

Evaluating an epidemiologically motivated surrogate model of a multi-model ensemble

Sam Abbott^{1,2}, Katharine Sherratt^{1,2}, Nikos Bosse^{1,2}, Hugo Gruson^{1,2}, Johannes Bracher³, and Sebastian Funk^{1,2}

¹The Centre for Mathematical Modelling of Infectious Diseases, London School of Hygiene & Tropical Medicine, London, UK

²Department of Infectious Disease Epidemiology, London School of Hygiene & Tropical Medicine, London, UK

³Chair of Statistical Methods and Econometrics, Karlsruhe Institute of Technology (KIT), Karlsruhe, Germany

Corresponding author:

Sam Abbott^{1,2}

Email address: sam.abbott@lshtm.ac.uk

ABSTRACT

Multi-model and multi-team ensemble forecasts have become widely used to generate reliable short-term predictions of infectious disease spread. Notably, various public health agencies have used them to leverage academic disease modelling during the COVID-19 pandemic. However, ensemble forecasts are difficult to interpret and require extensive effort from numerous participating groups as well as a coordination team. In other fields, resource requirements have been reduced by training simplified models that reproduce some of the observed behaviour of more complex models. Here we used observations of the behaviour of the European COVID-19 Forecast Hub ensemble combined with our own forecasting experience to identify a set of properties present in current ensemble forecasts. We then developed a parsimonious surrogate forecast model intending to mirror these properties. We assess forecasts generated from this model in real time over six months (the 15th of January 2022 to the 19th of July 2022) and for multiple European countries. We focused on forecasts of cases one to four weeks ahead and compared them to those by the European forecast hub ensemble. We find that the surrogate model behaves qualitatively similarly to the ensemble in many instances, though with increased uncertainty and poorer performance around periods of peak incidence (as measured by the Weighted Interval Score). However, the proposed model appears better probabilistically calibrated than the ensemble. We conclude that our surrogate forecast model may have captured some of the dynamics of the hub ensemble, but more work is needed to understand the implicit epidemiological model that it represents.

INTRODUCTION

Multi-model and multi-team ensembles have become increasingly popular as an approach to increase the robustness and performance of infectious disease forecasts over the last decade (Reich et al. 2022). The experience of other domains has inspired these approaches, for example, climate modelling (IPCC, n.d.), where ensembles of both multiple models and from multiple teams have a long history of providing forecasts that stakeholders trust. The trend towards large-scale multi-team ensemble forecasting in infectious diseases has accelerated during the COVID-19 pandemic due to a pressing need for reliable forecasts and a perception that many publicly available forecasts were low quality. Over 2020 and 2021, teams established COVID-19 Forecasting Hubs covering the US (Cramer et al. 2022), Germany and Poland (J. Bracher et al. 2021), and Europe (Sherratt et al. 2022) (all three including authors of this study). All of these collaborations ensembled contributions from multiple independent teams using a similar approach and have shown that their ensemble forecasts outperform most individually contributed forecasts whilst remaining generally robust to outliers in reporting. Both the US and European Forecast Hubs were supported and received funding from public health agencies (the Center for Disease Control, CDC, and European Center for Disease Prevention and Control, ECDC, respectively) with their forecasts

47 used in official communications by these agencies.

48 Whilst there is robust and consistent evidence that multi-team ensemble forecasts provide reliable
49 and performant forecasts across domains (Reich et al. 2022) they also have a range of downsides. The
50 most significant is the difficulty in interpreting them. This relates both to the underlying mechanisms for
51 the forecasts they produce and to understanding if and when their behaviour is desirable. This impacts
52 users' trust, how easily ensemble performance can be improved, and how easily contributor forecasts can
53 be improved. Forecasts from these ensembles also require considerable resource cost to produce as they
54 typically require contributions from multiple independent teams, the development of several models, and
55 a centralised group to run the ensembling project. Additional challenges with maintaining multi-team
56 collaborations can include providing detailed feedback to those contributing forecasts that would allow
57 them to improve their forecast approaches, providing incentives for forecasters to continue to contribute
58 and adjust their models to changing conditions, and difficulty improving the quality of the ensemble by
59 learning from past predictive performance (Sherratt et al. 2022). Each of these issues may impact the
60 long-term quality of the resulting forecasts and have implications for end-users. Little progress has so far
61 been made in mitigating these downsides or in improving access to the high-quality and robust forecasts
62 they seek to generate for geographies without coverage or for other infectious diseases. There has also
63 been limited critical feedback on the structure of forecasting ensembling projects for infectious disease
64 epidemiology and little evaluation of the effort required to produce them relative to their benefits for
65 improving forecast performance.

66 In climate forecasting (Castelletti et al. 2012; Edwards et al. 2021; Williamson et al. 2013), as
67 well as in other fields such as astrophysics (Vernon, Goldstein, and Bower 2014), emulation approaches
68 have been used to circumvent resource requirement issues for complex models by training a simplified
69 model, usually, a non-parametric statistical model, to replicate the behaviour of either the entire model or
70 sub-components. These approaches generally take the same inputs as the models they seek to emulate
71 and then are trained based on the output from those models. In the context of epidemiological models,
72 non-parametric emulation has been used to allow the rapid exploration of the parameter space of complex
73 models that would otherwise be resource-prohibitive (Iskauskas et al. 2022; Charles et al. 2022). These
74 methods may be less useful for resolving some of the issues of multi-team and multi-model forecasts as
75 they do not provide interpretability, key for stakeholder take-up. Additionally, it is not clear how these
76 methods perform out of sample, or how they would be applied to a quantile-based forecast.

77 In this work, we draw insights from ensemble forecasts produced and endorsed by the COVID-19
78 Forecast Hubs, as well as our forecasting work, to propose and evaluate a "surrogate" forecast model. This
79 surrogate model seeks to reproduce ensemble performance by mimicking its behaviour based on a minimal
80 set of easily communicated and epidemiologically justifiable assumptions, and limited computational
81 resources with an easily generalised implementation. The primary aim of this approach is to help highlight
82 the behaviour, and potential mechanisms behind this behaviour, of ensemble forecasts widely considered
83 the gold standard for COVID-19 forecasting. Our secondary aim is to provide the basis for a robust
84 forecasting system that others can easily reuse both in operational contexts and as a platform for future
85 research.

86 To achieve these aims, we evaluate an initial attempt at developing a surrogate model to replicate the
87 observed behaviour of current multi-team forecast ensembles based on a set of clear assumptions. We
88 submitted this model to the European Forecast Hub and here we evaluate its performance relative to the
89 Hub ensemble. In this work, we first define the model and summarise its implementation, with a focus on
90 minimal resource use and reproducibility as a GitHub Actions workflow ("About GitHub-hosted Runners"
91 2022).

92 We then evaluate the surrogate model's real-time performance in comparison to the European Forecast
93 Hub ensemble by visualising forecasts using the weighted interval score (Johannes Bracher et al. 2021), a
94 commonly used proper scoring rule, and quantifying the empirical coverage of the forecasts produced.
95 We highlight settings where this model performs well as a surrogate for the ensemble forecast and areas
96 where it performs less well. Finally, we summarise our findings, discuss their implications, and highlight
97 areas where more work is needed. We aim for this work to highlight some of the potential implicit
98 assumptions of current COVID-19 Forecast Hub ensembles, provide a sensible, low-resource, surrogate
99 model where large-scale collaborative forecasting efforts are not possible, and provide inspiration for
100 forecasters looking to make principled improvements to their models.

MATERIALS AND METHODS

Setting of the European COVID-19 Forecast Hub

```
#> Warning in !is.null(rmarkdown::metadata$output) && rmarkdown::metadata$output
#> %in% : 'length(x) = 2 > 1' in coercion to 'logical(1)'
```

To understand the behaviour of the Forecast Hub ensembles we need to first explore the structure of the COVID-19 Forecast Hubs (Cramer et al. 2022; J. Bracher et al. 2021; Sherratt et al. 2022). These collaborations share a similar design with a central team running the hub, vetting forecasts, and producing the ensemble forecast as well as teams of independent forecast contributors who design their forecast models and then use them to produce a weekly forecast that they then submit to the central hub team. Each hub targets a range of metrics, including test-positive reported cases, reported deaths, and hospitalisations; has a specific geographic focus, and asks for weekly forecasts (using MMWR epidemiological weeks i.e. Sunday to Saturday (Department of Health, n.d.)) over a time horizon of a few weeks. Observed data are available and updated daily, and most submitted forecasts use this dataset, along with potentially other sources of real-time information, to produce forecasts. Here, we focus on reported cases and primarily on the European Forecast Hub but our observations hold, in our view, across COVID-19 Forecast Hubs and to a lesser degree targets. We focus on reported cases as these represent the most common forecast target for COVID-19 forecast models (Nixon et al. 2022), they are often of the most direct interest due to being a leading indicator for other metrics such as hospitalisations (Meakin et al. 2022), and they are generally the most challenging to predict (Sherratt et al. 2022). In general, 5 main classes of forecast models are submitted (Johannes Bracher et al. 2022; Cramer et al. 2022), statistical forecasting models such as ARIMA models, mechanistic forecasting models based on the compartmental modelling framework and its generalisations (Srivastava, Xu, and Prasanna 2020; Li et al. 2021), semi-mechanistic approaches that blend both of these approaches (Castro et al. 2021; Nikos I. Bosse et al. 2022a), agent-based simulation models (Rakowski et al. 2010; Adamik et al. 2020), and human insight based forecast models that may also include elements of other methods (Karlen 2020; Nikos I. Bosse et al. 2022a). Real-time evaluation has shown that each of these classes of models may perform well depending on the context and specific implementation of the forecast model (Nikos I. Bosse et al. 2022a).

We extracted forecasts and data on notified weekly COVID-19 cases from the European forecasting hub (Sherratt et al. 2022; E. C. F. H. Team 2021) from the 15th of January 2022 to the 19th of July 2022 for the ensemble model (referred to as the `EuroCOVIDhub-ensemble` by the hub team) and the surrogate model (submitted as `epiforecasts-weeklygrowth` and defined in the following section). We included all locations covered by the European forecasting hub which were 32 European countries, including all countries of the European Union and European Free Trade Area, and the United Kingdom. Data on notified weekly cases was originally sourced from the Johns Hopkins University (JHU) curated data repository (Dong, Du, and Gardner 2020). We used the latest available observed data as of the 1st of September 2022 (commit `f6922c3e4bdcb055abcbbba8e73472afacac4cf40` from (E. C. F. H. Team 2022)). Incidence was aggregated by epidemiological week (i.e. Sunday through Saturday). As observations are subject to revisions this means that the data used to produce forecasts for a given date may not reflect the data used for evaluation. To account for this we followed the practice of the European forecasting hub project in excluding forecasts made using anomalous truth data in the week of the forecasts production and excluding forecasts for target weeks with anomalous data (Sherratt et al. 2022). We defined anomalous data based on the implementation used by (E. C. F. H. Team 2021) where a data point is considered anomalous if a future revision alters it by more than 5%.

The European Forecast Hub requests forecasts for one to four-week forecast horizon and requires forecasts to use a pre-specified format with 23 quantiles of the predictive probability distribution. No restrictions were placed on who could submit forecasts and the hub team actively invited participation from research groups known to be involved with COVID-19 forecasting projects. Teams submitted forecasts at the latest two days after the complete dataset for the forecast week became available and were allowed to use all data available at the time of submission (i.e. including up to two days of data for the current week). The ensemble forecast was constructed by taking the median of all forecasts for each predictive quantile without the exclusion of any validly submitted forecast (where validity was defined as passing minimal formatting checks by the hub team and timely submission) (Sherratt et al. 2022). An ensemble was only produced for locations with at least 3 independent forecast models including the hub baseline model. Submitted forecasts and target observations are available from the European Forecast Hub

Table 1. Observations on the relative performance of the Forecast Hub ensemble compared to our forecast submissions.

No.	Observation
1	Robust to daily reporting artefacts
2	Some ability to forecast future trend changes
3	Less reactive to apparent observed changes in trend
4	Sharper forecasts
5	A tendency towards underprediction
6	Modelling the reporting process appears to have little impact

155 GitHub repository (E. C. F. H. Team 2022). We provide code in the repository of this study to streamline
156 access.

157 **Observations based on previous forecasts**

158 We have contributed a range of forecasts to COVID-19 forecast collaborations, generally focused on
159 semi-mechanistic statistical methods and human insight-based forecasts. Our forecast submissions
160 have not systematically over- or under-performed other forecasts submitted to the forecasting hubs (see
161 *epiforecast* tagged models at (E. C. F. H. Team 2021) and (Nikos I. Bosse et al. 2022a; Cramer et al.
162 2022; J. Bracher et al. 2021; Sherratt et al. 2022)). The model-based forecasts we have contributed have
163 focussed on trying to carefully model the underlying infectious disease dynamics from infection through
164 to symptom onset, and test positivity using non-exponential delay distributions whilst also attempting
165 to model the complexity of daily, within the week, reporting periodicity (Nikos I. Bosse et al. 2022a;
166 Abbott, Hellewell, Thompson, et al. 2020; Abbott, Hellewell, Sherratt, et al. 2020). Based on our
167 observations our forecasts have generally captured the current trend relatively well but have not been
168 robust to reporting issues such as large outliers in reporting and changes to reporting patterns. Our
169 previous methodology also requires significant computational resources, running for an hour on a Azure D
170 v5-series 16-core machine, when producing forecasts for the European forecasting hub (“Pricing - Linux
171 Virtual Machines” 2022). This resource usage is likely beyond the capacity of many interested in having
172 access to state-of-the-art short-term forecasts of infectious diseases. In our model-based forecasts, we
173 did not attempt to capture potential future interventions or known interventions not currently observed
174 in the epidemiological data whereas in our human insight models these were implicitly included. We
175 found that our human insight-based forecasts outperformed our model-based forecasts on average. This
176 was particularly the case when forecasting cases and at longer forecast horizons. We hypothesised that
177 this may have been driven by including additional information not observed in the epidemiological data
178 (Nikos I. Bosse et al. 2022a).

179 Unlike our epidemiologically motivated forecast submissions, the hub ensemble forecasts were
180 typically robust to daily reporting artefacts. They also demonstrated some ability to forecast future
181 changes in trends that were not present in the observed data similarly to our human insight forecasts
182 indicating the likely inclusion of either human insight, or assumptions about future interventions. In
183 comparison to our submitted forecasts, the ensemble forecasts were less reactive to changes in trend
184 such as from stable or reducing case incidence to increasing incidence. On the other hand, this also
185 meant that the ensemble was less likely to adopt short-term changes in incidence and hence produced
186 better long-term forecasts. Finally, the ensemble forecast tended to produce sharper forecasts and have a
187 tendency toward under- vs overpredicting. Our observations are summarised in Table 1.

188 **Model**

189 ***Assumptions and simplifications***

190 Based on our observations of forecast performance (summarised in Table 1), here we define a model
191 with similar, but simplified, epidemiological characteristics to our previous approaches to model-based
192 forecasting (Nikos I. Bosse et al. 2022a) to produce an ensemble-like performance without sacrificing
193 interpretability and with a lower cost to produce. The first simplification we make is to model only weekly
194 data, rather than using daily data and then aggregating. This mitigates the impact of daily reporting
195 artefacts. It also serves to increase the auto-correlation of the forecasting model as there is an increased

Table 2. Assumptions/simplifications based on observations of the relative performance of Forecast Hub ensembles compared to our forecast submissions.

Assumption	Observation
Reported cases can be modelled using weekly data and a generative process discretised by week	1, and 2
Reported cases can be modelled as if they represented infections	6
The growth rate of infections can be represented as an auto-regressive process with an order of 1 week	3 and 4
Unobserved interventions and more general changes in transmission towards a stable state can be represented using a multiplicative decay parameter	2, and 5

196 lag before changes in daily observations gain significant weight in the model. This leads to the observed
 197 ensemble behaviour of being relatively auto-correlated and resistant to short-term changes in trend.

198 The second simplification we make is to ignore the underlying latent infection process and focus only
 199 on the observed reported cases. This removes the need for, potentially misspecified, external information
 200 on the delay from infection to report, and reduces computational requirements due to a reduction in model
 201 complexity. However, this sacrifices some of the interpretability of the forecast model as any transmission
 202 statistics we now calculate will be based on reported cases and not latent infections. As discussed in
 203 (Gostic et al. 2020) this leads to varying amounts of bias depending on the epidemic phase.

204 The final simplification is to model the growth rate as a differenced auto-regressive process with an
 205 order 1 rather than using a gaussian process-based method as we have done in other work (Nikos I. Bosse
 206 et al. 2022a; Abbott, Hellewell, Thompson, et al. 2020; Abbott, Hellewell, Sherratt, et al. 2020). This
 207 represents a parsimonious approach in that we encode our expectation that the growth rate should vary
 208 over time and allow this to influence the forecast but we include only a single lag term, reducing the
 209 computational overhead of the model. To model potential unobserved interventions and more general
 210 changes in transmission, we include an additional growth rate modifier restricted to be between 0 and 1
 211 that differs depending on if the growth rate is positive or negative (due to potential differing responses
 212 when cases are growing or increasing) and that acts in a multiplicative fashion (meaning that larger
 213 absolute growth rates are reduced to zero growth more rapidly). This reflects a simplified interpretation
 214 of how the ensemble appears to react to potential future changes by assuming a gradual return to stable
 215 incidence.

216 The only observation for which we do not make an adaptation is the apparent sharpness of the
 217 ensemble compared to our prior forecasting models. Instead, we make use of a negative binomial
 218 observation model allowing the inclusion of overdispersion. This choice is motivated by our belief that
 219 the underlying transmission process is an exponential discrete one and therefore a count error model with
 220 a log link function, where variance is linked to the mean, is a sensible choice. We suggest that part of the
 221 reason the hub ensembles exhibit such sharpness is due to the penalisation of overprediction compared
 222 to underprediction caused by the use of a generalised form of absolute error for the majority of forecast
 223 evaluations (Johannes Bracher et al. 2021). Our set of assumptions and simplifications are summarised in
 224 Table 2.

225 **Definition**

226 We model the expectation (λ_t) of reported cases (C_t) given past reported cases as an order 1 autoregressive
 227 (AR(1)) process by epidemiological week (t) on the log scale. The model is initialised by assuming that
 228 the initially reported cases are representative with a small amount of error (2.5%). We assume a negative
 229 binomial observation model with overdispersion ϕ for reported cases (C_t).

$$\lambda_0 \sim \text{LogNormal}(\log C_0, 0.025 \times \log C_0)$$

$$\lambda_t = C_{t-1} e^{r_t}, t > 0$$

$$C_t | \lambda_t \sim \text{NB}(\lambda_t, \phi)$$

where the mean and variance of the negative binomial are given by

$$\mathbb{E}[C_t | \lambda_t] = \lambda_t \quad \text{and} \quad \text{Var}[C_t | \lambda_t] = \lambda_t + \frac{\lambda_t^2}{\phi}.$$

230 Here r_t can be interpreted as the weekly growth rate. r_t is then modelled as a piecewise constant differenced
 231 AR(1) process modified such that the dependence of r_{t-1} is multiplied by a decay factor ($\xi_{+,-}$) that varies
 232 dynamically according to the sign of r_{t-1} . This assumes that the growth rate is non-stationary with a trend
 233 that is independent of the current growth rate (the differenced AR(1) process), the additional decay factor
 234 encodes the belief that larger absolute growth rates will tend more quickly towards no growth and that this
 235 process may work differently for positive or negative growth rates. This process can be defined as follows,

$$\begin{aligned} r_0 &\sim \text{Normal}(0, 0.25) \\ r_t &= (\mathbf{1}_{r_{t-1} > 0} \xi_+ + \mathbf{1}_{r_{t-1} \leq 0} \xi_-) r_{t-1} + \varepsilon_t \\ \varepsilon_t &= \mathbf{1}_{t > 0} \beta \varepsilon_{t-1} + \eta_t \end{aligned}$$

236 where ε_t and η_t are error terms. The following priors are used,

$$\begin{aligned} \xi_+ &\sim \text{Beta}(3, 1) \\ \xi_- &\sim \text{Beta}(3, 1) \\ \beta &\sim \text{Normal}(0, 0.25) \\ \eta_t &\sim \text{Normal}(0, \sigma) \\ \sigma &\sim \text{Half-Normal}(0, 0.2) \\ \frac{1}{\sqrt{\phi}} &\sim \text{Half-Normal}(0, 1) \end{aligned}$$

237 Where σ , and $\frac{1}{\sqrt{\phi}}$ are truncated to be greater than 0 and β is truncated to be between -1 and 1. The
 238 Beta priors for $\xi_{+,-}$ have been chosen to be weakly informative that the reduction towards 0 growth
 239 is relatively slow. Similarly the prior for β has been chosen to be weakly informative that there is
 240 weak auto-correlation in differenced growth rates. σ has also been made weakly informative under the
 241 assumption that the potential change in growth rates in a single time-step should be relatively small.

242 Forecast evaluation

243 We standardised the magnitude of observations and forecasts across forecast locations, in order to facilitate
 244 comparison, by scaling both weekly notified test positive cases and forecast test positive cases by the
 245 population in the forecast region to an incidence rate per 10,000 people. This differs from the approach
 246 typically taken by the Forecast Hubs where no population standardisation is used (Cramer et al. 2022;
 247 J. Bracher et al. 2021; Sherratt et al. 2022). We then visually evaluated forecasts from a subset of
 248 locations by forecast horizon (1 and 4 weeks) on both the natural and log scales. The countries in this
 249 subset were Germany, Greece, Italy, Poland, Slovakia, and the United Kingdom. These countries were
 250 selected to include forecasts based on different numbers and types of submitted forecast models, to be at
 251 least partially representative of the full sample of forecast locations, and to include nations for which the
 252 authors had a good understanding of local data and transmission dynamics in the study period.

253 We evaluate forecasts for all locations and horizons quantitatively using the absolute error (AE) of
 254 the median forecast and the weighted interval score (WIS) (Johannes Bracher et al. 2021). The WIS is a
 255 quantile-based proper scoring rule that approximates the continuous ranked probability score (CRPS).
 256 Both the WIS and CRPS are generalisations of the absolute error to evaluate probabilistic forecasts and
 257 are widely used to evaluate COVID-19 forecasts, including by the European Forecast Hub (Sherratt et
 258 al. 2022). We present WIS for the subset of forecasts we explore visually for both the ensemble and
 259 surrogate model by date and forecast horizon (1 and 4 weeks).

260 To understand the relative performance of the surrogate model compared to the ensemble model, we
 261 calculate the relative performance (rWIS and rAE) by dividing the WIS/AE for the surrogate model by
 262 the WIS/AE of the ensemble model for all locations and forecast horizons. To maintain the propriety

263 of this score, we do this after first taking the means of scores for the relevant stratification. We explore
264 relative performance by forecast horizon, by month and horizon, and by location and horizon.

265 In addition to presenting the WIS for a subset of locations and the relative WIS for all locations, we
266 also calculate and visualise the empirical coverage, which is the percentage of observed values within a
267 given interval or below a given quantile, of both the surrogate and ensemble model for the 30%, 60%, and
268 90% prediction intervals and by quantile (Nikos I. Bosse et al. 2022b). We also calculate the bias (see
269 (Nikos I. Bosse et al. 2022b) and (Funk et al. 2019) for a more detailed definition) of both forecasting
270 approaches, stratified by forecast horizon. This metric aims to capture the tendency for a forecast to under
271 or over-predict. It captures the average proportion of the mass of the forecast distribution that is above or
272 below the true value (and so can range from -1 to 1) with an unbiased forecast having an average bias
273 value of 0. Lastly, we calculate and visualise the relative weighted interval score by quantile, stratified by
274 forecast horizon, to assess the relative difference in performance across the predictive distribution.

275 **Implementation**

276 The model is implemented in `stan` (S. D. Team 2021) and `R` (4.2.0) (R Core Team 2019) as an
277 extension of the baseline model from the `forecast.vocs` `R` package (0.0.9.7000) (Abbott 2021).
278 We note that our use of an indicator function introduces a discontinuity to the posterior making it less
279 suited for use with `stan`. Other model formulations without this feature would be more efficient and
280 robust. The `cmdstanr` `R` package (0.5.2) (Gabry and Češnovar 2021) is used for model fitting with 2
281 MCMC chains each having 1000 warm-up and 1000 sampling steps each (Gabry and Češnovar 2021).
282 `cmdstanr` surfaces several settings that trade-off between sampling speed and the robustness of the
283 approach. Here we take a conservative approach, as the model fit is not manually inspected during
284 real-time usage and due to the expected complexity of the posterior (Betancourt 2017), and set the `adapt`
285 `delta` setting to 0.99, and the maximum tree depth setting to 15. For real-time usage, convergence was
286 not assessed, but during model development, the Rhat diagnostic was used alongside feedback from
287 `cmdstanr` about the number of divergent transitions and exceedance of the maximum tree depth (Gabry
288 and Češnovar 2021). During development, posterior predictions were also visually compared to observed
289 data.

290 To download and manipulate forecasts from the European forecasting hub (E. C. F. H. Team 2021) we
291 use the `data.table` (1.14.2) (Dowle and Srinivasan 2021) and `gh` (1.3.0) (Bryan and Wickham
292 2021) `R` packages. We make use of further functionality from the `forecast.vocs` `R` package (Abbott
293 2021) to prepare data for forecasting, visualise forecasts and summary measures, and summarise forecasts.
294 Forecast evaluation is implemented using the `scoringutils` `R` package (1.0.0) (Nikos I. Bosse et al.
295 2022b), and the `scoringRules` `R` package (1.0.1) (Jordan, Krüger, and Lerch 2019).

296 To ensure the reproducibility of this analysis dependencies are managed using the `renv`
297 `R` package (0.14.0) (Ushey 2021) and a Dockerfile file along with a built Docker image
298 (Boettiger 2015) (via GitHub Actions (“About GitHub-hosted Runners” 2022)) is provided in
299 the code repository. Weekly forecasts were made using `renv` and based on GitHub Actions
300 free tier as available in 2022 to ensure they require limited compute and that our implemen-
301 tation is independent of local resources facilitating democratised access. The free GitHub
302 Actions runner we used for all forecasts was Ubuntu 20.04 based with 2 cores (x86_64), 7
303 GB of RAM, and 14 GB of SSD space. The code for this analysis can be found here: <https://github.com/epiforecasts/simplified-forecaster-evaluation> The code for the
304 forecasting model defined above along with the infrastructure required to forecast using GitHub Actions
305 can be found here: <https://github.com/seabbs/ecdc-weekly-growth-forecasts>
306 Versions archived on Zenodo are available (Abbott and Bosse 2022) and (Abbott and Sherratt 2022).

308 **RESULTS**

309 **Summary of the European COVID-19 Forecast Hub Setting**

310 In our study period, incidence rates across European nations and in the UK were primarily driven by
311 the spread of novel subvariants of concern related to the Omicron variant and changes in population
312 susceptibility. Many countries, such as the UK, saw large BA.1 waves in January, resulting in declining
313 incidence rates through February (Figure 1). From late February through to the end of May, most nations
314 saw another wave driven by BA.2. This wave typically saw lower reported incidence rates, and was
315 characterised by a lower peak than the BA.1 wave with a more gradual decrease in incidence. The end of

316 our study period was dominated by the gradual take-over of the BA.4/BA.5 subvariants that again had
317 a lower peak and lower absolute growth rates. Unlike earlier periods in the pandemic, our study period
318 did not see the use of new non-pharmaceutical interventions (NPIs) in response to increasing COVID-19
319 incidence in most locations. In addition, ascertainment rates likely reduced over time in most locations
320 due to reductions in routine testing and test availability. Whilst both the reduced use of NPIs and testing
321 generally occurred across nations our study period also marked an increase in the heterogeneity of the
322 response to the COVID-19 pandemic with nations changing policy at different times and to different
323 degrees. This is in contrast to the early COVID-19 pandemic response for which most nations took similar
324 actions at similar times.

325 We extracted forecasts starting from the 15th of January until the 19th of July 2022 for all countries
326 covered by the European forecasting hub (nations of the European Union, the European Free Trade
327 Agreement, and the United Kingdom, making 32 unique locations). In total 8846 forecasts were made
328 across all locations, with 27 unique forecast dates and 32 independent forecast models (including the
329 European hub baseline model). Of these models, 10 forecasted in at least 30 locations including our
330 original submission (referred to as `epiforecasts-EpiNow2` by the hub), and our surrogate model.
331 Of the remaining models submitted 16 were submitted in only one location. Single-location models
332 were clustered in a few locations, particularly in Germany and Poland (likely due to the folding of the
333 German/Poland forecasting hub into the European forecasting hub project (Sherratt et al. 2022)). Italy
334 was also an outlier with 4 models that submitted nowhere else. 4 models were submitted for between 3
335 and 30 locations and all these models varied the number of locations they submitted forecasts for over
336 time, potentially indicating manual curation or models targeted at specific conditions.

337 Across all forecast dates and locations the minimum number of independent forecasts was 4 with the
338 maximum being 20. The median number of independent forecasts per location and forecast date was 10.
339 All locations received forecasts from at least 10 models with the median number of forecast models per
340 location being 12. Coverage of forecast dates varied across submitted models with 8 models submitting
341 for all dates, 16 models submitting for at least 90% of dates, and 6 models submitting for fewer than 50%
342 of forecast dates. In general, there was no clear difference in forecast date coverage between models that
343 submitted for all locations vs a small subset but models with partial coverage of locations all also had
344 partial coverage of forecast dates.

345 63 observations, stratified by week and location, were defined to be anomalous within the study period
346 by the European Forecast Hub (E. C. F. H. Team 2021). Forecasts for these observations were excluded
347 as were forecasts for forecast weeks where they were the latest available data. Data anomalies were not
348 randomly distributed with some locations being particularly prone to data revisions including Lithuania
349 (with 23 weeks with data anomalies), and Portugal (with 13 weeks with data anomalies). Anomalies
350 were also not evenly distributed over time with a higher proportion occurring earlier in the study period
351 (potentially due to our choice to extract data from the 1st of September which effectively truncated
352 anomalies). 7.3% of forecasts were excluded across all horizons due to anomalies in the observed data.
353 Aggregated across horizons 10.3% of forecasts included at least one week with anomalous data.

354 **Forecast evaluation**

355 ***Visualisation of forecasts by horizon***

356 In our example set of locations, the absolute performance of the ensemble and the surrogate model was
357 visually similar on the log scale in all locations at short forecast horizons though this varied by location
358 (Figure 1 b). On the natural scale the difference in performance was more marked, especially for periods
359 of peak incidence and at longer horizons (Figure 1 a). Performance was not homogeneous across our
360 set of example locations with the surrogate model performing similarly to the ensemble in Slovakia
361 whilst in the United Kingdom and Germany the surrogate model performed substantially worse for some
362 forecast dates (Figure 1). For both the ensemble and the surrogate, performance decreased as the forecast
363 horizon increased with this being particularly noticeable for the surrogate model during periods of peak
364 incidence. In general, in the study period, the ensemble appeared to be better able to forecast peak
365 incidence. Both models forecast large reductions in incidence in Poland during May that did not occur
366 whilst only the ensemble forecast spuriously forecast similar large reductions in Germany during June. In
367 comparison to the ensemble model the surrogate model appeared less likely to place weight on unfeasibly
368 large reductions in incidence during periods of declining incidence but on other hand was more likely to
369 forecast continuing increases in incidence (for example in February in Slovakia and Poland).

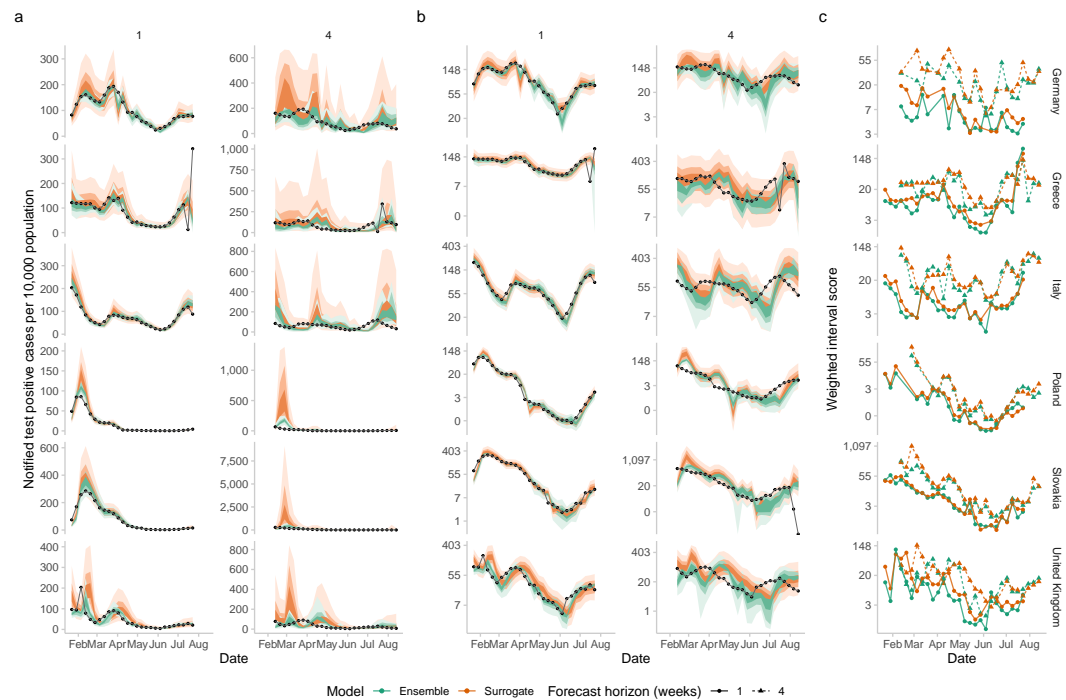


Figure 1. a.) Forecasts of notified test-positive cases (per 10,000 population) by epidemiological week in Germany, Greece, Italy, Poland, Slovakia, and the United Kingdom, by forecast horizon (one and four weeks) for the surrogate model (orange) and forecast ensemble (green). 30%, 60%, and 90% prediction intervals are shown. The black line and points are the notified cases as of the date of data extraction rather than those available at the time. b.) A replicate of a.) but with incidence rates on the log scale. c.) Weighted interval scores at the one-week and four-week forecast horizon by epidemiological week in Germany, Greece, Italy, Poland, Slovakia, and the United Kingdom on the log scale.

370 **Relative forecast evaluation**

371 Evaluating the ensemble and surrogate models using the WIS across all locations and forecast dates we
372 found that the mean relative performance of the surrogate model was 1.27 at the one-week horizon, 1.28
373 at the two-week horizon, 1.4 at the three-week horizon, and 1.69 at the four-week horizon, indicating
374 that the ensemble forecast outperformed the surrogate forecast for all horizons by at least 25% and that
375 the relative performance of the surrogate model degraded as the forecast horizon increased (Figure 2
376 c). Much of this outperformance, especially at longer forecast horizons, was driven by a small subset
377 of forecasts with relative performance having a heavy tail (Figure 2 a). If we instead consider median
378 relative performance (note this is not a proper scoring rule and should not be used to choose between
379 models) we find that, relative to the ensemble, the surrogate scored 1.21 at the one week horizon, 1.14
380 at the two week horizon, 1.2 at the three week horizon, and 1.28 at the four week horizon. This would
381 suggest that an increasingly skewed score distribution as the forecast horizon increased is responsible for
382 the increase in the mean relative score (Figure 2 a). 31% of individual surrogate forecasts scored better
383 than the comparable ensemble forecast, 68% performed within 50% of the comparable ensemble forecast,
384 and 17% had a more than 100% worse WIS than the comparable ensemble forecast.

385 If we consider only the median point forecast, using the absolute error, we see that the ensemble
386 forecast again outperformed the surrogate forecast (rAE for the surrogate compared to the ensemble 1.34).
387 If we instead consider the median of the absolute error we see that the difference in performance has
388 reduced indicating a similar skewed score distribution for point forecasts as for the whole predictive
389 distribution (rAE 1.11). Across forecast horizons the same pattern of outperformance holds. However, the
390 difference in relative performance was less than when the full probability distribution was accounted for,
391 with this becoming more marked as the forecast horizon increased (Figure 2 c).

392 The surrogate model's relative performance varied over time with substantially worse performance
393 from January to March compared to later in the year across all forecast horizons based on changes in
394 the relative score distribution and its summary statistics (Figure 2 b). The majority of the difference in
395 performance appeared to be driven by a thicker right tail with this being a particular feature of forecasts
396 at longer horizons. Forecast performance in March had a bimodal distribution at the four-week horizon
397 with a substantial fraction of surrogate forecasts outperforming the ensemble and a substantial fraction
398 substantially underperforming. This variation in performance may have been linked to the BA.2 wave
399 which peaked in most locations during this period if the surrogate model was more likely to overpredict
400 peak incidence than the ensemble forecast.

401 There was also substantial variation across forecast locations with the surrogate performing relatively
402 well in some locations at some forecast horizons, for example, the four-week horizon in the United
403 Kingdom, and badly in others, for example, the four-week forecast in Switzerland (Figure 2 c). In general,
404 across locations, as observed overall, relative forecast performance degraded across horizons with a
405 heavier right tail at longer horizons. Some locations showed less of this behaviour, for example, Spain,
406 and in some, it was very dominant, for example, Switzerland.

407 **Forecast calibration**

408 Overall the surrogate model was relatively well calibrated at the 30%, 60% and 90% prediction interval,
409 though with a tendency to be slightly underconfident, with empirical coverage of 30.5%, 62.5%, 92.3%
410 respectively. The ensemble model was less well calibrated, with a tendency to be overconfident with
411 empirical coverage of 24.8%, 51%, 79% respectively (Figure 3 a). When stratified by forecast horizon the
412 ensemble forecast was best calibrated at the one-week forecast horizon, and then became progressively
413 less well calibrated as the forecast horizon increased (Figure 3 a). In comparison, the surrogate forecast
414 was less well calibrated than the ensemble forecast at the one-week forecast horizon with a tendency to
415 have a larger empirical coverage than required (Figure 3 a). At longer horizons and narrower prediction
416 intervals, the surrogate forecast became better calibrated though with a tendency to be overconfident.
417 This was not the case for the 90% prediction interval where the surrogate model covered more than the
418 expected interval, for all horizons, indicating forecasts were overly uncertain for this interval regardless of
419 the horizon.

420 Stratifying calibration by quantile and forecast horizon the ensemble forecast was conservative at all
421 horizons for quantiles larger than the median whilst being comparably well calibrated for intervals below
422 the median (Figure 3 b). This behaviour became more prominent as the forecast horizon increased. In
423 contrast, the surrogate forecast was generally equally well calibrated across horizons with a tendency
424 to be under confident for intervals above the median. At longer horizons, however, quantiles below the

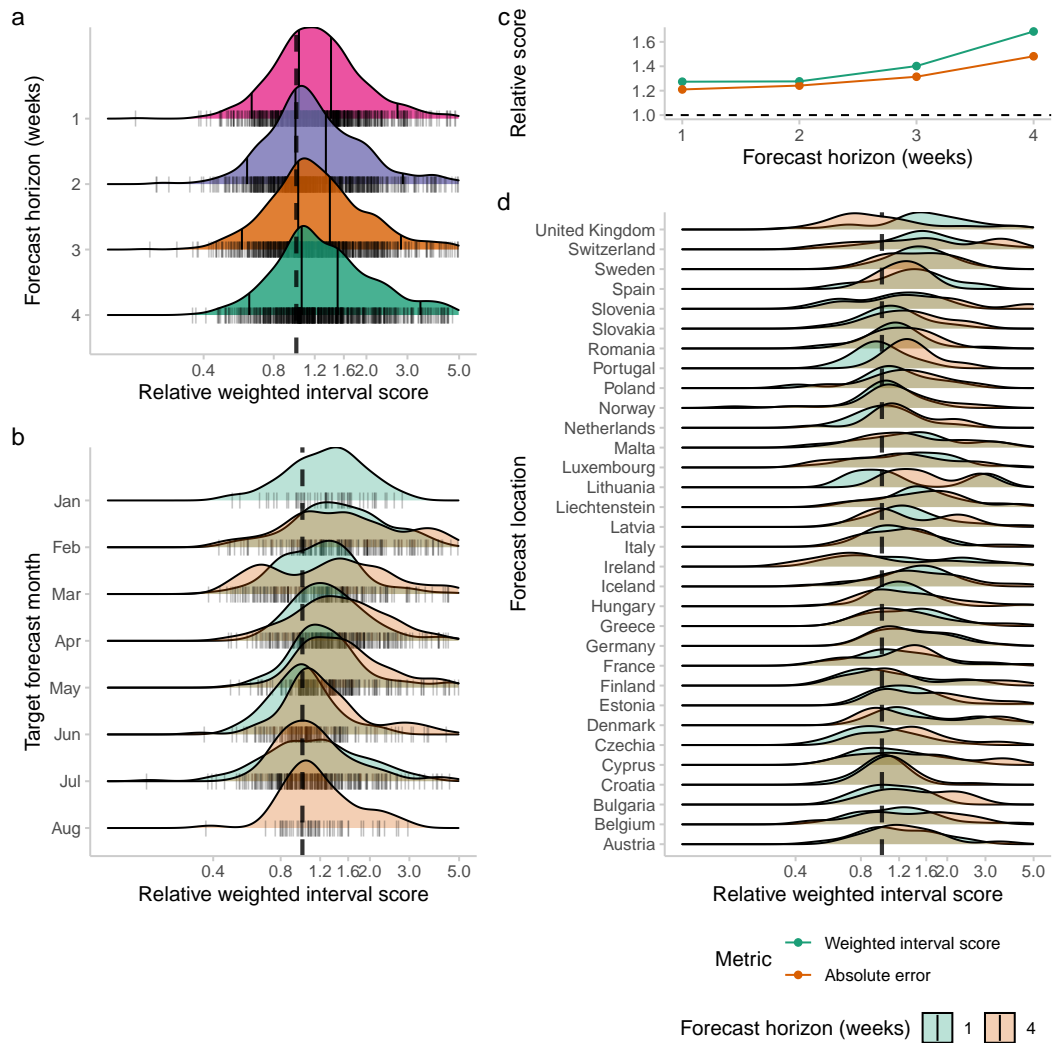


Figure 2. Relative weighted interval score by location, horizon, and forecast date for the surrogate forecast model compared to the ensemble forecast model on the log scale. a.) The density of the relative score by horizon. Horizontal black lines give the 5%, 35%, 65%, and 95% quantiles. b.) The density of the relative score by month for a given forecast horizon stratified by the one and four-week forecast horizon. c.) The average relative weighted interval score and absolute error for the surrogate model compared to the ensemble forecast by forecast horizon. d.) The density of the relative score by forecast location stratified by the one and four-week forecast horizon. The dashed line on all plots indicates when the ensemble forecast is equivalent to the surrogate forecast. The vertical black lines on the y-axis give individual relative scores.

425 median were over confident.

426 Breaking down the relative weighted interval score by forecast interval we observe that the surrogate
427 model produces forecasts that differ most from the ensemble in the outer intervals and in particular the
428 tails of the forecast (Figure 3 c). This is true across forecast horizons but the magnitude of the difference
429 increases.

430 Calculating the bias of the forecasts from each model we see that the (Figure 3 d) ensemble forecast
431 is initially biased towards underprediction but this bias reduces as the forecast horizon increases. In
432 comparison, the surrogate forecast model is biased towards overprediction for all forecast horizons with
433 the magnitude of this bias appearing to increase linearly with the forecast horizon.

434 **DISCUSSION**

435 **Summary**

436 In this study, we defined a surrogate model aiming to replicate some of the observed behaviour of the
437 European Forecast Hub multi-team ensemble for forecasting test-positive reported COVID-19 cases in
438 European nations. We first defined a set of assumptions for how the surrogate model should behave
439 based on our observations of the European Forecast Hub ensemble, and our experience submitting
440 forecasts to various Forecast Hubs. We aimed for a model that could be easily understood, that produced
441 epidemiologically meaningful summary statistics, and that could be run with low compute resources. We
442 further provide a fully reproducible workflow for running and evaluating this model using GitHub actions
443 facilitating others to do the same.

444 Over the 6 months of the study period, we found that our surrogate model produced forecasts that were
445 visually similar to those from the Forecast Hub ensemble on the log scale though with greater uncertainty.
446 Visual differences were more marked on the natural scale with the surrogate model forecasting spuriously
447 high peak incidence. In a subset of example locations, we observed some variation in performance
448 across locations, that the ensemble better-captured peak incidence, and that the surrogate model appeared
449 biased toward overprediction. Evaluating the relative performance of the surrogate model compared to
450 the European Forecast Hub ensemble we found that the mean performance was substantially worse and
451 that relative performance decreased with the forecast horizon. The median forecast performance of the
452 surrogate model was also worse when compared to forecasts from the ensemble though the majority of
453 surrogate forecasts were within 50% of the performance observed for the ensemble forecast. The difference
454 in mean and median relative performance suggested a skewed distribution in scores, which we confirmed
455 visually. This means that a relatively small fraction of forecasts were responsible for a substantial portion
456 of the difference in performance. Evaluating point forecast performance indicated a similar pattern of
457 performance as that observed using the full predictive distribution though the relative performance of
458 the surrogate model generally improved. Performance varied by location and forecast date with the
459 surrogate model performing worse in the first part of 2022 which may have been linked to incidence rates
460 peaking across forecast locations linked to the spread of BA.2. In general, the relative performance of the
461 surrogate model degraded as forecast horizons increased with the distribution of relative performance
462 having an increasingly heavy right tail as the forecast horizon increased indicating a greater share of
463 forecasts performing very poorly in comparison to the hub ensemble. The Forecast Hub ensemble was
464 poorly calibrated, particularly at longer forecast horizons and larger prediction intervals, compared to
465 the surrogate model though the surrogate model tended to be overly uncertain at large intervals. The
466 ensemble forecast was biased towards under-prediction at short to medium forecast horizons but unbiased
467 at longer horizons. In comparison, the surrogate model was biased towards overprediction and this bias
468 increased linearly with the forecast horizon.

469 **Strengths and Weaknesses**

470 Our study benefits from having been conducted using forecasts produced in real-time, rather than
471 retrospectively, and submitted to an independent forecast research hub (though we note the overlap
472 between authors on this study and the European Forecast Hub (Sherratt et al. 2022)). This means that our
473 results are not subject to hindsight bias. The downside of this approach is that it was not possible to update
474 the surrogate model over time in response to the initial evaluation or to explore other parameterisations that
475 might be more successful of which there are likely several. However, as our study has been conducted with
476 a focus on reproducibility and openness our findings can be replicated or extended by others regardless of
477 compute availability (due to our use of GitHub actions as a compute platform which is freely available to

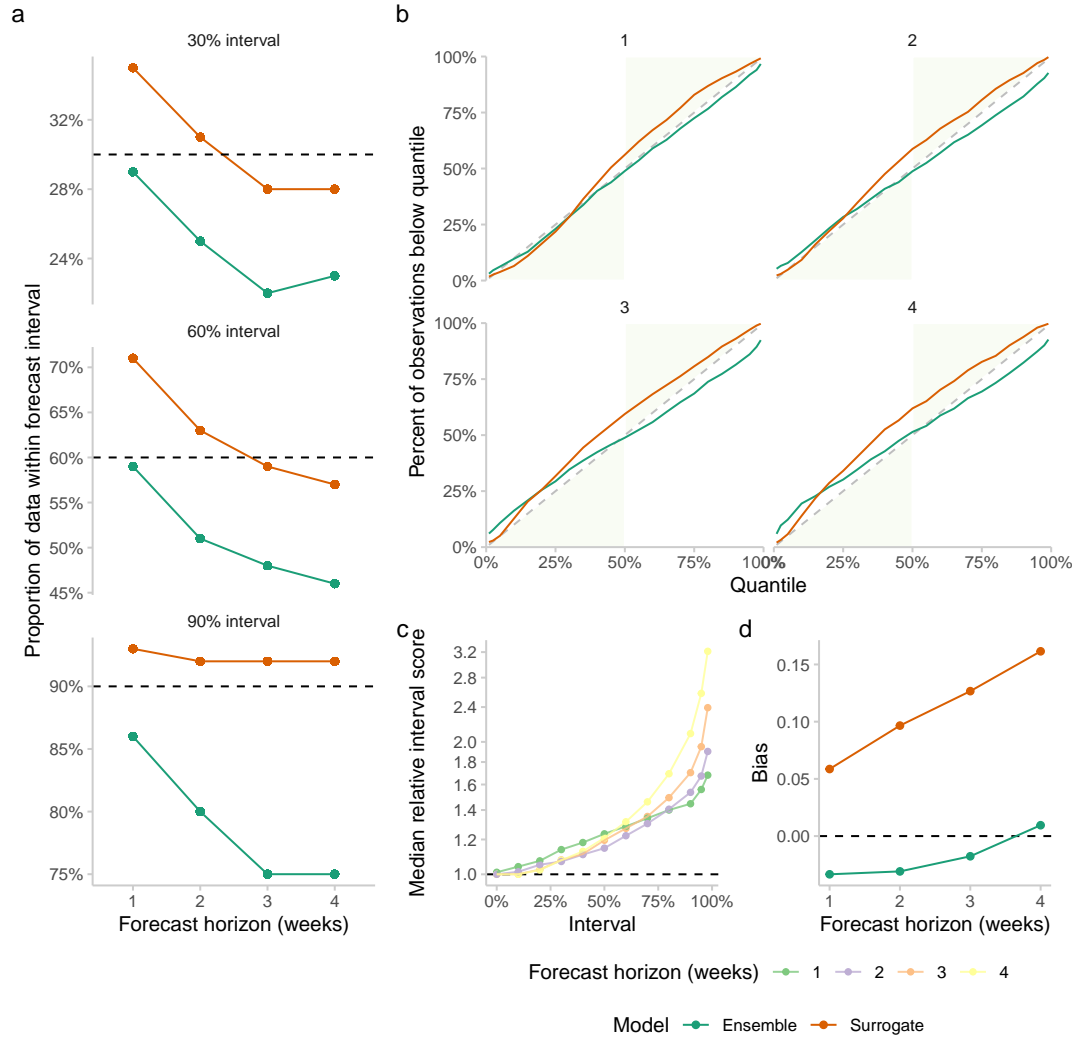


Figure 3. a.) Empirical coverage of the surrogate (orange) and ensemble (green) forecasts at the 90%, 60%, and 30% prediction intervals stratified by forecast horizon. Ideally, a well-calibrated forecast should have empirical coverage for a given prediction interval that equals the nominal level of the interval (i.e., 30%, 60% and 90%, respectively). b.) Empirical coverage by quantile for both the surrogate and ensemble forecasts. A well-calibrated forecast should have empirical quantiles that match the theoretical ones. The green area of this figure corresponds to conservative forecasts. c.) Median relative weighted interval score by quantile and forecast horizon for the surrogate forecast compared to the ensemble forecast. d.) The bias of the ensemble and surrogate forecasts stratified by horizon.

478 researchers). An additional downside to this approach is that the hub ensemble includes forecasts from our
479 surrogate model, increasing the similarity between the two approaches. This is difficult to avoid without
480 retrospectively re-calculating the ensemble using the same approach as taken by the hub which would
481 reduce the independence of the hub ensemble as a source of truth to compare our forecasts against. Given
482 the number of forecasts submitted in most locations and the European Forecast Hubs' practice of not
483 calculating an ensemble when fewer than 3 independent forecasts were available, the bias in our results
484 caused by this limitation should be relatively small. Notably in this study, we focussed on replicating
485 the Forecast Hub ensembles' observed behaviour rather than attempting to define an optimal forecast
486 for forecast consumers. It is possible that if we had instead aimed to develop a forecast methodology
487 that minimised the evaluation criteria we planned to use, especially if we relaxed our assumed compute
488 resource constraints, we would have produced forecasts that performed better relative to the hub ensemble.
489 However, if we start from the view that the Forecast Hub ensemble has traits that are desirable for use by
490 policy-makers (i.e robustness and good average performance), which can be found widely in the literature
491 (Cramer et al. 2022; J. Bracher et al. 2021; Sherratt et al. 2022), then our approach may make sense as a
492 way of producing a "good" forecast without sacrificing interpretability.

493 Developing forecast methodologies with limited resources is critical as whilst improving predictive
494 performance is a key goal of short-term forecasting it is also important that forecast models be accessible
495 as this makes it easier to iteratively improve them, and makes them more flexible when used in real-time
496 settings (Osthus 2022). An example of the lack of flexibility of the Forecast Hub ensemble, and its
497 constituent models, is the ensembles response to upswings linked to variant dynamics, with the growth
498 of one variant being temporally hidden by the decline of another. Rather than forecasting this ahead of
499 time the Forecast ensembles generally only reacted to changes in the observed data indicating that variant
500 information was not being used by most forecasters. Unlike the Forecast Hub ensemble the surrogate
501 model can be, and indeed has been (Abbott, Sherratt, and Funk 2021), easily extended to capture this.
502 Other examples where additional transient information is available to forecasters can be readily thought
503 of implying this is a general advantage of simpler methods.

504 Our focus on replicating the performance of the hub ensemble is also useful as the surrogate model
505 may highlight some of the emergent behaviour of the ensemble captured in our assumptions, such as auto-
506 correlation across time points, and the growth rate tending towards zero as the forecast horizon increases.
507 It also highlights some of the differences between our surrogate forecast model and the ensemble that
508 may lead to new insights into the mechanisms leading to the ensemble's behaviour, such as the generally
509 poor coverage of the ensemble that could not be explained by the assumptions we used in developing
510 our surrogate methodology. Whilst we normalised reported cases to be population-adjusted incidence
511 rates, and so can more easily compare across locations than using the approach commonly implemented
512 by the Forecast Hubs (Cramer et al. 2022; J. Bracher et al. 2021; Sherratt et al. 2022), our results are still
513 conditional on the use of the weighted interval score as an evaluation metric. As this proper scoring rule
514 scales with the order of magnitude of the predicted quantity this means that forecasts during periods of
515 higher incidence are given more weight than forecasts from periods of low incidence. It also means that
516 overprediction is penalised more than underprediction as incidence rates are bounded at zero but relatively
517 weakly bounded by populations at the upper bound (as incidence rates are typically only a small fraction
518 of the overall population). This bias could explain the relatively poor performance of the surrogate model,
519 compared to the ensemble, despite the surrogate model being comparably well-calibrated. We considered
520 alternative methods of forecast evaluation that would be robust to this potential source of bias but choose
521 to stick relatively closely to the methodology used by the European Forecast Hub (Sherratt et al. 2022),
522 aside from the use of population weighting to facilitate comparison between forecast locations, as these
523 choices inform the development of submitted models and so are key to our findings.

524 **Literature context**

525 There are no other studies in the epidemiology literature which we are aware of that attempt to develop a
526 forecasting model based on the observed behaviour of a multi-team, multi-model ensemble. Few studies
527 focus on delivering computationally feasible forecasting models in a reproducible framework backed
528 by an openly accessible compute platform. However, the US (Cramer et al. 2022), European (Sherratt
529 et al. 2022), and Germany/Poland (J. Bracher et al. 2021) forecasting hubs have published a range of
530 evaluations of forecasts submitted to their platforms and the relative performance of their ensembles.
531 In general, these studies have struggled to draw general conclusions about the structural assumptions

532 of forecast models they consider “good” (generally they have defined this as minimising the weighted
533 interval score, as in this study).

534 The poor calibration of the forecast ensembles produced by median Hub ensembles has been noted
535 repeatedly (Cramer et al. 2022; J. Bracher et al. 2021; Sherratt et al. 2022) but little progress has been
536 made in understanding the causes or suggesting alternatives. Progress in understanding which structural
537 model features lead to better infectious disease forecasts has been limited. The US Forecast Hub identified
538 the top 5 performing models and noted the structural assumptions they made, but couldn’t directly link
539 assumptions with performance (Cramer et al. 2022). They also did not extensively compare and contrast
540 these conclusions to arrive at a set of desired forecast assumptions (as done in this study to motivate the
541 surrogate model), or explore the performance of a forecasting model designed with these assumptions
542 in mind. Similarly, the Germany and Poland forecasting hubs were able to identify forecast models that
543 performed comparably as well as their ensemble forecasts but did not derive structural assumptions that
544 led to this out-performance or detail explicitly what the desirable performance characteristics would be,
545 aside from optimising the weighted interval score. All comparable Forecast Hub projects found that their
546 ensemble was often the best choice, had desirable characteristics such as robustness - though this was
547 rarely fully defined - and should be the output used by forecast consumers (Cramer et al. 2022; J. Bracher
548 et al. 2021; Sherratt et al. 2022). In general, during the study period, all projects used the same unweighted
549 median ensemble forecast of all submissions. The US (Ray et al. 2022), and European (Sherratt et al.
550 2022), forecasting hub also evaluated a range of other ensemble approaches, such as inverse weighted
551 interval score weighting, unweighted ensembles of a selection of models based on recent performance,
552 and mean ensembling. Work on this is still ongoing but these more complex ensembling approaches were
553 shown to outperform the median of all submitted forecasts in many cases in the case of the US forecasting
554 hub and did not outperform in the case of the European forecasting hub. No Forecast Hub has switched to
555 these alternative ensemble designs for their operational forecast of reported cases, though the US hub
556 has switched to a trained ensemble for death forecasts. This suggests that the hub teams do not think
557 the evidence base is strong enough for trained ensembles to be used by forecast consumers for reported
558 cases and hence the median of all submitted forecasts remains the community-suggested default ensemble
559 option and a sensible target for our study.

560 Other studies have been published evaluating single forecast models in comparison to ensemble
561 performance from the Forecast Hub. In general, these have not focussed on replicating ensemble
562 behaviour but rather optimising the target evaluation metric. Our previous work also highlighted the lack
563 of calibration in an ensemble forecast from the Germany/Poland forecasting hub compared to forecasts
564 from epidemiological models and noted the bias towards underprediction observed in the ensemble
565 forecasts and not in our model-based forecasts (Nikos I. Bosse et al. 2022a; J. Bracher et al. 2021).
566 Finally, our results are potentially sensitive to the definition used to define anomalous observations
567 (generally related to retrospective data revisions). Here we follow the practice of the European Forecast
568 Hub (E. C. F. H. Team 2021) of excluding forecasts for weeks with a data revision of more than 5% and
569 forecasts made based on data that is subsequently revised by more than 5%.

570 **Further work**

571 Whilst we derived our surrogate model from a range of assumptions based on observing ensemble
572 forecasts behaviour and the behaviour and structure of submitted models avenues for future improvement
573 remain in terms of improving the approach used to elicit these observations. In follow-up work, a more
574 rigorous approach to this could be taken to further refine this set of assumptions, in particular using
575 the input of a wider pool of researchers. The findings from our study may also be useful for informing
576 this improved set of assumptions. A particular focus should be on understanding why our surrogate
577 model was liable to overestimate peak incidence and what simple additional assumptions may be used to
578 mitigate this. In addition, the model we derived based on our assumptions was likely not optimal both
579 in terms of compute time and accuracy at reproducing ensemble-like behaviour. Models with a more
580 complex auto-correlation structure and more refined approaches to localised trends should be explored to
581 improve relative performance to ensemble forecasts. An example of a family of possible approaches are
582 structural time series models which have many of the characteristics implied by our assumptions for how
583 forecast ensembles typically operate. As we identified that the tails of our predictive distributions were
584 responsible for a large proportion of the difference in performance compared to the forecast ensemble
585 it may be the case that post-processing of forecasts from our surrogate model would enhance their

586 similarity to the forecast ensemble. This seems likely to improve out-of-sample performance but does not
587 help with understanding the implicit assumptions driving the performance of multi-model, multi-team
588 infectious disease forecast ensembles. As we have hypothesised that the use of absolute scoring measures
589 is inappropriate and leads to performance characteristics that are unlikely to be favoured by forecast
590 stakeholders more work should be done in this area. If new forecast ensemble methods are adopted as best
591 practice by Forecast Hubs then follow-up work attempting to create surrogate forecast models should also
592 use these approaches and this will likely alter the observed characteristics of the hub ensemble forecasts,
593 for example, the tendency to be poorly calibrated. In September 2022, GitHub announced support for
594 hosted GitHub Action runners with additional compute power (“GitHub Actions Larger Runners - Are
595 Now in Public Beta” 2022). Whilst a paid feature this may allow more compute-intensive models, with
596 fewer potential performance trade-offs, to be easily democratised though only if funds are available to
597 support the hosting costs. One potential research area is to explore forecasting methods that can be used
598 with a range of computing resources though this would require extensive evaluation and documentation to
599 make it clear to users what the trade-offs between compute usage and forecast performance are. More work
600 is needed to understand the best practice treatment of data revisions when evaluating forecasts and the
601 potential bias these may cause. Lastly, here we have only explored a surrogate for an ensemble for a single
602 disease, a limited set of locations, and a single target (incident cases), meaning our findings are difficult
603 to generalise. Follow-up work should explore whether this behaviour holds across diseases, locations, and
604 epidemiological targets where the behaviour of ensembles is notably different. However, this is limited to
605 infectious diseases with similar large-scale forecast ensembling projects. These projects remain relatively
606 rare despite them showing obvious promise to improve the forecasts available to stakeholders.

607 **CONCLUSIONS**

608 We conclude that our simplified forecast model may have captured some of the dynamics of the hub
609 ensemble but that more work needs to be done to understand the epidemiological model that represents its
610 behaviour and whether or not this is the optimal choice for stakeholders’ requirements. We also conclude
611 that our findings may be largely driven by the choice of evaluation measure used by the Forecast Hub.
612 While this measure has desirable mathematical properties and is routinely used in a similar form e.g., in
613 weather forecasting, it is subject to debate whether it appropriately reflects forecast users’ requirements
614 and perceptions as to what makes a good forecast. Our work is useful for forecast users to understand the
615 inherent assumptions of the forecasts they are making use of and to researchers thinking about how to
616 develop forecasts that perform similarly to current multi-model and multi-team forecast ensembles that
617 are trusted by stakeholders.

618 **ACKNOWLEDGMENTS**

619 We thank the ECDC for supporting the forecasting hub, and all forecasters who submitted forecasts for
620 making this study possible. We thank the Forecast Hub team for publishing all data in an accessible
621 format. We thank the Epiforecast group for a productive discussion of an early version of this analysis.
622 We thank Molly for being a good Labrador.

623 **ADDITIONAL INFORMATION AND DECLARATIONS**

624 **Competing interests**

625 SF, JB, JS, BB, and HG have coordinated Forecast Hub platforms. SF and KS received funding from the
626 European Center for Disease Prevention and Control to this end.

627 **Author contribution**

628 SA conceived the study, developed the initial set of assumptions for the surrogate model, implemented
629 the model into code, designed and conducted the forecast evaluation, and wrote the first draft of the
630 manuscript. All other authors provided feedback on the manuscript and analyses and contributed to
631 revisions. HG and KS reviewed the code and reproducibility of the analyses.

632 **Data availability**

633 All data and code are available here:

634 <https://github.com/epiforecasts/simplified-forecaster-evaluation>
635 And are archived here:
636 <https://doi.org/10.5281/zenodo.7189308>, <https://doi.org/10.5281/zenodo.7189620>

638 **Funding**

639 SA,SF, KS and HG were funded by a Wellcome senior fellowship to SF (210758/Z/18/Z), KS and HG
640 were further funded by an ECDC grant to SF. JB acknowledges support from the Helmholtz Foundation
641 via the SIM-610 CARD Information and Data Science Pilot Project.

642 **REFERENCES**

- 643 Abbott, Sam. 2021. "Forecast.vocs: Forecast Case and Sequence Notifications Using Variant of Concern
644 Strain Dynamics." *Zenodo*. <https://doi.org/10.5281/zenodo.5559016>.
- 645 Abbott, Sam, and Nikos Bosse. 2022. "Epiforecasts/Simplified-Forecaster-Evaluation." <https://doi.org/10.5281/zenodo.7189309>.
- 647 Abbott, Sam, Joel Hellewell, Katharine Sherratt, Katelyn Gostic, Joe Hickson, Hamada S. Badr, Michael
648 DeWitt, Robin Thompson, EpiForecasts, and Sebastian Funk. 2020. *EpiNow2: Estimate Real-Time
649 Case Counts and Time-Varying Epidemiological Parameters*. <https://doi.org/10.5281/zenodo.3957489>.
- 651 Abbott, Sam, Joel Hellewell, Robin N Thompson, Katharine Sherratt, Hamish P Gibbs, Nikos I Bosse,
652 James D Munday, et al. 2020. "Estimating the Time-Varying Reproduction Number of SARS-
653 CoV-2 Using National and Subnational Case Counts." *Wellcome Open Res.* 5 (December): 112.
654 <https://doi.org/10.12688/wellcomeopenres.16006.2>.
- 655 Abbott, Sam, and Kath Sherratt. 2022. "Seabbs/Ecdc-Weekly-Growth-Forecasts." <https://doi.org/10.5281/zenodo.7189621>.
- 657 Abbott, Sam, Katharine Sherratt, and Sebastian Funk. 2021. "Real-Time Estimation of the Time-
658 Varying Transmission Advantage of Omicron in England Using S-Gene Target Status as a Proxy."
659 <https://doi.org/10.5281/zenodo.5812298>.
- 660 "About GitHub-hosted Runners." 2022. [https://ghdocs-prod.azurewebsites.net/en/
661 actions/using-github-hosted-runners/about-github-hosted-runners](https://ghdocs-prod.azurewebsites.net/en/actions/using-github-hosted-runners/about-github-hosted-runners).
- 662 Adamik, Barbara, Marek Bawiec, Viktor Bezborodov, Wolfgang Bock, Marcin Bodych, Jan Pablo Burgard,
663 Thomas Götze, et al. 2020. "Mitigation and Herd Immunity Strategy for COVID-19 Is Likely to Fail."
664 *bioRxiv*. medRxiv. <https://doi.org/10.1101/2020.03.25.20043109>.
- 665 Betancourt, Michael. 2017. "Diagnosing Biased Inference with Divergences." *Stan Case Studies* 4.
666 [https://mc-stan.org/users/documentation/case-studies/divergences_
667 and_bias.html](https://mc-stan.org/users/documentation/case-studies/divergences_and_bias.html).
- 668 Boettiger, Carl. 2015. "An introduction to Docker for reproducible research." *ACM SIGOPS Operating
669 Systems Review* 49 (1): 71–79.
- 670 Bosse, Nikos I, Sam Abbott, Johannes Bracher, Habakuk Hain, Billy J Quilty, Mark Jit, Centre for the
671 Mathematical Modelling of Infectious Diseases COVID-19 Working Group, Edwin van Leeuwen,
672 Anne Cori, and Sebastian Funk. 2022a. "Comparing Human and Model-Based Forecasts of COVID-
673 19 in Germany and Poland." *PLoS Comput. Biol.* 18 (9): e1010405. [https://doi.org/10.
674 1371/journal.pcbi.1010405](https://doi.org/10.1371/journal.pcbi.1010405).
- 675 Bosse, Nikos I, Hugo Gruson, Anne Cori, Edwin van Leeuwen, Sebastian Funk, and Sam Abbott. 2022b.
676 "Evaluating Forecasts with Scoringutils in r." *arXiv*. [https://doi.org/10.48550/ARXIV.
677 2205.07090](https://doi.org/10.48550/ARXIV.2205.07090).
- 678 Bracher, Johannes, Evan L. Ray, Tilmann Gneiting, and Nicholas G. Reich. 2021. "Evaluating Epidemic
679 Forecasts in an Interval Format." *PLoS Computational Biology* 17 (2): e1008618. [https://doi.
680 org/10.1371/journal.pcbi.1008618](https://doi.org/10.1371/journal.pcbi.1008618).
- 681 Bracher, Johannes, Daniel Wolfram, Jannik Deuschel, Konstantin Görden, Jakob L Ketterer, Alexander
682 Ullrich, Sam Abbott, et al. 2022. "National and Subnational Short-Term Forecasting of COVID-19 in
683 Germany and Poland During Early 2021." *medRxiv*. [https://doi.org/10.1101/2021.11.
684 05.21265810](https://doi.org/10.1101/2021.11.05.21265810).
- 685 Bracher, J, D Wolfram, J Deuschel, K Görden, J L Ketterer, A Ullrich, S Abbott, et al. 2021. "A Pre-
686 Registered Short-Term Forecasting Study of COVID-19 in Germany and Poland During the Second

687 Wave." *Nat. Commun.* 12 (1): 5173. <https://doi.org/10.1038/s41467-021-25207-0>.

688 Bryan, Jennifer, and Hadley Wickham. 2021. *Gh: 'GitHub' 'API'*. <https://CRAN.R-project.org/package=gh>.

689

690 Castelletti, A, S Galelli, M Ratto, R Soncini-Sessa, and P C Young. 2012. "A General Framework for
691 Dynamic Emulation Modelling in Environmental Problems." *Environmental Modelling & Software*
692 34 (June): 5–18. <https://doi.org/10.1016/j.envsoft.2012.01.002>.

693 Castro, Lauren, Geoffrey Fairchild, Isaac Michaud, and Dave Osthus. 2021. "COFFEE: COVID-19
694 Forecasts Using Fast Evaluations and Estimation," October. <https://arxiv.org/abs/2110.01546>.

695

696 Charles, Giovanni, Timothy M Wolock, Peter Winskill, Azra Ghani, Samir Bhatt, and Seth Flaxman.
697 2022. "Seq2Seq Surrogates of Epidemic Models to Facilitate Bayesian Inference," September.
698 <https://arxiv.org/abs/2209.09617>.

699 Cramer, Estee Y, Evan L Ray, Velma K Lopez, Johannes Bracher, Andrea Brennen, Alvaro J Castro
700 Rivadeneira, Aaron Gerding, et al. 2022. "Evaluation of Individual and Ensemble Probabilistic
701 Forecasts of COVID-19 Mortality in the United States." *Proc. Natl. Acad. Sci. U. S. A.* 119 (15):
702 e2113561119. <https://doi.org/10.1073/pnas.2113561119>.

703 Department of Health, NM-IBIS. n.d. "MMWR Week Description and Corresponding Cal-
704 endar Dates (2006–2025)." [https://ibis.health.state.nm.us/resource/
705 MMWRWeekCalendar.html](https://ibis.health.state.nm.us/resource/MMWRWeekCalendar.html).

706 Dong, Ensheng, Hongru Du, and Lauren Gardner. 2020. "An Interactive Web-Based Dashboard to Track
707 COVID-19 in Real Time." *Lancet Infect. Dis.* 20 (5): 533–34. [https://doi.org/10.1016/
708 S1473-3099\(20\)30120-1](https://doi.org/10.1016/S1473-3099(20)30120-1).

709 Dowle, Matt, and Arun Srinivasan. 2021. *Data.table: Extension of 'Data.frame'*. [https://CRAN.
710 R-project.org/package=data.table](https://CRAN.R-project.org/package=data.table).

711 Edwards, Tamsin L, Sophie Nowicki, Ben Marzeion, Regine Hock, Heiko Goelzer, H el ene Seroussi,
712 Nicolas C Jourdain, et al. 2021. "Projected Land Ice Contributions to Twenty-First-Century Sea Level
713 Rise." *Nature* 593 (7857): 74–82. <https://doi.org/10.1038/s41586-021-03302-y>.

714 Funk, Sebastian, Anton Camacho, Adam J Kucharski, Rachel Lowe, Rosalind M Eggo, and W John
715 Edmunds. 2019. "Assessing the Performance of Real-Time Epidemic Forecasts: A Case Study of
716 Ebola in the Western Area Region of Sierra Leone, 2014–15." *PLoS Comput. Biol.* 15 (2): e1006785.
717 <https://doi.org/10.1371/journal.pcbi.1006785>.

718 Gabry, Jonah, and Rok  e snovar. 2021. *Cmdstanr: R Interface to 'CmdStan'*.

719 "GitHub Actions Larger Runners - Are Now in Public Beta." 2022. [https://github.blog/
720 changelog/2022-09-01-github-actions-larger-runners-are-now-in-public-beta/;](https://github.blog/changelog/2022-09-01-github-actions-larger-runners-are-now-in-public-beta/)
721 GitHub.

722 Gostic, Katelyn M, Lauren McGough, Edward B Baskerville, Sam Abbott, Keya Joshi, Christine Tedijanto,
723 Rebecca Kahn, et al. 2020. "Practical Considerations for Measuring the Effective Reproductive Num-
724 ber, Rt." *PLoS Comput. Biol.* 16 (12): e1008409. [https://doi.org/10.1371/journal.
725 pcbi.1008409](https://doi.org/10.1371/journal.pcbi.1008409).

726 IPCC. n.d. "AR6 Synthesis Report: Climate Change 2022." [https://www.ipcc.ch/report/
727 sixth-assessment-report-cycle/](https://www.ipcc.ch/report/sixth-assessment-report-cycle/).

728 Iskauskas, Andrew, Ian Vernon, Michael Goldstein, Danny Scarponi, Nicky McCreesh, Trevelyan J
729 McKinley, and Richard G White. 2022. "Emulation and History Matching Using the Hmer Package,"
730 September. <https://arxiv.org/abs/2209.05265>.

731 Jordan, Alexander, Fabian Kr uger, and Sebastian Lerch. 2019. "Evaluating Probabilistic Forecasts with
732 scoringRules." *Journal of Statistical Software* 90 (12): 1–37. [https://doi.org/10.18637/
733 jss.v090.i12](https://doi.org/10.18637/jss.v090.i12).

734 Karlen, D. 2020. "Characterizing the Spread of CoViD-19."

735 Li, Michael Lingzhi, Hamza Tazi Bouardi, Omar Skali Lami, Thomas A Trikalinos, Nikolaos K Trichakis,
736 and Dimitris Bertsimas. 2021. "Forecasting COVID-19 and Analyzing the Effect of Government
737 Interventions." *medRxiv*. <https://doi.org/10.1101/2020.06.23.20138693>.

738 Meakin, Sophie, Sam Abbott, Nikos Bosse, James Munday, Hugo Gruson, Joel Hellewell, Katharine
739 Sherratt, CMMID COVID-19 Working Group, and Sebastian Funk. 2022. "Comparative Assessment
740 of Methods for Short-Term Forecasts of COVID-19 Hospital Admissions in England at the Local
741 Level." *BMC Med.* 20 (1): 86. <https://doi.org/10.1186/s12916-022-02271-x>.

742 Nixon, Kristen, Sonia Jindal, Felix Parker, Nicholas G Reich, Kimia Ghobadi, Elizabeth C Lee, Shaun
743 Truelove, and Lauren Gardner. 2022. “An Evaluation of Prospective COVID-19 Modelling Studies
744 in the USA: From Data to Science Translation.” *Lancet Digit Health* 4 (10): e738–47. [https://doi.org/10.1016/S2589-7500\(22\)00148-0](https://doi.org/10.1016/S2589-7500(22)00148-0).

745 Osthus, Dave. 2022. “Fast and Accurate Influenza Forecasting in the United States with Inferno.” *PLoS*
746 *Comput. Biol.* 18 (1): e1008651. <https://doi.org/10.1371/journal.pcbi.1008651>.

747 “Pricing - Linux Virtual Machines.” 2022. [https://azure.microsoft.com/en-gb/pricing/
748 details/virtual-machines/linux/](https://azure.microsoft.com/en-gb/pricing/details/virtual-machines/linux/).

749 R Core Team. 2019. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R
750 Foundation for Statistical Computing. <https://www.R-project.org/>.

751 Rakowski, Franciszek, Magdalena Gruziel, Łukasz Bieniasz-Krzywiec, and Jan P Radomski. 2010.
752 “Influenza Epidemic Spread Simulation for Poland — a Large Scale, Individual Based Model Study.”
753 *Physica A: Statistical Mechanics and Its Applications* 389 (16): 3149–65. [https://doi.org/
754 10.1016/j.physa.2010.04.029](https://doi.org/10.1016/j.physa.2010.04.029).

755 Ray, Evan L, Logan C Brooks, Jacob Bien, Matthew Biggerstaff, Nikos I Bosse, Johannes Bracher, Estee
756 Y Cramer, et al. 2022. “Comparing Trained and Untrained Probabilistic Ensemble Forecasts of
757 COVID-19 Cases and Deaths in the United States.” *Int. J. Forecast.*, July. [https://doi.org/
758 10.1016/j.ijforecast.2022.06.005](https://doi.org/10.1016/j.ijforecast.2022.06.005).

759 Reich, Nicholas G, Justin Lessler, Sebastian Funk, Cecile Viboud, Alessandro Vespignani, Ryan J
760 Tibshirani, Katriona Shea, et al. 2022. “Collaborative Hubs: Making the Most of Predictive Epidemic
761 Modeling.” *Am. J. Public Health*, April, e1–4. [https://doi.org/10.2105/ajph.2022.
762 306831](https://doi.org/10.2105/ajph.2022.306831).

763 Sherratt, Katharine, Hugo Gruson, Helen Johnson, Rene Niehus, Bastian Prasse, Frank Sandman,
764 Jannik Deuschel, et al. 2022. “Predictive Performance of Multi-Model Ensemble Forecasts of
765 COVID-19 Across European Nations.” *medRxiv*. [https://doi.org/10.1101/2022.06.16.
766 22276024](https://doi.org/10.1101/2022.06.16.22276024).

767 Srivastava, Ajitesh, Tianjian Xu, and Viktor K Prasanna. 2020. “Fast and Accurate Forecasting of
768 COVID-19 Deaths Using the $SlkJ\alpha$ Model,” July. <https://arxiv.org/abs/2007.05180>.

769 Team, European COVID-19 Forecast Hub. 2021. “Forecasts of New Cases and Deaths
770 Due to Covid-19 over the Next Four Weeks in Countries Across Europe and the UK.”
771 <https://covid19forecasthub.eu/>.

772 ———. 2022. “Covid19-Forecast-Hub-Europe: European Covid-19 Forecast Hub.” [https://github.
773 com/covid19-forecast-hub-europe/covid19-forecast-hub-europe](https://github.com/covid19-forecast-hub-europe/covid19-forecast-hub-europe); Github.

774 Team, Stan Development. 2021. *Stan Modeling Language Users Guide and Reference Manual*, 2.28.1.

775 Ushey, Kevin. 2021. *Renv: Project Environments*. <https://rstudio.github.io/renv/>.

776 Vernon, Ian, Michael Goldstein, and Richard Bower. 2014. “Galaxy Formation: Bayesian History
777 Matching for the Observable Universe.” *Stat. Sci.* 29 (1): 81–90.

778 Williamson, Daniel, Michael Goldstein, Lesley Allison, Adam Blaker, Peter Challenor, Laura Jackson, and
779 Kuniko Yamazaki. 2013. “History Matching for Exploring and Reducing Climate Model Parameter
780 Space Using Observations and a Large Perturbed Physics Ensemble.” *Clim. Dyn.* 41 (7): 1703–29.
781 <https://doi.org/10.1007/s00382-013-1896-4>.

782