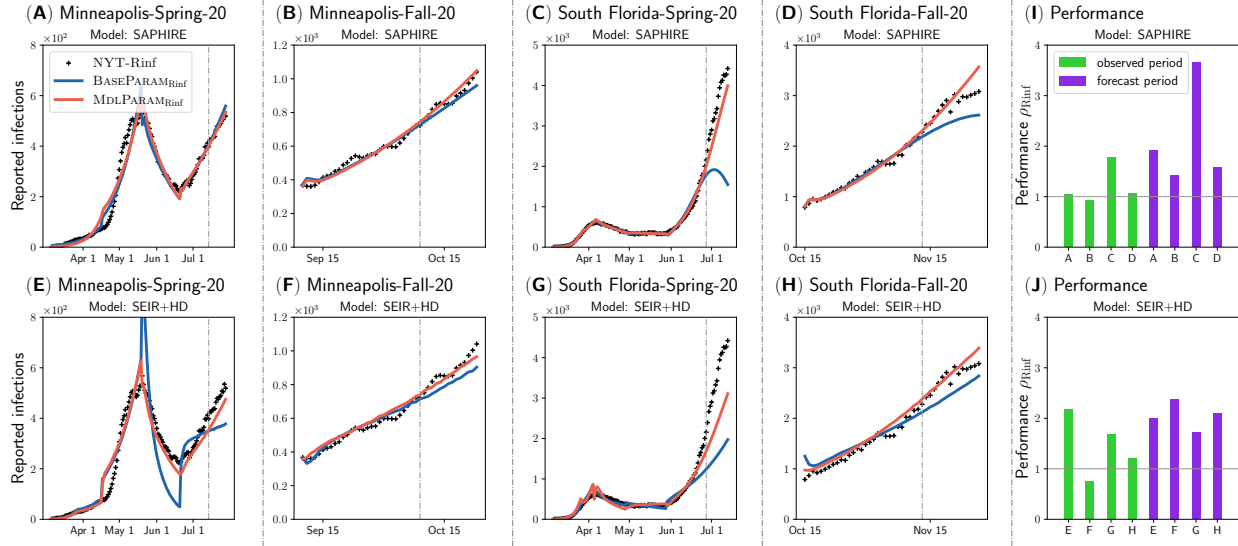


Figure 3: MDLINFER (red) gives a closer estimation of reported infections (black) than BASEINFER (blue) on various geographical regions and time periods. We use the reported infections in the observed period as inputs and try to forecast the future reported infections (forecast period). (A)-(H) The vertical grey dash line divides the observed period (left) and forecast period (right). The red and blue curves represent MDLINFER’s estimation of reported infections, MDLPARAM_{Rinf}, and BASEINFER’s estimation of reported infections, BASEPARAM_{Rinf}, respectively. The black plus symbols represent the reported infections collected by the New York Times (NYT-Rinf). (A)-(D) use SAPHIRE model and (E)-(H) use SEIR + HD model. (I)-(J) The performance metric, ρ_{Rinf} , comparing MDLPARAM_{Rinf} against BASEPARAM_{Rinf} in fitting reported infections is shown for each region. (I) is for SAPHIRE model in (A)-(D), and (J) is for SEIR + HD model in (E)-(H). Note that ρ_{Rinf} larger than 1 means that MDLPARAM_{Rinf} is closer to NYT-Rinf than BASEPARAM_{Rinf}. We show more experiments in the Supplementary Information.

(B) Estimating reported infections: MDLINFER leads to better fit and projection than BASEINFER at different stages of the COVID-19 epidemic

Here, we first use the observed period to learn the parameterizations. We then *forecast* the future reported infections (i.e., forecast periods), which were *not* accessible to the model while training. The results are summarized in Fig. 3. In Fig. 3A to Fig. 3H, the vertical grey dash line divides the observed and forecast period. The black plus symbols represent reported infections collected by the New York Times, NYT-Rinf. The red curve represents MDLINFER’s estimation of reported infections, MDLPARAM_{Rinf}. Similarly, the blue curve represents BASEINFER’s estimation of reported infections, BASEPARAM_{Rinf}. Note that the curves to the right of the vertical grey line are future predictions. As seen in Fig. 3, MDLPARAM_{Rinf} aligns more closely with NYT-Rinf than BASEPARAM_{Rinf}, indicating the superiority of MDLINFER in fitting and forecasting reported infections.

We define a performance metric ρ_{Rinf} as $\frac{RMSE(BASEPARAM_{Rinf}, NYT-Rinf)}{RMSE(MDLPARAM_{Rinf}, NYT-Rinf)}$ to compare MDLPARAM_{Rinf} against BASEPARAM_{Rinf} in a manner similar to ρ_{Tinf} . In Fig. 3I and Fig. 3J, we plot the ρ_{Rinf} for the observed and forecast period. In both periods, we notice that the ρ_{Rinf} is close to or greater than 1. This further shows that MDLINFER has a better or at least closer fit for reported infections than BASEINFER. Additionally, the ρ_{Rinf} for the forecast period is even greater than ρ_{Rinf}

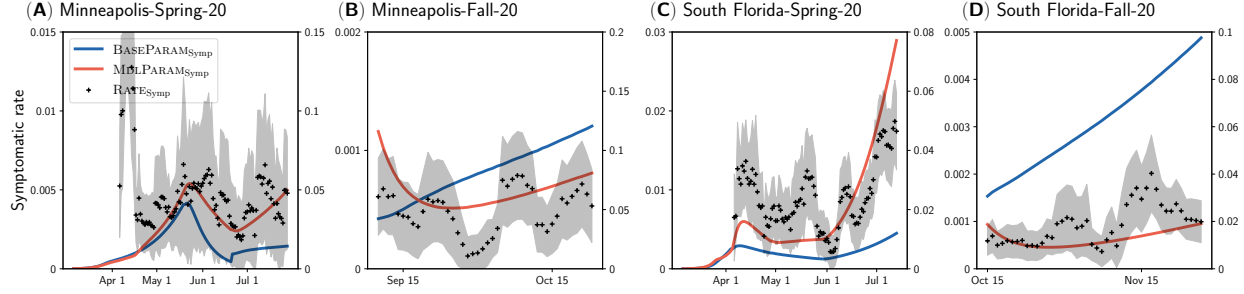


Figure 4: MDLINFER (red) gives a closer estimation of the trends of symptomatic rate (black) than BASEINFER (blue) on various geographical regions and time periods. (A)-(D) The red and blue curves represent MDLINFER’s estimation of symptomatic rate, MDLPARAM_{Symp}, and BASEINFER’s estimation of symptomatic rate, BASEPARAM_{Symp}, respectively. They use the y-scale on the left. The black points and the shaded regions are the point estimate with standard error for RATE_{Symp} (the COVID-related symptomatic rates derived from the symptomatic surveillance dataset [51, 53]). They use the y-scale on the right. Note that we focus on trends instead of the exact numbers, hence MDLPARAM_{Symp}/BASEPARAM_{Symp}, and RATE_{Symp} may scale differently. We show more experiments in the Supplementary Information.

for the observed period, which shows that MDLINFER performs even better than BASEINFER while forecasting.

Note that Fig. 3A, C, E, G correspond to the early state of the COVID-19 epidemic in spring and summer 2020, and Fig. 3B, D, F, H correspond to fall 2020. We can see that MDLINFER performs well in estimating temporal patterns at different stages of the COVID-19 epidemic. We show more experiments in the Supplementary Information.

(C) Estimating symptomatic rate trends: MDLINFER estimates the symptomatic rate trends more accurately than BASEINFER

We validate this observation using Facebook’s symptomatic surveillance dataset [51]. We plot MDLINFER’s and BASEINFER’s estimated symptomatic rate over time and overlay the estimates and standard error from the symptomatic surveillance data in Fig. 4. The red and blue curves are MDLINFER’s and BASEINFER’s estimation of symptomatic rates, MDLPARAM_{Symp} and BASEPARAM_{Symp} respectively. Note that SAPHIRE model does not contain states corresponding to the symptomatic infections. Therefore, we only focus on SEIR + HD model. We compare the trends of the MDLPARAM_{Symp} and BASEPARAM_{Symp} with the symptomatic surveillance results. We focus on trends rather than actual values because the symptomatic rate numbers could be biased [51] (see Methods section for a detailed discussion) and therefore cannot be compared directly with model outputs like what we have done for serological studies. As seen in Fig. 4, MDLPARAM_{Symp} captures the trends of the surveyed symptomatic rate RATE_{Symp} (black plus symbols) better than BASEPARAM_{Symp}. We show more experiments in the Supplementary Information.

To summarize, these three sets of experiments in (A), (B) and (C) together demonstrate that BASEINFER fail to accurately estimate the total infections including unreported ones. On the other hand, MDLINFER estimates total infections closer to those estimated by serological studies and better fits reported infections and symptomatic rate trends.

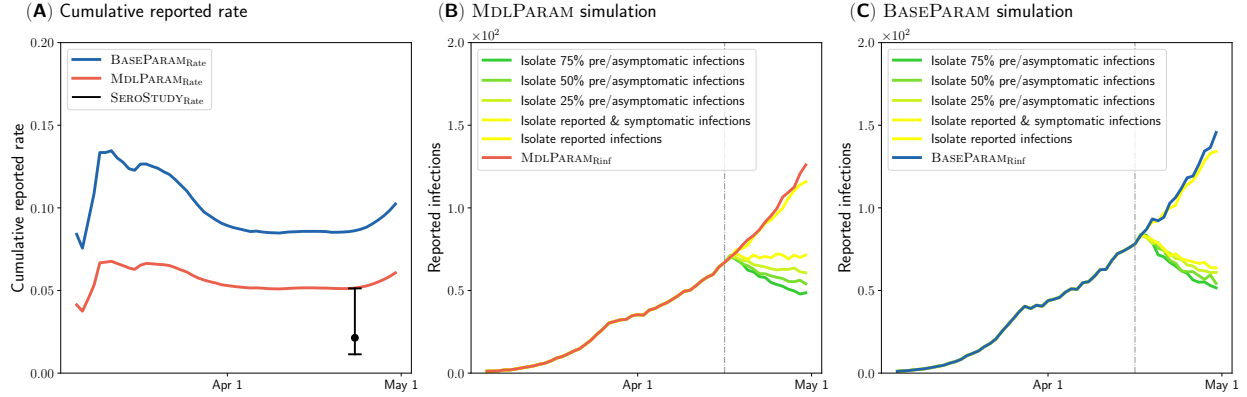


Figure 5: **(A)** MDLINFER estimates cumulative reported rate more accurately than BASEINFER: The blue and red curve represent the cumulative reported rate estimated by BASEINFER, $\text{BASEPARAM}_{\text{Rate}}$, and by MDLINFER, $\text{MDLPARAM}_{\text{Rate}}$, respectively. The black point estimate and its confidence interval represent the cumulative reported rate $\text{SEROSTUDY}_{\text{Rate}}$ estimated by serological studies [2, 26]. Note that both approaches try to fit the $\text{SEROSTUDY}_{\text{Rate}}$ without being informed with them. The results reveal that a large majority of COVID-19 infections were unreported. **(B)** MDLINFER reveals that non-pharmaceutical interventions (NPI) on asymptomatic and presymptomatic infections are essential to control the COVID-19 epidemic. Here, the red curve and other five curves represent the MDLINFER’s estimation of reported infections for no NPI scenario and 5 different NPI scenarios described in the Results section. The vertical grey dash line divides the observed period (left) and forecast period (right). **(C)** Inaccurate estimation by BASEINFER may lead to wrong non-pharmaceutical intervention conclusions. The blue curve and other five curves represent the BASEINFER’s estimation of reported infections for no NPI scenario and the same 5 scenarios in **(B)**.

Evaluating the effect of non-pharmaceutical Interventions

We have already shown that MDLINFER is able to estimate the number of total infections accurately. In the following three observations, we show that such accurate estimations are important for evaluating the effect of non-pharmaceutical interventions.

(D) MDLINFER reveals that a large majority of COVID-19 infections were unreported

We compute the *cumulative* reported rate $\text{MDLPARAM}_{\text{Rate}}$ measured by the ratio of the cumulative value of reported infections to the total infections estimated by MDLINFER over time and plotted it for Minneapolis-Spring-20 in Fig. 5A. The figure shows that the $\text{MDLPARAM}_{\text{Rate}}$ increases in early March, and then gradually decreases. This observation is explained by the community spread-driven COVID-19 outbreaks that were not reported until early March, which fits earlier studies [40].

(E) Non-pharmaceutical interventions on asymptomatic and presymptomatic infections are essential to control the COVID-19 epidemic

Our simulations show that non-pharmaceutical interventions on asymptomatic and presymptomatic infections are essential to control COVID-19. Here, we plot the simulated reported infections of MDLPARAM in Fig. 5B (red curve). We then repeat the simulation of reported infections for 5 different scenarios: (i) isolate just the reported infections, (ii) isolate just the symptomatic infec-

tions, and isolate symptomatic infections in addition to (iii) 25%, (iv) 50%, and (v) 75% of both asymptomatic and presymptomatic infections. In our setup, we assume that the infectivity reduces by half when a person is isolated. As seen in Fig. 5B, when only the reported infections are isolated, there is almost no change in the “future” reported infections. However, when we isolate both the reported and symptomatic infections, the reported infections decreases significantly. Even here, the reported infections are still not in decreasing trend. On the other hand, non-pharmaceutical interventions for some fraction of asymptomatic and presymptomatic infections make reported infections decrease. Thus, we can conclude that non-pharmaceutical interventions on asymptomatic infections are essential in controlling the COVID-19 epidemic.

(F) Accuracy of non-pharmaceutical intervention simulations relies on the good estimation of parameterization

Next, we also plot the simulated reported infections generated by BASEINFER in Fig. 5C (blue curve). As seen in the figure, based on BASEINFER, we can infer that only non-pharmaceutical interventions on symptomatic infections are enough to control the COVID-19 epidemic. However, this has been proven to be incorrect by prior studies and real-world observations [41]. Therefore, we can conclude that the accuracy of non-pharmaceutical intervention simulation relies on the quality of the learned parameterization.

Discussion and Future Work

This study proposes MDLINFER, a data-driven model selection approach that automatically estimates the number of total infections using epidemiological models. Our approach leverages the information theoretic Minimum Description Length (MDL) principle to select total infections that “best describe” the observed outbreak. Our approach addresses several gaps in current practice including the long-term infeasibility of serological studies [26], and ad-hoc assumptions in epidemiological models [33, 39, 44, 25].

Overall, our results show that MDLINFER estimates total infections at various geographical locations and different epidemiological models more accurately than BASEINFER from both directions, i.e., it corrects both over- and under-estimates. For example, compared to BASEINFER, we correctly estimate 55719 more infections by April 1 for the SEIR + HD model in Fig. 2F, and 87636 fewer infections for the SAPHIRE model in Fig. 2B for South Florida-Spring-20. We also show that MDLINFER leads to a better fit of the reported infections in the observed period and more accurate forecasts for the forecast period than BASEINFER. We reveal that a large majority of COVID-19 infections were unreported, where non-pharmaceutical interventions on unreported infections can help to mitigate the COVID-19 outbreak. We also show that MDLINFER estimates more accurate symptomatic rate trends than BASEINFER. Additionally, our results show consistent performance with respect to the reported infections and serological studies on both SAPHIRE and SEIR + HD model. We also show that MDLINFER identifies the ground truth parameters better than BASEINFER (see Supplementary Information section for details). As an aside, BASEINFER may also give uncertainty estimates for their calibrated parameterizations. Our framework MDLINFER can be adapted to generate such estimates as well (see Supplementary Information section for a demonstration).

The MDLINFER framework is likely to be helpful in the surveillance of COVID-19 in the near future, and for future epidemics. Even with the U.S. returning to normalcy, surveillance of the pandemic is still essential for public health. The daily incidence of COVID-19 has decreased from early 2021 to summer 2021, according to the CDC COVID Data Tracker portal [7, 35]. However, new variants of the SARS-CoV-2 (e.g., the Delta and Omicron variants) have been spreading rapidly [37,

48, 21]. Testing for these new variants and large-scale surveillance via laboratory tests may be limited and less systematic than what was done for COVID-19 before. In such settings, using our MDLINFER framework, epidemiologists and policymakers can improve the accuracy of estimates of total infections (without large-scale serological studies), as well as forecasts of their models.

One of the limitations of our work is that the benefits of using MDLINFER depends on the suitability of the epidemiological model. If the epidemiological model is not expressive enough for the observed data, then the gains from MDLINFER may not be significant. As a future work, it may be helpful to adapt MDLINFER to measure the quality of an epidemiological model. We also note that MDLINFER is built on ODE-based epidemiological models; other kinds of epidemic models, e.g., agent-based models [42, 28, 59, 18, 45], are more suitable in some settings. It would be interesting to extend MDLINFER to incorporate such models. Finally, there is significant population or spatial heterogeneity in disease outcomes [15, 31], e.g., differences in severity rate or mortality rate, when infected with COVID-19, for different age groups [22, 27], which has not been considered in our work.

To summarize, MDLINFER is a robust data-driven method to accurately estimate total infections, which will help data scientists, epidemiologists, and policy-makers to further improve existing ODE-based epidemiological models, make accurate forecasts, and combat the ongoing COVID-19 pandemic. More generally, MDLINFER opens up a new line of research in epidemic modeling using information theory.

Materials and Methods

Data

Datasets

We use the following publicly available datasets for our study:

1. **New York Times reported infections [3]:** This dataset (NYT-Rinf) consists of the daily time sequence of reported COVID-19 infections D_{reported} and the mortality $D_{\text{mortality}}$ (cumulative values) for each county in the US starting from January 21, 2020 to current.
2. **Serological studies [26, 2]:** This dataset consists of the point and 95% confidence interval estimates of the prevalence of antibodies to SARS-CoV-2 in 10 US locations every 3–4 weeks from March to July 2020. For each location, CDC works with commercial laboratories to collect the blood specimens in the population and test them for antibodies to SARS-CoV-2. Each specimen collection period ranges from 6 to 14 days. As suggested by prior work [32, 47], these serological studies have high sensitivity to antibodies for 6 months after infections. Hence, using the prevalence and total population in one location, we can compute the estimated total infections $\text{SEROSTUDY}_{\text{Tinf}}$ for the past 6 months (i.e., from the beginning of the pandemic since January 2020). However, this $\text{SEROSTUDY}_{\text{Tinf}}$ can not be compared with the epidemiological model estimated total infection numbers directly. The reasons are (i) the antibodies may take 10 to 14 days delay to be detectable after infection [65, 54] and (ii) the 6-14 range period for specimen collection as mentioned before. To account for this, we compare the $\text{SEROSTUDY}_{\text{Tinf}}$ numbers with the MDLINFER and BASEINFER estimated total infections of 7 days prior to the first day of specimen collection period as suggested by the CDC serological studies work [26].

3. **Symptomatic surveillance [51, 53]:** This dataset consists of point estimate $\text{RATE}_{\text{Symp}}$ and standard error of the COVID-related symptomatic rate for each county in the US starting from April 6, 2020 to date. The survey asks a series of questions on randomly sampled social media (Facebook) users to estimate the percentage of people who have a COVID-like symptoms such as the fever along with cough or shortness of breath or difficulty breathing on a given day. However, there are several caveats such as they could not cover all symptoms of COVID-19 and these symptoms can be also caused by many other conditions, due to which they are not expected to be unbiased estimates for the true symptomatic rate [51]. Besides, as the original symptomatic surveillance data is at a county level, we sum up the numbers to compute the $\text{RATE}_{\text{Symp}}$ and focus on trends instead of the exact numbers.

Our Approach

Two-part sender-receiver framework

In this work, we use two-part sender-receiver framework. The conceptual goal of the framework is to transmit the DATA from the possession of the hypothetical sender S to the hypothetical receiver R . We assume the sender does this by first sending a MODEL and then sending the DATA under this MODEL. In this MDL framework, we want to minimize the number of bits for this process. We do this by identifying the MODEL that encodes the DATA such that the total number of bits needed to encode both the MODEL and the DATA is minimized. Hence our cost function in the total number of bits needed is composed of two parts: (i) model cost $L(\text{MODEL})$: The cost in bits of encoding the MODEL and (ii) data cost $L(\text{DATA}|\text{MODEL})$: The cost in bits of encoding the DATA given the MODEL. Intuitively, the idea is that a good MODEL will lead to a fewer number of bits needed to encode both MODEL and DATA. We formulate the general MDL optimization problem as follows: Given the DATA, $L(\text{MODEL})$, and $L(\text{DATA}|\text{MODEL})$, find MODEL^* such that

$$\text{MODEL}^* = \arg \min_{\text{MODEL}} L(\text{MODEL}) + L(\text{DATA}|\text{MODEL}) \quad (1)$$

In our situation, the DATA is the reported COVID-19 infections D_{reported} : it is the only real-world data given to us. Note that total infections are not directly observed. As described in the introduction section, the MODEL is intuitively $(D, \alpha'_{\text{reported}})$. Here D refers to a candidate total infections time series, and $\alpha'_{\text{reported}}$ is the corresponding reported rate. Specifically, we calibrate O_M on (D, D_{reported}) using CALIBRATE to get the "candidate" parameterization Θ' , and then compute $\alpha'_{\text{reported}}$ from Θ' . Further, we choose to also add $\hat{\Theta}$ estimated by BASEINFER, making our MODEL to be $(D, \Theta', \hat{\Theta})$. There are alternative MODELS that can be considered, but we choose this MODEL = $(D, \Theta', \hat{\Theta})$ and explain more in the Supplementary Information. Note that as two-part MDL (and MDL in general) does not assume the nature of the DATA or the MODEL, our MDLINFER can be applied to any ODE model. We have also discussed intuitive advantages of the MDLINFER over BASEINFER briefly in the introduction section (see Supplementary Information for more details). Next, we give more details how to formulate our problem of estimating total infections D .

MDL formulation

First, we need to introduce some notations. Given an epidemiological model O_M and the parameterization $\hat{\Theta}$ estimated by BASEINFER, we can compute the reported infections. However, this is only an estimate of the reported infections rather than the exact D_{reported} . This is because even though we have already calibrated O_M using D_{reported} , the calibration cannot be perfect, and there will be

differences between these estimated reported infections and D_{reported} . Here, we term this estimated reported infections as $D_{\text{reported}}(\hat{\Theta})$. We can also estimate the total infections $D(\hat{\Theta})$ for O_M in the same way. Similarly, we have the $D_{\text{reported}}(\Theta')$ and $D(\Theta')$ for Θ' . As described in the introduction section, we can also calculate the reported rate $\hat{\alpha}_{\text{reported}}$ and $\alpha'_{\text{reported}}$ using $\hat{\Theta}$ and Θ' . With these notations, next we will formulate the space of all possible MODELS and give the equation for the cost in bits of encoding MODEL and DATA.

MODEL space

We have MODEL = $(D, \Theta', \hat{\Theta})$ as described above. Hence our MODEL space will be all possible daily sequences for D and all possible parameterizations for Θ' and $\hat{\Theta}$. The MDL framework will search in this space to find the MODEL*.

Model cost

With MODEL = $(D, \Theta', \hat{\Theta})$, we conceptualize the model cost by imagining that the sender S will send the MODEL = $(D, \Theta', \hat{\Theta})$ to the receiver R in three parts: (i) first send the $\hat{\Theta}$ by encoding $\hat{\Theta}$ directly (ii) next send the Θ' given $\hat{\Theta}$ by encoding $\Theta' - \hat{\Theta}$ and (iii) then send D given Θ' and $\hat{\Theta}$ by encoding $\alpha'_{\text{reported}} \times D - D_{\text{reported}}(\hat{\Theta})$. Intuitively, both $\alpha'_{\text{reported}} \times D$ and $D_{\text{reported}}(\hat{\Theta})$ should be close to D_{reported} , and the receiver could recover the D using $\hat{\Theta}$, $\alpha'_{\text{reported}}$, and $D_{\text{reported}}(\hat{\Theta})$ as they have already been sent. We term the model cost as $L(D, \Theta', \hat{\Theta})$ with three components: $\text{COST}(\hat{\Theta})$, $\text{COST}(\Theta'|\hat{\Theta})$, and $\text{COST}(D|\Theta', \hat{\Theta})$. Hence,

$$L(D, \Theta', \hat{\Theta}) = \text{COST}(\hat{\Theta}) + \text{COST}(\Theta' - \hat{\Theta}|\hat{\Theta}) + \text{COST}(\alpha'_{\text{reported}} \times D - D_{\text{reported}}(\hat{\Theta})|\Theta', \hat{\Theta}) \quad (2)$$

For Equation 2, the $\text{COST}(\cdot)$ function gives the total number of bits we need to spend in encoding each term. The details of the encoding method can be found in the Supplementary Information.

Data cost

We need to send the DATA = D_{reported} next given the MODEL. Given MODEL = $(D, \Theta', \hat{\Theta})$, we send DATA by encoding $\frac{D - D_{\text{reported}}}{1 - \alpha'_{\text{reported}}} - D(\Theta')$. Intuitively, $D - D_{\text{reported}}$ corresponds to the unreported infections, and $1 - \alpha'_{\text{reported}}$ is the unreported rate. Therefore, $\frac{D - D_{\text{reported}}}{1 - \alpha'_{\text{reported}}}$ should be close to the total infections D and $D(\Theta')$. The receiver could also recover the D_{reported} using D , $\alpha'_{\text{reported}}$, and $D(\Theta')$ as they have already been sent. We term data cost as $L(D_{\text{reported}}|D, \Theta', \hat{\Theta})$ and formulate it as Equation 3.

$$L(D_{\text{reported}}|D, \Theta', \hat{\Theta}) = \text{COST}\left(\frac{D - D_{\text{reported}}}{1 - \alpha'_{\text{reported}}} - D(\Theta')|D, \Theta', \hat{\Theta}\right) \quad (3)$$

Total cost

With $L(D, \Theta', \hat{\Theta})$ as in Equation 2 and $L(D_{\text{reported}}|D, \Theta', \hat{\Theta})$ as in Equation 3 above, the total cost $L(D_{\text{reported}}, D, \Theta', \hat{\Theta})$ is:

$$\begin{aligned}
L(D_{\text{reported}}, D, \Theta', \hat{\Theta}) &= L(D, \Theta', \hat{\Theta}) + L(D_{\text{reported}}|D, \Theta', \hat{\Theta}) \\
&= \text{COST}(\hat{\Theta}) + \text{COST}(\Theta' - \hat{\Theta}|\hat{\Theta}) \\
&+ \text{COST}(\alpha'_{\text{reported}} \times D - D_{\text{reported}}(\hat{\Theta})|\Theta', \hat{\Theta}) \\
&+ \text{COST}\left(\frac{D - D_{\text{reported}}}{1 - \alpha'_{\text{reported}}} - D(\Theta')|D, \Theta', \hat{\Theta}\right)
\end{aligned} \tag{4}$$

Problem statement

Note that our main objective is to estimate the total infections D . With $L(D_{\text{reported}}, D, \Theta', \hat{\Theta})$ formulated in Equation 4, we can state the problem as: Given the time sequence D_{reported} , epidemiological model O_M , and a calibration procedure CALIBRATE, find D^* that minimizes the MDL total cost i.e.

$$D^* = \arg \min_D L(D_{\text{reported}}, D, \Theta', \hat{\Theta}) \tag{5}$$

Algorithm

Next, we will present our algorithm to solve the problem in Equation 5. Note that directly searching D^* naively is intractable since D^* is a daily sequence not a scalar. Instead, we propose first finding a ‘‘good enough’’ reported rate $\alpha^*_{\text{reported}}$ quickly with the constraint $D = \frac{D_{\text{reported}}}{\alpha^*_{\text{reported}}}$ to reduce the search space. Then with this $\alpha^*_{\text{reported}}$, we can search for the optimal D^* in Equation 5. Hence we propose a two-step algorithm: (i) do a linear search to find a good reported rate $\alpha^*_{\text{reported}}$ (ii) given the $\alpha^*_{\text{reported}}$ found above, use an optimization method to find the D^* that minimizes $L(D_{\text{reported}}, D, \Theta', \hat{\Theta})$ with $\alpha^*_{\text{reported}}$ constraints.

Step 1: Find the $\alpha^*_{\text{reported}}$

In step 1, we do a linear search on α_{reported} to find the $\alpha^*_{\text{reported}}$. As stated before, we use $\frac{D_{\text{reported}}}{\alpha_{\text{reported}}}$ as D in $L(D_{\text{reported}}, D, \Theta', \hat{\Theta})$ to help reduce the search space. Here, we formulate step 1 algorithm as Equation 6.

$$\alpha^*_{\text{reported}} = \arg \min_{\alpha_{\text{reported}}} L\left(D_{\text{reported}}, \frac{D_{\text{reported}}}{\alpha_{\text{reported}}}, \Theta', \hat{\Theta}\right) \tag{6}$$

Step 2: Find the D^* given $\alpha^*_{\text{reported}}$

With the $\alpha^*_{\text{reported}}$ found in step 1, we next find the D^* that minimizes the $L(D_{\text{reported}}, D, \Theta', \hat{\Theta})$. Note that we have already found a good $\alpha^*_{\text{reported}}$, we can constrain the D^* to ensure that the sum of D^* equals to the sum of $\frac{D_{\text{reported}}}{\alpha^*_{\text{reported}}}$. We use the Nelder-Mead method [20] to solve this constrained optimization problem for D^* . Here, we formulate step 2 algorithm as Equation 7.

$$D^* = \arg \min_D L(D_{\text{reported}}, D, \Theta', \hat{\Theta}) \tag{7}$$

We describe the two-step algorithm in more detail in the Supplementary Information.

BASEINFER and MDLINFER formulation

Here, we also give the mathematical formulations for BASEINFER and MDLINFER. As described in the introduction section, given an epidemiological model O_M , a typical approach is to calibrate the O_M to D_{reported} using the calibration procedure CALIBRATE. We call this methodology as BASEINFER(O_M , CALIBRATE, D_{reported}). As in Equation 8, the output of BASEINFER is the baseline parameterization (BASEPARAM) $\hat{\Theta}$.

$$\begin{aligned}\hat{\Theta} &= \text{BASEINFER}(O_M, \text{CALIBRATE}, D_{\text{reported}}) \\ &= \text{CALIBRATE}(O_M, D_{\text{reported}})\end{aligned}\tag{8}$$

As for the MDLINFER, it also takes the same input (O_M , CALIBRATE, D_{reported}) as BASEINFER. Assume we are given the total infections D , we calibrate the O_M on (D, D_{reported}) to get a "candidate" parameterization Θ' in Equation 9.

$$\Theta' = \text{CALIBRATE}(O_M, (D, D_{\text{reported}}))\tag{9}$$

However, we are not given the D . Hence, we use the MDL framework to find such D^* as in Equation 7. With such D^* , we could finally calibrate the O_M on $(D^*, D_{\text{reported}})$ and gets another parameterization Θ^* . As in Equation 10, we call Θ^* as the MDL parameterization, or MDLPARAM.

$$\begin{aligned}\Theta^* &= \text{MDLINFER}(O_M, \text{CALIBRATE}, D_{\text{reported}}) \\ &= \text{CALIBRATE}(O_M, (D^*, D_{\text{reported}}))\end{aligned}\tag{10}$$

where $D^* = \arg \min_D L(D_{\text{reported}}, D, \Theta', \hat{\Theta})$. Intuitively, if $\hat{\Theta}$ estimated by BASEINFER is perfect, MDLINFER will also give the same Θ^* as $\hat{\Theta}$.

Epidemiological models

Next, we describe the two epidemiological models we use in our experiments: SEIR + HD and SAPHIRE model. SEIR + HD [33] consists of 10 states: Susceptible S , exposed E , pre-symptomatic I_P , severe symptomatic I_S , mild symptomatic I_M , asymptomatic I_A , hospitalized (eventual death) H_D , hospitalized (eventual recover) H_R , recovered R , and dead D . The parameters to be calibrated are the transmission rate β_0 (the transmission rate in the absence of interventions), σ (the proportional reduction on β_0 under shelter-in-place), and E_0 (number of initial infections). The other parameters are fixed and given. They assume the importations only happen at the beginning of the pandemic (captured by E_0), and the total population N remains constant. We also extend SEIR + HD model to infer two more parameters: α (proportion of asymptomatic infections) and α_1 (proportion of new symptomatic infections that are reported). We compute the new reported infections and unreported infections as follows:

1. New reported infections = $\alpha_1 \times (N_{I_P I_S} + N_{I_P I_M})$: Here $N_{I_P I_S} + N_{I_P I_M}$ is the number of new symptomatic infections everyday. $N_{I_P I_S}$ is the number of patients switching their state from I_P to I_S (and similarly for $N_{I_P I_M}$). We assume α_1 proportion of new symptomatic infections every day are reported.
2. New unreported infections = $(1 - \alpha_1) \times (N_{I_P I_S} + N_{I_P I_M}) + N_{E I_A}$.

SAPHIRE [25] consists of 7 states: Susceptible S , exposed E , pre-symptomatic P , ascertained infectious I , unascertained infectious A , hospitalized H , and recovered R . The parameters to be calibrated are the transmission rate β and reported rate r while keeping other parameters fixed as given values. We also compute the new reported infections and unreported infections as follows:

1. New reported infections = $\frac{rP}{D_p}$: Here $\frac{P}{D_p}$ is the number of new infections from pre-symptomatic every day. D_p is the parameter for the presymptomatic infectious period and is fixed. r is the reported rate estimated by the epidemiological model.
2. New unreported infections = $\frac{(1-r)P}{D_p}$.

Estimating infections using BASEPARAM and MDLPARAM

Here, we describe how we get the estimations in the results section using BASEPARAM and MDLPARAM. Here we use the BASEPARAM from BASEINFER as the example (this can also be repeated for MDLPARAM for MDLINFER). Using the epidemiological model O_M , we can calculate the BASEPARAM's estimation of total infections $\text{BASEPARAM}_{\text{Tinf}}$ as the cumulative values of $D(\hat{\Theta})$ from pandemic's beginning. $D_{\text{reported}}(\hat{\Theta})$ can be directly used as the BASEPARAM's estimation of reported infections. For the cumulative reported rate $\text{BASEPARAM}_{\text{Rate}}$, we calculate it as the cumulative values of NYT-Rinf divided by $D(\hat{\Theta})$. For the symptomatic rate, SEIR + HD model [33] could estimate the number of symptomatic rate $\text{BASEPARAM}_{\text{Symp}}$ by dividing the number of infections in state I_S and I_M by the population number. However, SAPHIRE model [25] does not contain states that correspond to the symptomatic cases, so we cannot estimate the symptomatic rate using this model.

Acknowledgements

This paper was supported in part by the NSF (Expeditions CCF-1918770 and CCF-1918656, CAREER IIS-2028586, RAPID IIS-2027862, Medium IIS-1955883, Medium IIS-2106961, IIS-1931628, IIS-1955797, IIS-2027848, PIPP CCF-2200269), NIH 2R01GM109718, CDC MInD program U01CK000589, ORNL and funds/computing resources from Georgia Tech and GTRI. B. A. was in part supported by the CDC MInD-Healthcare U01CK000531-Supplement. A.V.'s work is also supported in part by grants from the UVA Global Infectious Diseases Institute (GIDI). J.V. is institutionally funded by CISPA.

References

- [1] 4q gdp: Economy expands at a 4.0% annualized rate. <https://finance.yahoo.com/news/4q-gdp-2020-us-economy-coronavirus-pandemic-180133456.html>.
- [2] Commercial laboratory seroprevalence survey data. <https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/commercial-lab-surveys.html>.
- [3] Coronavirus in the u.s.:latest map and case count. <https://www.nytimes.com/interactive/2020/us/coronavirus-us-cases.html>.
- [4] Covid-19 test price information, <https://www.questdiagnostics.com/business-solutions/health-plans/covid-19/pricing>.
- [5] Hidden outbreaks spread through u.s. cities far earlier than americans knew, estimates say. <https://www.nytimes.com/2020/04/23/us/coronavirus-early-outbreaks-cities.html>.
- [6] Interim economic projections for 2020 and 2021. <https://www.cbo.gov/publication/56351>.

- [7] United states covid-19 cases, deaths, and laboratory testing (naats) by state, territory, and jurisdiction. <https://covid.cdc.gov/covid-data-tracker/#cases>.
- [8] Us covid cases likely more than double official count, experts say. <https://www.cidrap.umn.edu/news-perspective/2021/07/us-covid-cases-likely-more-double-official-count-experts-say>, 2021.
- [9] ACCORSI, E. K., QIU, X., RUMPLER, E., KENNEDY-SHAFFER, L., KAHN, R., JOSHI, K., GOLDSTEIN, E., STENSRUD, M. J., NIEHUS, R., CEVIK, M., ET AL. How to detect and reduce potential sources of biases in studies of sars-cov-2 and covid-19. *European Journal of Epidemiology* (2021), 1–18.
- [10] ADHIKARI, B., RANGUDU, P., PRAKASH, B. A., AND VULLIKANTI, A. Near-optimal mapping of network states using probes. In *Proceedings of the 2018 SIAM International Conference on Data Mining* (2018), SIAM, pp. 108–116.
- [11] AGUILAR, J. B., FAUST, J. S., WESTAFER, L. M., AND GUTIERREZ, J. B. Investigating the impact of asymptomatic carriers on covid-19 transmission. *MedRxiv* (2020).
- [12] ANGELOPOULOS, A. N., PATHAK, R., VARMA, R., AND JORDAN, M. I. On identifying and mitigating bias in the estimation of the covid-19 case fatality rate. *arXiv preprint arXiv:2003.08592* (2020).
- [13] BAI, Y., YAO, L., WEI, T., TIAN, F., JIN, D.-Y., CHEN, L., AND WANG, M. Presumed asymptomatic carrier transmission of covid-19. *JAMA* *323*, 14 (2020), 1406–1407.
- [14] BEDFORD, T., GRENINGER, A. L., ROYCHOUDHURY, P., STARITA, L. M., FAMULARE, M., HUANG, M.-L., NALLA, A., PEPPER, G., REINHARDT, A., XIE, H., ET AL. Cryptic transmission of sars-cov-2 in washington state. *Science* *370*, 6516 (2020), 571–575.
- [15] BI, Q., WU, Y., MEI, S., YE, C., ZOU, X., ZHANG, Z., LIU, X., WEI, L., TRUELOVE, S. A., ZHANG, T., ET AL. Epidemiology and transmission of covid-19 in 391 cases and 1286 of their close contacts in shenzhen, china: a retrospective cohort study. *The Lancet infectious diseases* *20*, 8 (2020), 911–919.
- [16] BUDHATHOKI, K., AND VREEKEN, J. Origo: causal inference by compression. *Knowledge and Information Systems* *56*, 2 (2018), 285–307.
- [17] CAO, Q., AND HEYDARI, B. Micro-level social structures and the success of covid-19 national policies. *Nature Computational Science* *2*, 9 (2022), 595–604.
- [18] CHANG, S., PIERSON, E., KOH, P. W., GERARDIN, J., REDBIRD, B., GRUSKY, D., AND LESKOVEC, J. Mobility network models of covid-19 explain inequities and inform reopening. *Nature* *589*, 7840 (2021), 82–87.
- [19] DONG, E., DU, H., AND GARDNER, L. An interactive web-based dashboard to track covid-19 in real time. *The Lancet infectious diseases* *20*, 5 (2020), 533–534.
- [20] GAO, F., AND HAN, L. Implementing the nelder-mead simplex algorithm with adaptive parameters. *Computational Optimization and Applications* *51*, 1 (2012), 259–277.

- [21] GEERS, D., SHAMIER, M. C., BOGERS, S., DEN HARTOG, G., GOMMERS, L., NIEUWKOOP, N. N., SCHMITZ, K. S., RIJSBERGEN, L. C., VAN OSCH, J. A., DIJKHUIZEN, E., ET AL. Sars-cov-2 variants of concern partially escape humoral but not t-cell responses in covid-19 convalescent donors and vaccinees. *Science Immunology* 6, 59 (2021).
- [22] GOLDSTEIN, J. R., AND LEE, R. D. Demographic perspectives on the mortality of covid-19 and other epidemics. *Proceedings of the National Academy of Sciences* 117, 36 (2020), 22035–22041.
- [23] GOPALAKRISHNAN, V., PETHE, S., KEFAYATI, S., SRINIVASAN, R., HAKE, P., DESHPANDE, A., LIU, X., HOANG, E., DAVILA, M., BIANCO, S., ET AL. Globally local: Hyper-local modeling for accurate forecast of covid-19. *Epidemics* 37 (2021), 100510.
- [24] GRÜNWARD, P. D. *The minimum description length principle*. MIT press, 2007.
- [25] HAO, X., CHENG, S., WU, D., WU, T., LIN, X., AND WANG, C. Reconstruction of the full transmission dynamics of covid-19 in wuhan. *Nature* 584, 7821 (2020), 420–424.
- [26] HAVERS, F. P., REED, C., LIM, T., MONTGOMERY, J. M., KLENA, J. D., HALL, A. J., FRY, A. M., CANNON, D. L., CHIANG, C.-F., GIBBONS, A., ET AL. Seroprevalence of antibodies to sars-cov-2 in 10 sites in the united states, march 23-may 12, 2020. *JAMA internal medicine* 180, 12 (2020), 1576–1586.
- [27] HO, F. K., PETERMANN-ROCHA, F., GRAY, S. R., JANI, B. D., KATIKIREDDI, S. V., NIEDZWIEDZ, C. L., FOSTER, H., HASTIE, C. E., MACKAY, D. F., GILL, J. M., ET AL. Is older age associated with covid-19 mortality in the absence of other risk factors? general population cohort study of 470,034 participants. *PloS one* 15, 11 (2020), e0241824.
- [28] HOERTEL, N., BLACHIER, M., BLANCO, C., OLFSON, M., MASSETTI, M., RICO, M. S., LIMOSIN, F., AND LELEU, H. A stochastic agent-based model of the sars-cov-2 epidemic in france. *Nature medicine* 26, 9 (2020), 1417–1421.
- [29] IONIDES, E. L., NGUYEN, D., ATCHADÉ, Y., STOEV, S., AND KING, A. A. Inference for dynamic and latent variable models via iterated, perturbed bayes maps. *Proceedings of the National Academy of Sciences* 112, 3 (2015), 719–724.
- [30] IRONS, N. J., AND RAFTERY, A. E. Estimating sars-cov-2 infections from deaths, confirmed cases, tests, and random surveys. *Proceedings of the National Academy of Sciences* 118, 31 (2021).
- [31] JAY, J., BOR, J., NSOESIE, E. O., LIPSON, S. K., JONES, D. K., GALEA, S., AND RAIFMAN, J. Neighbourhood income and physical distancing during the covid-19 pandemic in the united states. *Nature human behaviour* 4, 12 (2020), 1294–1302.
- [32] 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. Estimated us infection-and vaccine-induced sars-cov-2 seroprevalence based on blood donations, july 2020-may 2021. *JAMA* 326, 14 (2021), 1400–1409.
- [33] KAIN, M. P., CHILDS, M. L., BECKER, A. D., AND MORDECAI, E. A. Chopping the tail: How preventing superspreading can help to maintain covid-19 control. *Epidemics* 34 (2021), 100430.

- [34] KOUTRA, D., KANG, U., VREEKEN, J., AND FALOUTSOS, C. Vog: Summarizing and understanding large graphs. In *Proceedings of the 2014 SIAM international conference on data mining* (2014), SIAM, pp. 91–99.
- [35] KRAEMER, M. U., SCARPINO, S. V., MARIVATE, V., GUTIERREZ, B., XU, B., LEE, G., HAWKINS, J. B., RIVERS, C., PIGOTT, D. M., KATZ, R., ET AL. Data curation during a pandemic and lessons learned from covid-19. *Nature Computational Science* 1, 1 (2021), 9–10.
- [36] KRAEMER, M. U., YANG, C.-H., GUTIERREZ, B., WU, C.-H., KLEIN, B., PIGOTT, D. M., DU PLESSIS, L., FARIA, N. R., LI, R., HANAGE, W. P., ET AL. The effect of human mobility and control measures on the covid-19 epidemic in china. *Science* 368, 6490 (2020), 493–497.
- [37] KUSTIN, T., HAREL, N., FINKEL, U., PERCHIK, S., HARARI, S., TAHOR, M., CASPI, I., LEVY, R., LESHCHINSKY, M., KEN DROR, S., ET AL. Evidence for increased breakthrough rates of sars-cov-2 variants of concern in bnt162b2-mrna-vaccinated individuals. *Nature medicine* 27, 8 (2021), 1379–1384.
- [38] LAI, S., RUKTANONCHAI, N. W., ZHOU, L., PROSPER, O., LUO, W., FLOYD, J. R., WESOLOWSKI, A., SANTILLANA, M., ZHANG, C., DU, X., ET AL. Effect of non-pharmaceutical interventions to contain covid-19 in china. *Nature* 585, 7825 (2020), 410–413.
- [39] LI, R., PEI, S., CHEN, B., SONG, Y., ZHANG, T., YANG, W., AND SHAMAN, J. Substantial undocumented infection facilitates the rapid dissemination of novel coronavirus (sars-cov-2). *Science* 368, 6490 (2020), 489–493.
- [40] LU, F. S., NGUYEN, A. T., LINK, N. B., MOLINA, M., DAVIS, J. T., CHINAZZI, M., XIONG, X., VESPIGNANI, A., LIPSITCH, M., AND SANTILLANA, M. Estimating the cumulative incidence of covid-19 in the united states using influenza surveillance, virologic testing, and mortality data: Four complementary approaches. *PLOS Computational Biology* 17, 6 (2021), e1008994.
- [41] MOGHADAS, S. M., FITZPATRICK, M. C., SAH, P., PANDEY, A., SHOUKAT, A., SINGER, B. H., AND GALVANI, A. P. The implications of silent transmission for the control of covid-19 outbreaks. *Proceedings of the National Academy of Sciences* 117, 30 (2020), 17513–17515.
- [42] NANDE, A., SHEEN, J., WALTERS, E. L., KLEIN, B., CHINAZZI, M., GHEORGHE, A. H., ADLAM, B., SHINNICK, J., TEJEDA, M. F., SCARPINO, S. V., ET AL. The effect of eviction moratoria on the transmission of sars-cov-2. *Nature communications* 12, 1 (2021), 1–13.
- [43] PADMANABHAN, P., DESIKAN, R., AND DIXIT, N. M. Modeling how antibody responses may determine the efficacy of covid-19 vaccines. *Nature Computational Science* 2, 2 (2022), 123–131.
- [44] PEI, S., KANDULA, S., AND SHAMAN, J. Differential effects of intervention timing on covid-19 spread in the united states. *Science advances* 6, 49 (2020), eabd6370.
- [45] PEI, S., TENG, X., LEWIS, P., AND SHAMAN, J. Optimizing respiratory virus surveillance networks using uncertainty propagation. *Nature communications* 12, 1 (2021), 1–10.
- [46] PEI, S., YAMANA, T. K., KANDULA, S., GALANTI, M., AND SHAMAN, J. Burden and characteristics of covid-19 in the united states during 2020. *Nature* 598, 7880 (2021), 338–341.

- [47] PELUSO, M. J., TAKAHASHI, S., HAKIM, J., KELLY, J. D., TORRES, L., IYER, N. S., TURCIOS, K., JANSON, O., MUNTER, S. E., THANH, C., ET AL. Sars-cov-2 antibody magnitude and detectability are driven by disease severity, timing, and assay. *Science advances* 7, 31 (2021), eabh3409.
- [48] PLANAS, D., VEYER, D., BAIDALIUK, A., STAROPOLI, I., GUIVEL-BENHASSINE, F., RAJAH, M. M., PLANCHAIS, C., PORROT, F., ROBILLARD, N., PUECH, J., ET AL. Reduced sensitivity of sars-cov-2 variant delta to antibody neutralization. *Nature* 596, 7871 (2021), 276–280.
- [49] PRAKASH, B. A., VREEKEN, J., AND FALOUTSOS, C. Spotting culprits in epidemics: How many and which ones? In *2012 IEEE 12th International Conference on Data Mining* (2012), IEEE, pp. 11–20.
- [50] PRESS, W. H., AND LEVIN, R. C. Modeling, post covid-19. *Science* 370, 6520 (2020), 1015–1015.
- [51] REINHART, A., BROOKS, L., JAHJA, M., RUMACK, A., TANG, J., AGRAWAL, S., AL SAEED, W., ARNOLD, T., BASU, A., BIEN, J., ET AL. An open repository of real-time covid-19 indicators. *Proceedings of the National Academy of Sciences* 118, 51 (2021).
- [52] RUSSELL, T. W., GOLDING, N., HELLEWELL, J., ABBOTT, S., WRIGHT, L., PEARSON, C. A., VAN ZANDVOORT, K., JARVIS, C. I., GIBBS, H., LIU, Y., ET AL. Reconstructing the early global dynamics of under-ascertained covid-19 cases and infections. *BMC medicine* 18, 1 (2020), 1–9.
- [53] SALOMON, J. A., REINHART, A., BILINSKI, A., CHUA, E. J., LA MOTTE-KERR, W., RÖNN, M. M., REITSMA, M. B., MORRIS, K. A., LARocca, S., FARAG, T. H., ET AL. The us covid-19 trends and impact survey: Continuous real-time measurement of covid-19 symptoms, risks, protective behaviors, testing, and vaccination. *Proceedings of the National Academy of Sciences* 118, 51 (2021).
- [54] SETHURAMAN, N., JEREMIAH, S. S., AND RYO, A. Interpreting diagnostic tests for sars-cov-2. *JAMA* 323, 22 (2020), 2249–2251.
- [55] SHAMAN, J. An estimation of undetected covid cases in france. *Nature* 590 (2020), 38–39.
- [56] SOOD, N., SIMON, P., EBNER, P., EICHNER, D., REYNOLDS, J., BENDAVID, E., AND BHATTACHARYA, J. Seroprevalence of sars-cov-2-specific antibodies among adults in los angeles county, california, on april 10-11, 2020. *JAMA* 323, 23 (2020), 2425–2427.
- [57] STOCKMAIER, S., STROEYMEYER, N., SHATTUCK, E. C., HAWLEY, D. M., MEYERS, L. A., AND BOLNICK, D. I. Infectious diseases and social distancing in nature. *Science* 371, 6533 (2021).
- [58] SUBRAMANIAN, R., HE, Q., AND PASCUAL, M. Quantifying asymptomatic infection and transmission of covid-19 in new york city using observed cases, serology, and testing capacity. *Proceedings of the National Academy of Sciences* 118, 9 (2021).
- [59] TIAN, Y., SRIDHAR, A., YAĞAN, O., AND POOR, H. V. Analysis of the impact of mask-wearing in viral spread: Implications for covid-19. In *2021 American Control Conference (ACC)* (2021), IEEE, pp. 3132–3137.

- [60] TIWARI, S., VYASARAYANI, C., AND CHATTERJEE, A. Data suggest covid-19 affected numbers greatly exceeded detected numbers, in four european countries, as per a delayed seiqr model. *Scientific reports* 11, 1 (2021), 1–12.
- [61] WELLS, C. R., SAH, P., MOGHADAS, S. M., PANDEY, A., SHOUKAT, A., WANG, Y., WANG, Z., MEYERS, L. A., SINGER, B. H., AND GALVANI, A. P. Impact of international travel and border control measures on the global spread of the novel 2019 coronavirus outbreak. *Proceedings of the National Academy of Sciences* 117, 13 (2020), 7504–7509.
- [62] WILDER, B., CHARPIGNON, M., KILLIAN, J. A., OU, H.-C., MATE, A., JABBARI, S., PERRAULT, A., DESAI, A. N., TAMBE, M., AND MAJUMDER, M. S. Modeling between-population variation in covid-19 dynamics in hubei, lombardy, and new york city. *Proceedings of the National Academy of Sciences* 117, 41 (2020), 25904–25910.
- [63] WU, S. L., MERTENS, A. N., CRIDER, Y. S., NGUYEN, A., POKPONGKIAT, N. N., DJAJADI, S., SETH, A., HSIANG, M. S., COLFORD, J. M., REINGOLD, A., ET AL. Substantial underestimation of sars-cov-2 infection in the united states. *Nature communications* 11, 1 (2020), 1–10.
- [64] ZHANG, W., GOVINDAVARI, J. P., DAVIS, B. D., CHEN, S. S., KIM, J. T., SONG, J., LOPATEGUI, J., PLUMMER, J. T., AND VAIL, E. Analysis of genomic characteristics and transmission routes of patients with confirmed sars-cov-2 in southern california during the early stage of the us covid-19 pandemic. *JAMA network open* 3, 10 (2020), e2024191.
- [65] ZHAO, J., YUAN, Q., WANG, H., LIU, W., LIAO, X., SU, Y., WANG, X., YUAN, J., LI, T., LI, J., ET AL. Antibody responses to sars-cov-2 in patients with novel coronavirus disease 2019. *Clinical infectious diseases* 71, 16 (2020), 2027–2034.