

RESPONSE TO REVIEWERS

Editor

Editor: The comments and recommendations from two expert reviewers are now available for your manuscript. These reviewers judged the reported discoveries to be of medium significance, and the potential scientific impact of your work to be medium/high. They found that the manuscript needs improvement in text and additional data analysis. Editors generally agree with their concerns and recommendations, which led to a designation of medium/high priority.

***Response:** Thank you very much for handling our manuscript and providing these two constructive reviews that have helped us to improve the paper.*

Editor: Note that many manuscripts receiving medium priority based on reviewer comments are not accepted by the Board of Editors. If appropriate, the board may invite a resubmission following the rejection, which is intended to enable you to improve the manuscript towards receiving a high or top priority. The authors should pay close attention to the detailed review comments and address each comment with significant improvements.

***Response:** We have used the revision to include additional analysis and improve the presentation. We hope that you and the reviewers now feel that our manuscript makes a major contribution towards understanding what aspects of influenza virus evolution are predictable, and which ones are not. Given the major public health importance of influenza and the fact that influenza virus evolution is one of the very few systems where adaptive evolution can be observed in real time, we feel that our results are of very high relevance to the community.*

Reviewer: 1

R1: The manuscript by Barrat-Charlaix et al. discusses the problem of predictability of mutations in seasonal influenza viruses. It presents data analysis of frequency trajectories in H1N1 and H3N2 lineages and simulations of different selection scenarios. It is an interesting revisit to predictions in influenza, however I have some major concerns about the analysis and formulation of the conclusions that should be addressed:

***Response:** We address the individual points by the reviewer below.*

R1: 1. The prediction problem is defined for each mutation: based on the tracked frequency trajectory, can the future (fixation, loss, or polymorphism) of the mutation be predicted. Such formulation has been previously proposed by Illingworth and Mustonen, (eg. Genetics 2011, Plos Pathogens 2012). By averaging over all amino-acid substitutions, the authors show that the fates of mutations are not determined by the value of the starting frequency.

***Response:** We thank the reviewer for pointing out the work by Illingworth and Mustonen, which is indeed very relevant and complementary to our approach. We discuss their relation below.*

R1: 1.1. Mutations in influenza are highly nested, with a substantial hitch-hiking, and no attempt is made to disentangle such dependencies when counting the mutations.

***Response:** The reviewer is correct that our original manuscript did not attempt to account for nested mutations and we agree that this should have been done.*

We added a new section in the Supplementary Material (“Correcting for nested trajectories”) where investigate the extent to which nesting affects our results. We cluster trajectories of mutations that partly appear on the same strains to various degrees, and trajectories corresponding to mutations always or often appearing on the same strains are then counted as one “effective” trajectory. Our analysis shows that nesting of mutations does not change our results significantly. We therefore left the figures of main text unchanged and refer to the supplement for a discussion of this issue.

R1: 1.2 Despite the more general formulation in the beginning, this approach makes use only of the last time

point in the trajectory, rather than the full trace (contrary to the work of Illingworth& Mustonen). I find the observed neutral-like statistics not surprising for such a limited data input, which doesn't capture past frequency dynamics. Therefore, these conclusions should be revisited and benchmarked against the more general implementation of Illingworth&Mustonen-like approach.

Response: Illingworth and Mustonen (I&M) presented a careful analysis to what extent trajectories can be fit using a model assuming a single fitness coefficient, or a model that differentiates different background fitnesses. Our results are consistent with their finding that a single parameter model (which only allows monotonic trajectories) does not capture the majority of trajectories. Beyond that, our approach and aim differ from I&M. We investigate a priori features of trajectories that potentially have information on whether the mutation fixes, not whether the trajectories can be parametrized by a model. Fitting a model to a partial trajectory to predict fitness is not expected to work: I&M showed that subsequent interference has a substantial effect and our results indicate that the speed at which the mutation rises does not predict fitness.

Furthermore, I&M estimate fixing mutations to have a much higher intrinsic fitness than dying mutations while we show that very little such information is available a priori. The retrospective fit to trajectories tends to assign a positive fitness effect to fixing (and hence on average rising) mutations. This is why we resort to a model free approach to ask what a priori information can help to determine which mutations fix or die.

ResponseDraft: *Should we mention all the changes we made in the text? We spend more time explaining that even though there is linkage, and even though we're looking at just the last point, we expect something above diagonal. (Last paragraph before "Results", last paragraph before "Predicting future frequencies"), and first paragraph of "Simulations of models of adaptation"*

R1: 2. The authors examined the predictive power of one predictive method, the LBI. However, they did not do a systematic comparison of the different methods that they cite (Morris et al, 2018), which differ in prediction targets and methods. Therefore, the general conclusions about predictability, eg. on page 8, in line 47, and on page 9, line 48 are too sweeping and should be made precisely for those methods looked at in detail.

Response: *The reviewer is right to point out that our original manuscript did not perform a systematic comparison of different prediction methods.*

In a recent model-driven effort (Huddleston et al), we investigate the power of sequence based, antigenically informed, experimentally informed, or phylogenetic (LBI) scores to predict future influenza virus populations. The LBI was among the highest performing models and we hence used this score in the original manuscript. But we took the suggestion of the referee (and referee 2) to heart and now include all highly performing models from Huddleston et al. in a new section in the Supplementary Material ("Ability of fitness models to predict fixation"). These additional results show that only mutational load has some power to differentiate between fixing and non-fixing trajectories. We believe this addition has made the manuscript stronger and we thank both referees for the suggestion.

We also agree with the reviewer that these different methods differ in prediction targets: our results indicate that most of them are not predictive of frequency dynamics, but they remain predictive of the composition of the future population. In consequence, we have toned down our conclusions about predictability. Pierre: I say that because I changed p8 l47. I have not touched p9 l48 because the new appendix shows it is a correct statement.

R1: 3. Going through the previous and cited literature, I think the authors should cite some of these works in a more careful way. Specifically, the very related work by Illingworth and Mustonen is not mentioned. The distribution $P_{\Delta\text{elta}_t}(f|f_0)$ has been at the core of the method of Strelkowa and Lassig, 2012 (which is cited, but at other parts of the text), and the same paper also uses a similar simulation model, which should be acknowledged at the appropriate points of the text.

ResponseDraft: *I cited (Strelkowa and Lassig) in the first paragraph of the section "Predicting future frequencies". Did not cite it close to the simulation part, will do. We cite Illingworth and Mustonen plenty of times now.*

R1: It would help if the plots with red-green-blue lines were also distinguished by differing markers.

Response: *We have now added markers in most figures.*

Reviewer: 2

R2: The paper does a retrospective study of amino acid substitutions in seasonal Influenza to determine what properties of these substitutions could help predict their fate in the future. The authors find that future frequency trajectories are surprisingly unpredictable. Even predicting which mutations fix in the population is hard. The authors find that the current frequency of a mutation is the best predictor for the probability of fixation, which would be expected under neutrality but not in a model with selection. I appreciate this study, I think it will be of interest to many readers and it is generally well done and fairly easy to follow.

Response: *We thank the reviewer for this assessment and the constructive suggestions below.*

R2: 1. The authors focus on one feature at a time (frequency, epitope status, LBI). It is interesting to see that each of these is not very predictive of future frequency / prob of fixation. However, I think the obvious next step is to see whether a combination of many features could do a better job of predicting. I am not sure why the authors don't try to fit a model that takes into account all information they have about sites (say, type of AA change, location in the gene, current frequency etc) and see if a ML model is able to make predictions.

Response: *A similar point was also brought up by reviewer 1. In addition to the LBI, epitopes, age, and geographic distribution we now also include the fitness scores of the best performing models in a recent machine learning effort to predict the composition of future populations (as opposed to predicting fixation). These models include antigenic, experimental, and sequence or phylogenetic information. Interestingly, all these models do similarly poorly with the exception of mutational load. These results are now included as supplementary information and discussed.*

R2: 2. Fig 2A and 2B look quite different to me. In the text it appears to me as if these are very similar to the authors. What are the characteristics of the AAs that fix in H1N1?

Response: *Whether 2A or 2B are perceived as similar or not probably depends on how big a deviation from the diagonal one expected a priori. We now discuss and highlight the differences in the text and discuss potential reasons for the different behaviour.*

R2: 3. Question: what do these results mean for vaccines? How to decide which strains to use for vaccines? This may be known to those who work on Influenza, but for a relative outsider it is not clear.

Response: *Computational prediction of which circulating strains will dominate future influenza populations have been attempted for about 2 decades and have recently become a regular input to the biannual consultations on the influenza vaccine composition. Since performance of predictive models can only be assessed on limited historical data, we believe that our model-free approach clarifies which aspect of influenza evolution are predictable and which ones are not.*

ResponseDraft: *We should maybe make clear that this is not directly related to predictions for vaccines (at least it's not the aim). That's John's paper. The main implications are about our understanding of the evolutionary dynamics of influenza. This could also help us to reply to the first comment: doing ML is more useful to predict the future population (for vaccines) than for predicting fixation of a single mutation.*

R2: 4. Question: at what point is it expected that the vaccine itself will influence frequency trajectories? See Wen et al Biorxiv 2020

Response: *The degree to which the vaccine drives evolution is a very interesting question! In the past, global influenza vaccine uptake was unlikely sufficient to generate much pressure and fairly thorough mixing of variants across countries with very different vaccination rates have been used to argue against a strong effect. But this might have changed in recent years and the paper by Wen et al in Viruses, 2018 picks up some signal. Overall, the case for strong effect is still weak to our knowledge.*

R2: Fig 1B is very hard to read. Maybe it should get more space (fig 1C could work with less space).

Response: *We have reworked this figure to allow more space for panel B.*