Dear Editor:

Thank you for the chance to revise and resubmit our manuscript. Below please find our responses to reviewers.

# Reviewer 1

Authors aim to evaluate how using the wrong model in terms of the distributions of the infectious and latent periods affect estimates of $\mathcal{R}_0$.

Authors did some work in this direction, but I believe more work is needed to address their question. In particular, authors will need to use a model that incorporates realistic distributions of latent and infectious periods. They can accomplish this using the linear chain tricks or alternative integro-differential equation models. Then they can generate simulated data from that "true/realistic" model and go ahead and estimate $\mathcal{R}_0$ using simple models as they do here.

Having said this, I would like to point out that I am sure there are a number of papers out there that have look at similar questions in the context of different infectious diseases, etc. It'd be good to make sure the relevant literature is considered and appropriately cited.

We agree that this is an important point. Although our study considers the effect of making incorrect assumptions about the reporting time distributions, these assumptions implicitly depend on latent and infectious period distributions. We have added the following paragraph in the Discussion to address this point:

"Delay distributions in reporting time implicitly depend on latent and infectious period distributions. For example, when cases are assumed to be reported at the end of infectious periods, the delay distribution is equivalent to the convolution of the latent and infectious period distributions. These distributions play important roles in shaping the epidemiological dynamics, including the probability and size of an outbreak, the stability and persistence of a system, and the prediction of dynamical transitions. Assumptions about latent and infectious period distributions can also have large effects on

the estimates of the basic reproductive number."

# Reviewer 2

This manuscript demonstrates that the timing of when disease cases are reported can influence the performance of statistical models fit to time series of incidence data. Essentially, the results show that generating data with one reporting scenario and fitting a model that assumes a different reporting scenario leads to sub-optimal inference on the basic reproductive number, $\mathcal{R}_0$. In the case when the infectious period is estimated, the degraded inference is not too bad (bias on the order of 5 - 10%, actual coverage rates of nominally 95% intervals is off by a few percent). Coverage rates plummet when the infectious period is assumed fixed, but the bias is still not awful, so it is not clear just how much of a worry the loss of coverage will be for real data, given all the other strong assumptions of the simulation study.

Overall, it seems to me that this is a publishable piece of work. It may make epidemic modelers think twice about the consequences of their particular assumptions about the process generating incidence data, and such an outcome would propel the field forward. Technically, the study is soundly executed. Given the modest bias of the estimates from the mismatched models, I think it's possible that readers may take the reverse message from this study than the one the authors intend. Namely, readers might interpret these results to indicate that faulty assumptions about the process generating the incidence data are tolerable, especially in light of all the other complications that attend real data, and that have been assumed away here. But, such a demonstration would be a worthy contribution to the literature as well.

I have one comment that the authors may wish to consider, although this doesn't rise to the level of a critical flaw that must be corrected. It seemed to me that the Discussion lost sight of the fact that statistical models often have to make simplifying assumptions. Despite these simulation results, it seems to me that making a simplifying assumption about the process generating the incidence data may not be a terrible thing to do with real data. When it comes to fitting statistical models to noisy data, models that are "more realistic", in the sense of incorporating more of the underlying mechanism, do not necessarily lead to

improved inference. One has to choose their battles with respect to what one tries to capture in a model, and what one simplifies away. With real data, it is entirely unclear whether model fit would be improved by fitting a model with a more complicated description of the process generating the incidence data, and it may very well be that such a question can only be answered on a case-by-case basis. In the absence of such understanding, it seems to me that the Discussion could strike a somewhat more measured tone regarding the merits of simpler vs. more complicated models.

We agree. We did not intend to say that modelers should always thrive to build the more complex models. We have revised the final paragraph of the Discussion to better reflect our intention as well the reviewer's comment:

"All statistical models must rely on simplifying assumptions. Given a range of choices of assumptions one can make, modelers have to decide which aspect of a system they wish to capture at the cost of simplifying other aspects. However, some assumptions are often more carefully tested than others. In the case of epidemic models, assumptions about the process of transmission tend to be the most important ones; processes through which cases are reported are often overlooked and simplified. Our study shows that a seemingly negligible change in the assumptions of an epidemic model can affect the inference of infectious disease transmission. Although the amount of bias introduced from the misspecification of the observation model can be small (e.g., less than 5% in our examples), modelers trying to make serious predictions may want to control for potential biases. We caution disease modelers to be more careful about modeling decisions as they may have unexpected consequences on the inference."

As trivia, perhaps I missed something, but the distinction between the boxplots within each pair in the left panels of Fig 2 was not immediately clear to me. I guessed at the distinction by comparison with the parallel structure of the panels on the right of the figure.

We have modified the figure to make it clearer.

I congratulate the authors on a tidy piece of work, and a well written manuscript.

Thank you very much!