

1 We thank the reviewers for their constructive feedback. We appreciate the comments that our “careful empirical study”  
 2 [R4] and “sensitivity analyses [...] are of extreme importance” [R2]. Further, they are “critical to understand[ing]  
 3 whether assumptions, data or both are providing evidence” [R3] about the effectiveness of different nonpharmaceutical  
 4 interventions (NPIs) against COVID19 transmission. Given the importance and time-sensitivity of these results, and the  
 5 minor criticisms raised by the reviewers, we hope that our clarifications below will allow the reviewers to increase their  
 6 scores. Additionally, in line with reviewer comments, **we have run a number of additional experiments that will**  
 7 **also be included in the camera-ready version.**

8 [R2] [R3] [R4] **Contribution.** Our work is the first that performs  
 9 *structural* sensitivity analysis and compares the robustness of data-  
 10 driven NPI effectiveness models. Our findings are policy relevant;  
 11 the high sensitivity of the model used in [8], subsequently published  
 12 in Nature, raises concerns (though the authors do not claim to distin-  
 13 guish individual NPI effects). A recent preprint (concurrent work to  
 14 us) also finds [8] has high sensitivity [Soltesz *et al*, *On the sensitivity*  
 15 *of non-pharmaceutical intervention models for SARS-CoV-2 spread*  
 16 *estimation, 2020*]. We highlight that neither [2] nor [8] test structural  
 17 assumptions<sup>1</sup>, and [8] *never reports the sensitivity of NPI effective-*  
 18 *ness estimates in the tests they perform.* [R2] correctly points out  
 19 that these models make “assumptions that we now know are vio-

20 lated”, exactly why our novel mathematical results (§5 *Effectiveness*  
 21 *in Context*) are important steps forward. Our results show that when  
 22 commonly made assumptions are violated, estimates must be inter-  
 23 preted as averages, taken over contexts of the dataset, and expert  
 24 judgement is required to adjust them to local, unique circumstances.

25 [R2] **Implementation.** Our implementation of the model of [8] is  
 26 correct; we only model latent infections as a discrete renewal process  
 27 while deaths are modelled as in [8] and [2] i.e., produced via discrete  
 28 convolutions. **We believe this misunderstanding is due to typo-**  
 29 **graphical errors in the supplement, Eqs. (128), (129).** The correct  
 30 equation, modelling only deaths, is  $N_{t,c}^{(D)} = R_{t,c} \sum_{\tau=1}^t N_{t-\tau,c}^{(D)} \cdot$   
 31  $\pi_{SI}[\tau]$  where  $\pi_{SI}[\tau]$  is the discretised serial interval distribution,  
 32  $N_{t,c}^{(D)}$  is the daily number of infections that result in fatalities.  $R_{t,c}$   
 33 is the instantaneous reproduction number at time  $t$  in country  $c$ . We seed  
 34 this with a latent variable  $N_{0,c}^{(D)}$  that incorporates the country-specific

35 infection fatality rate,  $IFR_c$ . Other than truncation and naming, this is identical to  $c_{t,m} = R_{t,m} \sum_{\tau=0}^{t-1} c_{\tau,m} g_{t-\tau}$  [8]  
 36 where the convolution has been rewritten indexing over the other variable. Since  $c$  represents the *total* number of  
 37 infections, we have  $c_{t,c} = N_{t,c}^{(D)} / IFR_c$ . We compute the expected number of deaths as  $\bar{D}_{t,c} = \sum_{\tau=1}^{63} N_{t-\tau,c}^{(D)} \pi_D[\tau]$ ,  
 38 where  $\pi_D[\tau]$  is the discretised infection-to-death delay. We implemented all models ourselves to minimise discrepancies  
 39 between models and make fair comparisons.

40 [R3] [R4] **Confounding.** Thank you for pointing out the no-confounder assumption. We agree that this assumption is  
 41 critical, and will update the tone of the conclusion to reflect the assumptions we tested. For clarity, the NPI leave-out test  
 42 assesses how much the effect of unobserved interventions are attributed to observed NPIs [R2] **thereby testing this**  
 43 **assumption.** We apologise for not clarifying the purpose of this test. In light of your feedback, we have run additional  
 44 experiments finding low sensitivity when previously unobserved NPIs from the OxCGRT NPI dataset [Thomas Hale  
 45 *et al*. *Oxford COVID-19 Government Response Tracker. (2020)*] are observed (Fig. 1, bottom). These tests are  
 46 imperfect but considered best practice [Rosenbaum *et al.*, *Assessing Sensitivity to an Unobserved Binary Covariate in*  
 47 *an Observational Study with Binary Outcome, 1983*]. We highlight that our results show that confounding is the key  
 48 limitation of such NPI effectiveness models. For example, if we had found that effectiveness estimates fluctuate widely  
 49 under different epidemiological parameters, we would not have been able to make strong conclusions regardless of  
 50 whether we observe all relevant factors.

51 [R2] **Effectiveness Prior.** Thank you for your comment. We take our effectiveness prior from [2], and it reflects the  
 52 belief that NPIs have moderate effects. Low posterior correlation  $r < 0.4$  between NPI effectiveness estimates and low  
 53 sensitivity suggests that collinearity is manageable. Furthermore, we have added run a test using the suggested prior  
 54 from [8] (Fig. 1, top).

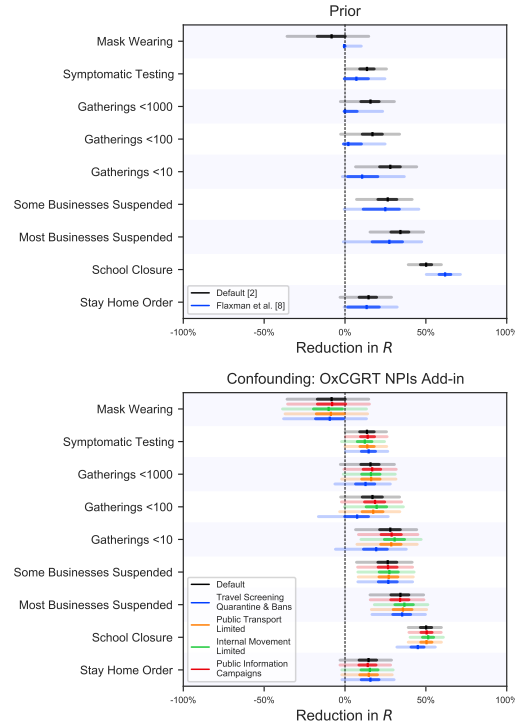


Figure 1: Additional experiments using the baseline model. Top: [R2] prior from [8]. Bottom: [R3] [R4] additional confounding tests. The NPIs from the OxCGRT dataset (as labeled) are now observed.

<sup>1</sup>The most recent version of [2] reproduces our structural sensitivity analysis from the preprint corresponding to this submission.