360 **COVID-19 infection data encode a dynamic reproduction number in response to**
361 **policy decisions with secondary wave implications**
362
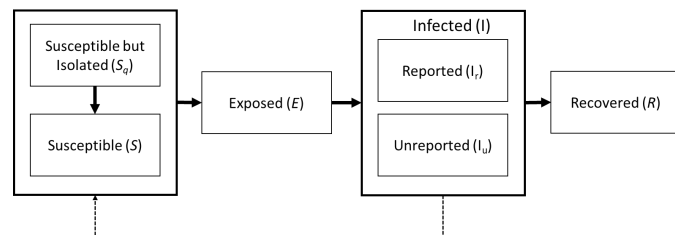363 **Supplemental Material**
364
365 Michael A. Rowland, Todd Swannack, Michael L. Mayo, Matthew Parno, Matthew Farthing, Ian
366 Dettwiller, Glover George, Molly Reif, Jeffrey Cegan, Benjamin Trump, Igor Linkov, Brandon Laf-
367 ferty, Todd Bridges
368
369 **Model overview**
370 The ERDC SEIR model simulates the number of infections in a state or metropolitan area. The
371 model is a modification of classic SEIR models that distinguishes between reported and unre-
372 ported cases, much like the approach employed by the Columbia University team (Li et al, 2020).
373 The ERDC SEIR model however, also employees another compartment to represent people that
374 are isolated from larger population and will not be exposed to the virus. The model parameters
375 are fit independently for each state and metropolitan area to match time-series data of con-
376 firmed infection reports. This calibration process is repeated daily as new data becomes available
377 and is performed on ERDC's high performance computing cluster. More details are provided be-
378 low.

379
380 **Model formulation**

381 The progression of the disease by the
382 transition of individuals in a population
383 through 5 states, as illustrated in Figure
384 S1. These states are: Susceptible (S), Ex-
385 posed (E), Infected (I) or Recovered/Re-
386 moved (R). The model makes a number
387 of assumptions. First, it assumes that all
388 individuals of a population can be
389 treated identically; that is, there is no



**Figure S1.** Conceptual model of the disease states. Healthy individuals are exposed to COVID-19 through infected individuals, and only a fraction of symptomatic individuals receive a test. All infected individuals "recover," which we use to account for those who become immune from further infection or die from the symptoms.

390 population variation in the transition rates between infection states. Another assumption of this
391 type of model is that individuals are "well mixed," by which we assume that rates of disease
392 progression do not depend upon the positions of any individuals within the subpopulations.
393 Therefore, we do not account for any spatial heterogeneity in population density (e.g., cities vs.
394 rural), which would otherwise directly affect the frequency of individual interactions that drive
395 disease spread. Another important assumption we make is that populations are large enough
396 that fluctuations driven by some of the individual variation in disease progression are small and
397 can be ignored, which is common in population and epidemiological modeling. As a consequence,
398 the individuals move between disease subpopulations (e.g., susceptible transitions to exposed)
399 at deterministic rates. Finally, we have imposed a constraint in which individuals cannot return
400 to a previous disease state, which reflects our presumption that those who recover are immune
401 from further infection. Despite these various approximations and simplifying assumptions, the

402 model is flexible enough for adaption to future scenarios. For example, population density can
403 be accounted for using this conceptual approach by creating additional compartments to de-
404 scribe subpopulations at a scale under which further variation is ignored. However, such changes
405 come at a cost of introducing additional model parameters, which is typically undesirable due to
406 concerns of overfitting and non-uniqueness.
407
408 More specifically, these assumptions lead to four ordinary differential equations (ODEs), each of
409 which define evolution of the four disease states:
410

$$\frac{dS_q}{dt} = \left[ \sum_{m=1}^{M} \gamma_m S \delta(t, t_m) \right]$$

$$\frac{dS}{dt} = -\frac{\beta S I_r}{N} - \frac{\mu \beta S I_u}{N} - \left[ \sum_{m=1}^{M} \gamma_m S \delta(t, t_m) \right]$$

$$\frac{dE}{dt} = \frac{\beta S I_r}{N} + \frac{\mu \beta S I_u}{N} - \frac{E}{Z}$$

$$\frac{dI_r}{dt} = \alpha \frac{E}{Z} - \frac{I_r}{D}$$

$$\frac{dI_u}{dt} = (1 - \alpha) \frac{E}{Z} - \frac{I_u}{D},$$

411
412 Collectively, Eqs. [1]-[4] describe evolution of the number of individuals susceptible to the dis-
413 ease, $S$, exposed to the disease after contact with infected individuals, $E$, those individuals with
414 reported/tested infections, $I_r$, and individuals with unreported infections, $I_u$. The overall number
415 of ACTIVE cases being tracked by authorities each day is given by $I_r(t - T_{delay})$; i.e., the number
416 of infected individuals being tracked (Eq. [3]) is shifted deterministically by an amount equal to
417 the delay time, $T_{delay}$. The current model is entirely deterministic, which may be stochastically
418 adjusted in the future to account for uncertainties in the values for some of the fitted parameters
419 (e.g., number of exposed individuals at t=0, which is not measured).
420
421 **Model Calibration and Uncertainty Quantification**
422 There are 9 unknown quantities in Eqs. [1]-[4]: the initial conditions $[S_0, E_0, I_{r0}, I_{u0}]$ and the
423 parameters $[\alpha, \beta, \mu, Z, D]$. To characterize these 9 variables, we use daily observations of the
424 number of active cases, denoted here by $\{d_1, d_2, \ldots, d_T\}$, where $T$ is the total number of days
425 with observations. These observations are related to the value of reported infections $I_r(t)$ in the
426 model equations. We employ a Bayesian formulation of this model calibration problem and
427 define a posterior probability distribution over the 9 model variables given the observations.
428 Maximizing the density of this posterior distribution is akin to nonlinear least squares and results
429 in a single point estimate of the most likely model variables. These parameters are called the
430 Maximum aposteriori (MAP) point. The posterior distribution however, also contains information
431 about uncertainty in parameters that can result from observation noise or ill-posedness.[1] The
432 posterior distribution depends on the nonlinear model in Eqs. [1]-[4] and does not therefore fall
433 into a canonical family of distributions that can be sampled directly (e.g., Gaussian). We therefore

---

[1] Ill-posedness in this context means that multiple parameter settings can match the data equally well.

434    employ Markov chain Monte Carlo (MCMC) to generate samples of the posterior.[2] These samples
435    are then propagated through the system of differential equations in Eqs. [1]-[4] to characterize
436    uncertainty in our predictions.

437            Defining the posterior distribution requires us to define two things: (1) a prior probability
438    distribution over the model variables that describes knowledge we may have about the variables
439    before observing the data (e.g., known bounds, positivity, etc...) and (2) a statistical error model
440    for the difference between the model predictions and observations. For the parameters
441    $[\alpha, \beta, \mu, Z, D]$ we adapt the prior distribution used in the Columbia model. For the initial
442    conditions $[S_0, E_0, I_{r0}, I_{u0}]$ we employ a combination of log-normal and uniform distributions to
443    represent our prior knowledge. Notice that even though the prior distribution is constructed from
444    canonical distribution families, the posterior will not be of a standard form because it depends
445    on the model in Eqs. [1]-[4]. For the statistical error model, we make the common assumption
446    that at any time $t_i$, the errors between the model predictions and observations, $e_i = I_r(t_i) - $
447    $d_i$, are normally distributed with some constant variance $\sigma^2$. This variance accounts for both
448    observation noise and the difference between our model and reality. We employ a hierarchical
449    formulation and estimate $\sigma^2$ along with the model variables. An inverse-Gamma hyper prior is
450    used.

451

---

[2] We employ a simple random walk Metropolis MCMC algorithm with a proposal based on the Laplace approxima-
tion of the posterior density at the MAP point.