# 1 Predictive performance of international COVID-19
# 2 mortality forecasting models

3  Joseph Friedman*, Patrick Liu*, Christopher E. Troeger, Austin Carter, Robert C. Reiner JR, Ryan M.

4  Barber, James Collins, Stephen S. Lim, David M. Pigott, Theo Vos, Simon I. Hay[†], Christopher J.L.

5  Murray[†], Emmanuela Gakidou[†**]

6

7  *These authors contributed equally to the analysis and are listed in alphabetical order.

8  [†]These authors jointly supervised the work.

9  **Correspondence to: Emmanuela Gakidou (gakidou@uw.edu).

10

## 11 Abstract

12  Forecasts and alternative scenarios of COVID-19 mortality have been critical inputs into a range of
13  policies and decision-makers need information about predictive performance. We identified n=386
14  public COVID-19 forecasting models and included n=8 that were global in scope and provided public,
15  date-versioned forecasts. For each, we examined the median absolute percent error (MAPE) compared
16  to subsequently observed mortality trends, stratified by weeks of extrapolation, world region, and
17  month of model estimation. Models were also assessed for ability to predict the timing of peak daily
18  mortality. The MAPE among models released in July rose from 1.8% at one week of extrapolation to
19  24.6% at twelve weeks. The MAPE at six weeks were the highest in Sub-Saharan Africa (34.8%), and the
20  lowest in high-income countries (6.3%). At the global level, several models had about 10% MAPE at six
21  weeks, showing surprisingly good performance despite the complexities of modelling human
22  behavioural responses and government interventions. The framework and publicly available codebase
23  presented here (https://github.com/pyliu47/covidcompare) can be routinely used to compare
24  predictions and evaluate predictive performance in an ongoing fashion.

25  [†]Correspondence to: Emmanuela Gakidou (gakidou@uw.edu).

26  *These authors contributed equally to the analysis and are listed in alphabetical order.

## Introduction

27 Forecasts and alternative scenarios of COVID-19 have been critical inputs into a range of important
28 decisions by healthcare providers, local and national government agencies and international
29 organizations and actors[1–4]. For example, hospitals need to prepare for potential surges in the demand
30 for hospital beds, ICU beds and ventilators[1]. National critical response agencies such as the US Federal
31 Emergency Management Agency have scarce resources including ventilators that can be moved to
32 locations in need with sufficient notice[5,6]. Longer range forecasts are important for decisions such as the
33 potential to open schools, universities and workplaces, and under what circumstances[7]. Much longer-
34 range forecasts—six months to a year—are important for a wide range of policy choices, where efforts
35 to reduce disease transmission must be balanced against economic outcomes such as unemployment
36 and poverty[8]. Furthermore, vaccine and new therapeutic trialists need to select locations that will have
37 sufficient transmission to test new products in the time frame when phase three clinical trials are ready
38 to be launched. Nevertheless, hundreds of forecasting models have been published and/or publicly
39 released, and it is often not immediately clear which models have had the best performance, or are
40 most appropriate for predicting a given aspect of the pandemic.

42 Existing COVID-19 forecasting models differ substantially in methodology, assumptions, range of
43 predictions, and quantities estimated. Furthermore, mortality forecasts for the same location have often
44 differed substantially, in many cases by more than an order of magnitude, even within a six-week
45 forecasting window. The challenge for decision-makers seeking input from models to guide decisions,
46 which can impact many thousands of lives, is therefore not the availability of forecasts, but guidance on
47 which forecasts are likely to be most accurate. Out-of-sample predictive validation—checking how well
48 past versions of forecasting models predict subsequently observed trends—provides insight into future
49 model performance[9]. Although some comparisons have been conducted for models describing the
50 epidemic in the United States[10–13], to our knowledge similar analyses have not been undertaken for
51 models covering multiple countries, despite the growing global impact of COVID-19.

52 This paper introduces a publicly available dataset and evaluation framework
53 (https://github.com/pyliu47/covidcompare) for assessing the predictive validity of COVID-19 mortality
54 forecasts. The framework and associated open-access software can be routinely used to track model
55 performance. This will, overtime, serve as a reference for decision-makers on historical model
56 performance, and provide insight into which models should be considered for critical decisions in the
57 future.

## Results

59 Eight models which fit all inclusion criteria were evaluated (Table 1). These included those modelled by:
60 DELPHI-MIT (Delphi)[14,15], Youyang Gu (YYG)[10], the Los Alamos National Laboratory (LANL)[16], Imperial
61 College London (Imperial)[17],the SIKJ-Alpha model from the USC Data Science Lab (SIKJalpha)[18], and three
62 models produced by the Institute for Health Metrics and Evaluation (IHME)[19] (see methods section for
63 more details). Results are presented in the main text for two main predictive tasks: 1) predicting the
64 magnitude of mortality, and 2) predicting the timing of peak mortality (see methods). Magnitude results
65 are presented in the main text for models that continued to produce forecasts at the time of publication
66 of this article, while peak timing results are presented for models released early enough to capture the
67 first peak in most locations. Results for all historical models are shown in the appendix. Magnitude of

68 mortality results in the main text are presented according to two main analytical approaches. In the
69 "most current" approach, used to select data shown in Figure 3, the most recent 4-week period allowing
70 for the calculation of errors is selected for each extrapolation length. In the "month stratified" approach,
71 used to select data for Figures 4 and 5, models from July were used to calculate errors at each length of
72 extrapolation, with all months shown in the appendix. In each case errors were assessed from one to
73 twelve weeks of forecasting (see methods and Figure 2 for more details).

74 The evaluation framework developed here for assessing how well models predicted the total number of
75 cumulative deaths is shown in Figure 1 for an example country—the United States—and similar figures
76 for all locations included in the study can be found in the appendix. Figure 1, and similar figures in the
77 appendix, also highlight the direction of error for each model in each location. When looking across
78 iterations of forecasts, a wide range of variation can be observed for nearly all of the models.
79 Nevertheless, in many locations, models largely reached consensus regarding trajectories in the summer
80 of 2020. Models diverged again when predicting trajectories for Fall 2020 and Winter 2021, as some
81 models predicted upticks related to seasonality, while others projected continued slow declines in
82 mortality.

83 Figure 2 highlights the most recent errors for each length of extrapolation. For all models, the most
84 recent 1-week errors, reflecting forecasts created in October, ranged from 1% to 2%. The 12-week
85 median absolute percent errors (MAPE), reflecting models produced in July and August, ranged from
86 22.4% for the SIK-J Alpha model, to 79.9% for the Imperial model. At the global level pooling across
87 models, the most recent 6-week MAPE value was 7.2%.

88 Systematic assessments of bias for all models produced in July are shown in Figure 4, and Supplemental
89 Figure 2. The Delphi and LANL models from July underestimated mortality, with median percent errors
90 of -5.6% and -8.3% at 6 weeks respectively, while Imperial tended to overestimate (+47.7%), and the
91 remaining models were relatively unbiased.

92 Overall model performance for models produced in July is shown for cumulative deaths by week in
93 Figure 5. As one might expect, MAPE tends to increase by the number of weeks of extrapolation. Across
94 models released in July the MAPE rose from 1.8% at one week to 24.6% at twelve weeks. Decreases in
95 predictive ability with greater periods of extrapolation were similarly noted for errors in weekly deaths
96 (Supplemental Figure 3). At the global level, MAPE at six weeks was less than 15% for LANL (10.6%),
97 IHME-MS-SEIR (10.6%), SIKJalpha (12.3%) and Delphi (13.6%). The Imperial model had larger errors,
98 about 5-fold higher than other models by six weeks. This appears to be largely driven by the
99 aforementioned tendency to overestimate mortality. At twelve weeks, MAPE values were lowest for the
100 IHME-MS-SEIR (23.7%) model, while the Imperial model had the most elevated MAPE (98.8%). Predictive
101 performance between models was generally similar for median absolute errors (MAEs) (see
102 supplemental figure 4). Global MAE values at 12 weeks, among models released in July varied from 204
103 for the IHME-MS-SEIR model to 1,264 for the Imperial model.

104 Figure 5 also shows that model performance varies substantially by region. The lowest errors across
105 models were observed among high-income countries with a 6-week MAPE values of 6.3%. In contrast,
106 the largest errors were seen in sub-Saharan Africa, with a 6-week MAPE of 34.8%, and Latin America and
107 the Caribbean, with a MAPE of 22.4%. Individual model performance and availability also varied by
108 region.

109    The evaluation framework for exploring the ability of models to predict the timing of peak mortality
110    accurately—a matter of paramount importance for health service planning—is shown in Figure 6 for an
111    example location, Massachusetts. Similar figures for all locations are shown in the appendix. Median
112    absolute errors (MAE) for peak timing also rose with increased forecasting weeks, from 13 days at one
113    week to 30 days at eight weeks (Figure 7). The MAE at eight weeks ranged from 27 days for the IHME
114    Curve Fit and SIKJ-Alpha models to 54 days for the LANL model, with an overall error across models of 30
115    days (Figure 7). Models were generally biased towards predicting peak mortality too early
116    (Supplemental Figure 5).

## Discussion

118    Eight COVID-19 models were identified that covered more than five countries, were regularly updated,
119    publicly released and provide archived results for past forecasts. Taken together at twelve weeks, the
120    models released in July had a median average percent error of 24.6% percent. Errors tend to increase
121    with longer forecasts, rising from 1.8% at one week to 24.6% at 12 weeks. At twelve weeks of
122    extrapolation, the best predictive performance among models considered at the global level was
123    observed for the IHME-MS-SEIR models, with a MAPE of 23.7%, although the best performing model
124    varied by region. The projections provided by Imperial had considerably higher error (98.8%) and the
125    SIKJalpha and Delphi models had intermediate performance for the same period. In the most current
126    models, the 6-week MAPE across models was 7.2%.

127    Although models largely converged in their predictions for the summer of 2020 period, forecasts began
128    to diverge again among predictions for Fall 2020 and Winter 2021. These later divergences are likely due
129    to differences in model assumptions related to the effects of seasonality. Although the top performing
130    models are currently performing in a highly comparable fashion, the updated results presented in this
131    framework in an ongoing fashion may highlight major predictive performance differences as the validity
132    of these assumptions are born out in the coming months.

133    A forecast of the trajectory of the COVID-19 epidemic for a given location depends on three sets of
134    factors: 1) attributes of the virus itself, and characteristics of the location, such as population density
135    and the use of public transport; 2) individual behavioural responses to the pandemic such as avoiding
136    contact with others or wearing a mask; and 3) the actions of governments, such as the imposition of a
137    range of social distancing mandates. Given the complexity of forecasting human and governmental
138    behaviours, especially in the context of a new pandemic, performance of most of the models evaluated
139    here was encouraging. Nevertheless, errors were observed to grow with greater extrapolation time,
140    indicating that governments and planners should recognize the wide uncertainty that comes with longer
141    range forecasts, and strategize accordingly. Hospital administrators may want to hedge on the higher
142    end of the forecast range, while government policymakers may elect to use the mean forecast,
143    depending on their risk tolerance.

144    We also observed substantial differences in average model predictive performance between regions,
145    which can likely be explained by several factors. Data quality has been shown to vary substantially
146    between countries, and many models were initially calibrated on data from early epidemics in China,
147    Europe, and the United States. Furthermore, different regions are at different stages of their epidemic
148    at any given time. For many of the countries in Sub-Saharan Africa for example, the challenge is
149    predicting if, and when, large outbreaks will occur. It is therefore easier for a model to demonstrate

150    large magnitude errors when it predicts a completely different epidemic trajectory. Contrastingly, in
151    some of the more established epidemics, it is easier to predict the nature of more stabilized, ongoing
152    transmission dynamics.

153    We also note that the vast majority of COVID-19 forecasting models did not provide sufficient
154    information to be included in this framework, given that publicly available and date-version forecasts
155    were not made available. We would encourage all research groups forecasting COVID-19 mortality to
156    consider providing historical versions of their models in a public platform for all locations, to facilitate
157    ongoing model comparisons. This will improve reproducibility, the speed of development for modelling
158    science, and the ability of policy makers to discriminate between a burgeoning number of models[20].
159    Many of the models featured in this analysis were generally unbiased, or tended to underestimate
160    future mortality, while other models, such as the Imperial model, as well as many other published
161    models that did not meet our inclusion criteria, tend to substantially overestimate transmission, even
162    within the first four weeks of a forecast. This tendency towards over-estimation among SEIR and other
163    transmission-based models is easy to understand given the potential for the rapid doubling of
164    transmission. Nevertheless, sustained exponential growth in transmission is not often observed, likely
165    due to the behavioural responses of individuals and governments; both react to worsening
166    circumstances in their communities, modifying behaviours and imposing mandates to restrict activities.
167    This endogenous behavioural response is commonly included in economic analyses, however, it has not
168    been routinely featured in transmission dynamics modelling of COVID-19. More explicit modelling of the
169    endogenous response of individuals and governments may improve future model performance for a
170    range of models.

171    Modelling groups are increasingly providing both reference forecasts, describing likely future trends,
172    and alternative scenarios describing the potential effects of policy choices, such as school openings,
173    timing of mandate re-imposition, or planning for hospital surges. For these scenarios, the error in the
174    reference forecast—which we describe in this manuscript—is actually less important than the error in
175    the effect implied by the difference between the reference forecast and policy scenario. Unfortunately,
176    evaluating the accuracy of these counterfactual scenarios is an extremely difficult task. The validity of
177    such claims depends on the supporting evidence for the assumptions about a policy's impact on
178    transmission. The best option for decision-makers is likely to examine the impact of these policies as
179    portrayed by a range of modelling groups, especially those that have historically had reasonable
180    predictive performance in their reference forecasts.

181    Given that a number of very different models demonstrated recent six-week errors for cumulative
182    deaths below 10%, it would likely be worthwhile to construct an ensemble of these models and evaluate
183    the performance the ensemble compared to each component. Although from a logistical standpoint,
184    creating an ensemble of the forecasts would be relatively straightforward, it would be more challenging
185    to integrate such a model pool with scenarios assessing policy options, given that the models have
186    highly different underlying structures. Nevertheless, the inclusion of the models shown here, and future
187    models meeting criteria into an ensemble framework, is an important area for future research.

188    This analysis of the performance of publicly released COVID-19 forecasting models has limitations. First,
189    we have focused only on forecasts of deaths, as they are available for all models included here. Hospital
190    resource use is also of critical importance, however, and deserves future consideration. Nevertheless,
191    this will be complicated by the heterogeneity in hospital data reporting; many jurisdictions report

192   hospital census counts, others report hospital admissions, and still others do not release hospital data
193   on a regular basis. Without a standardized source for these data, assessment of performance can only
194   be undertaken in an *ad hoc* way. Second, many performance metrics exist which could have been
195   computed for this analysis. We have focused on reporting median absolute percent error, as the metric
196   is frequently used, quite stable, and provides an easily interpreted number that can be communicated to
197   a wide audience. Relative error is an exacting standard, however. For example, a forecast of three
198   deaths in a location that observed only one may represent a 200% error, yet it would be of little policy
199   or planning significance. Conversely, focusing on absolute error would create an assessment dominated
200   by a limited number of locations with large epidemics. Future assessment could consider different
201   metrics that may offer new insights, although the relative rank of performance by model is likely to be
202   similar.

203   When taking an inclusive approach to including forecasts from various modelling groups, including
204   estimates from a wide range of time periods and geographies, extra care must be taken to ensure
205   comparability between models. We use various techniques to construct fair companions, such as
206   stratifying by region, month of estimation, and weeks of forecasting, and masking summary statistics
207   representing a small number of values. Nevertheless, other researchers may prefer distinct methods of
208   maximizing comparability over a complex and patchy estimate space. Furthermore, the domains
209   assessed here —magnitude of total mortality and peak timing—are not an exhaustive list of all possible
210   dimensions of model performance. By providing an open-access framework to compile forecasts and
211   calculate errors, other researchers can build on the results presented here to provide additional
212   analyses.

213   COVID-19 mortality forecasts have been used in myriad ways by policymakers as they make difficult
214   decisions about resource management under unprecedented circumstances. Examples include
215   prospectively managing or moving resources between regions such as hospital beds, ICU beds,
216   ventilators, masks and other personal-protective equipment, as well as decisions about social distancing
217   measures, stay-at-home orders, and closing schools, universities and workplaces[1,7]. It is therefore of
218   paramount importance that decision-makers can quickly assess how robust each modelling groups
219   predictions have been historically. Furthermore, we believe a similar approach could be adopted in
220   future pandemics, and for modelling other infectious diseases such as influenza.

221   Ultimately, policymakers would benefit from considering a multitude of forecasting models as they
222   consider resource planning decisions related to the response to the ongoing COVID-19 pandemic. This
223   study provides a publicly available framework and codebase, which will be updated in an ongoing
224   fashion, to continue to monitor model predictions in a timely manner, and contextualize them with prior
225   predictive performance. It is our hope that this spurs conversation and cooperation amongst
226   researchers, which might lead to more accurate predictions, and ultimately aid in the collective
227   response to COVID-19. As the pandemic continues worldwide and resurges in Europe and North America
228   become more evident, regularly updating models, and continually assessing their predictive validity, will
229   be important in order to provide stakeholders with the best tools for COVID-19 decision-making.

230

## Methods

### Systematic Review

A total of 386 published and unpublished COVID-19 forecasting models were reviewed (see appendix). Models were excluded from consideration if they did not 1) produce estimates for at least five different countries, 2) did not extrapolate at least four weeks out from the time of estimation, 3) did not estimate mortality, 4) did not provide downloadable, publicly available results, or 5) did not provide date-versioned sets of previously estimated forecasts, which are required to calculate subsequent out-of-sample predictive validity. Eight models which fit all inclusion criteria were evaluated (Table 1). These included those modelled by: DELPHI-MIT (Delphi)[14,15], Youyang Gu (YYG)[10], the Los Alamos National Laboratory (LANL)[16], Imperial College London (Imperial)[17],the SIKJ-Alpha model from the USC Data Science Lab (SIKJalpha)[18], and three models produced by the Institute for Health Metrics and Evaluation (IHME)[19]. Beginning March 25th, IHME initially produced COVID forecasts using a statistical curve fit model (IHME-CF), which was used through April 29th for publicly released forecasts[1]. On May 4th, IHME switched to using a hybrid model, drawing on a statistical curve fit first stage, followed a second-stage epidemiological model with susceptible, exposed, infectious, recovered compartments (SEIR)[21]. This model—referred to herein as the IHME-CF SEIR model—was used through May 26th. On May 29th, the curve fit stage was replaced by a spline fit to the relationship between log cumulative deaths and log cumulative cases, while the second stage SEIR model remained the same[22]. This model, referred to as the IHME-MS SEIR model, is the basis for recently published work on US State level scenarios of COVID-19 projections in the fall and winter of 2020/2021[23]and was still in use at the time of this publication. The three IHME models rely upon fundamentally different assumptions and core methodologies, and therefore are considered separately. They were also released during different windows of the pandemic, and are therefore compared to models released during similar time periods.

In some cases, numerous scenarios were produced by modelling groups, to describe the potential effects of interventions, or future trajectories under different assumptions. In each case the baseline or status quo scenario was selected to evaluate model performance as that represents the modelers' best estimate about the most probable course of the pandemic. Table 1 summarizes information about each model assumptions, methodologies, input data, modelled outputs, and forecasting range.

### Model Comparison Framework

In order to conduct a systematic comparison of the out-of-sample predictive validity of international COVID-19 forecasting models, a number of issues must be addressed. Looking across models, a high degree of heterogeneity can be observed in numerous dimensions, including sources of input data, frequency of public releases of model estimates, geographies included in the results, and how far into the future predictions are made available for. Differences in each of these areas must be taken into account, in order to provide a fair and relevant comparison.

Input data: A number of sources of input data—describing observed epidemiological trends in COVID-19—exist, and they often do not agree for a given country and time point[24–26]. We chose to use mortality data collected by the Johns Hopkins University Coronavirus Resource Center as the in-sample data against which forecasts were validated at the national level, and data from the New York Times for state-level data for the United States[25,26]. We chose to mainly rely on the Hopkins data as 1) it was the most common input data source used in the different models considered, 2) it covered all countries for

272    which modelling groups produced forecasts, 3) although some quality issues were noted, and managed
273    in our analysis, largely quality was deemed acceptable, and 4) data were made publicly available on a
274    GitHub page and updated daily, which facilitates the maintenance of a timely comparison framework.
275    Locations were excluded from the evaluation (including Ecuador and Peru) where models used
276    alternative data sources, such as excess mortality, in settings with known marked under-registration of
277    COVID-19 deaths and cases[27,28]. We adjusted for differences in model input data using intercept shifts,
278    whereby all models where shifted to perfectly match the in-sample data for the date in which the model
279    was released (see supplemental methods).

280    Frequency of public releases of model estimates: Most forecasting models are updated regularly, but at
281    different intervals, and on different days. Specific days of the week have been associated with a greater
282    number of reported daily deaths. Therefore, previous model comparison efforts in the United States—
283    such as those conducted by the US Centers for Disease Control and Prevention—have required modelers
284    to produce estimates using input data cut-offs from a specific day of the week[29]. For the sake of
285    including all publicly available modelled estimates, we took a more inclusive approach, considering each
286    publicly released iteration of each model. To minimize the effect of day-to-day fluctuations in death
287    reporting, we focus on errors in cumulative and weekly total mortality, which are less sensitive to daily
288    variation.

289    Geographies and time periods included in the results: Each model produces estimates for a different set
290    of national and subnational locations, and extrapolates a variable amount of time from the present.
291    Each model was also first released on a different date, and therefore reflects a different window of the
292    pandemic. Here, we also took an inclusive but stratified approach, and included estimates from all
293    possible locations and time periods. To increase comparability, summary error statistics were stratified
294    by super-region used in the Global Burden of Disease Study[30], weeks of extrapolation, and month of
295    estimation, and we masked summaries reflecting a small number of locations or time points. Models
296    were included in the global predictive validity results only when they were present for all regions.
297    Estimates were included at the national level for all countries, except the United States, where they
298    were also included at the admin-1 (state) level, as they were available for most models. In order to be
299    considered for inclusion, models were required to forecast at least four weeks into the future.

300    Outcomes: Finally, each model also includes different estimated quantities, including daily and
301    cumulative mortality, number of observed or true underlying cases, and various dimensions of hospital
302    resource utilization. The focus of this analysis is on mortality, as it was the most widely reported
303    outcome, and it also has a high degree of societal, epidemiological and public health importance. We did
304    not focus on forecasts of confirmed cases for several reasons. Certain models we wished to include did
305    not provide an estimate of confirmed cases to subsequently assess predictive performance. Mortality,
306    on the other hand, was available for all models. Furthermore, confirmed cases also depend on testing
307    rates, which vary widely over time and across locations. Modelling confirmed cases, therefore,
308    represents different and perhaps larger challenges. Of course, death numbers also have limitations, but
309    they are generally more reliable than case numbers, at least in the early stages of the pandemic, and in
310    locations with limited capacity to test.

311    **Comparison of Cumulative Mortality Forecasts**

312    The total magnitude of COVID-19 deaths is a key measure for monitoring the progression of the
313    pandemic. It represents the most commonly produced outcome of COVID-19 forecasting models, and

314   perhaps the most widely debated measure of performance. The main quantity that is considered is
315   errors in total cumulative deaths—as opposed to other metrics such as weekly or daily deaths—as it has
316   been most commonly discussed measure, to-date, in academic and popular press critiques of COVID-19
317   forecasting models. Nevertheless, alternate measures are presented in the appendix. Errors were
318   assessed for systematic upward or downward bias, and errors for weekly, rather than cumulative
319   deaths, were also assessed. In calculating summary statistics, percent errors were used to control for the
320   large differences in the scale of the epidemic between locations. Medians, rather than means, are
321   calculated due to a small number of large magnitude outliers present in a few time-series. Errors from
322   all models were pooled to calculate overall summary statistics, in order to comment on overarching
323   trends by geography and time.

324   Results are presented using two analytical strategies in the main text. Both strategies are highlighted in
325   Figure 2. The "most current" approach is used to select the data shown in Figure 3. The "month
326   stratified" approach is used for Figures 3 and 4. In the "most current" approach, the most recent 4
327   weeks of model dates are used for each extrapolation length. Therefore, for 1-week errors, models from
328   October were used, whereas for 12-week errors, models from July and August were used. This allows for
329   the assessment of the most recent evidence possible for each set of errors displayed. 4-week periods
330   are used to ensure that the results are not unduly biased by featuring only a small number of runs for
331   each model.

332   In the "month stratified" approach, models from July are used in all cases. This strategy allows for more
333   reliable assessment of certain aspects of predictive validity, as the same models are being compared
334   over time and geographies. For example, the month stratified approach may provide a more comparable
335   assessment of how errors grow with increased length of extrapolation. Models are shown for July in the
336   main text—the most recent month allowing for assessment of errors at twelve weeks of forecasting—
337   and errors stratified for all months are shown in the appendix.

**Comparison of Peak Daily Mortality Forecasts**

339   Each model was also assessed on how well it predicted the timing of peak daily deaths—an additional
340   aspect of COVID-19 epidemiology with acute relevance for resource planning. Peak timing may be better
341   predicted by different models than those best at forecasting the magnitude of mortality, and therefore
342   deserves separate consideration as an outcome of predictive performance. In order to assess peak
343   timing predictive performance, the observed peak of daily deaths in each location was estimated first—
344   a task complicated by the highly volatile nature of reported daily deaths values. Each timeseries of daily
345   deaths was smoothed, and the date of the peak observed in each location, as well as the predicted peak
346   for each iteration of each forecasting model was calculated (see supplemental methods). A LOESS
347   smoother was used, as it was found to be the most robust to daily fluctuations. Results shown here
348   reflect only those locations for which the peak of the epidemic had passed at the time of publication,
349   and for which at least one set of model results was available seven days or more ahead of the peak date.
350   Predictive validity statistics were stratified by the number of weeks in advance of the observed peak that
351   the model was released, as well as the month in which the model was released. Results shown in the
352   main text were pooled across months, as there was little evidence of dramatic differences over time
353   (see appendix). There was insufficient geographic variation to stratify results by regional groupings,
354   although that remains an important topic for further study, which will become feasible as the pandemic
355   peaks in a greater number of countries globally.

## Data and Code Availability

All data and versioned code required to reproduce this analysis its included visualizations are publicly available at ( https://github.com/pyliu47/covidcompare).

## References

1  Team IC-19 health service utilization forecasting, Murray CJ. Forecasting COVID-19 impact on hospital bed-days, ICU-days, ventilator-days and deaths by US state in the next 4 months. *medRxiv* 2020; : 2020.03.27.20043752.

2  Lu FS, Nguyen AT, Link NB, Lipsitch M, Santillana M. Estimating the Early Outbreak Cumulative Incidence of COVID-19 in the United States: Three Complementary Approaches. *medRxiv* 2020; : 2020.04.18.20070821.

3  Weinberger D, Cohen T, Crawford F, *et al.* Estimating the early death toll of COVID-19 in the United States. *medRxiv* 2020; : 2020.04.15.20066431.

4  Epidemic Model Guided Machine Learning for COVID-19 Forecasts in the United States | medRxiv. https://www.medrxiv.org/content/10.1101/2020.05.24.20111989v1 (accessed June 23, 2020).

5  Critical Supply Shortages — The Need for Ventilators and Personal Protective Equipment during the Covid-19 Pandemic | NEJM. *New England Journal of Medicine* http://www.nejm.org/doi/full/10.1056/NEJMp2006141 (accessed July 26, 2020).

6  FEMA Administrator March 27, 2020, letter to Emergency Managers Requesting Action on Critical Steps | FEMA.gov. https://www.fema.gov/news-release/2020/03/27/fema-administrator-march-27-2020-letter-emergency-managers-requesting-action (accessed July 26, 2020).

7  Viner RM, Russell SJ, Croker H, *et al.* School closure and management practices during coronavirus outbreaks including COVID-19: a rapid systematic review. *The Lancet Child & Adolescent Health* 2020; **4**: 397–404.

8  Atkeson A. What Will Be the Economic Impact of COVID-19 in the US? Rough Estimates of Disease Scenarios. National Bureau of Economic Research, 2020 DOI:10.3386/w26867.

9  Tashman LJ. Out-of-sample tests of forecasting accuracy: an analysis and review. *International Journal of Forecasting* 2000; **16**: 437–450.

10 Gu Y. COVID-19 Projections Using Machine Learning. https://covid19-projections.com/ (accessed June 23, 2020).

11 Reich Lab COVID-19 Forecast Hub. https://reichlab.io/covid19-forecast-hub/ (accessed June 23, 2020).

389  12 Project Score Data: COVID-19 Forecasts - Zoltar.
390     https://zoltardata.com/project/44/score_data (accessed June 23, 2020).

391  13 UCLAML Combating COVID-19. http://covid19.uclaml.org/compare (accessed June 23, 2020).

392  14 MIT DELPHI Epidemiological Case Predictions COVIDAnalytics.
393     https://www.covidanalytics.io/projections (accessed June 23, 2020).

394  15 Li ML, Bouardi HT, Lami OS, Trikalinos TA, Trichakis NK, Bertsimas D. Forecasting COVID-19
395     and Analyzing the Effect of Government Interventions. *medRxiv* 2020; :
396     2020.06.23.20138693.

397  16 Los Alamos Natinoal Laboratory COVID-19 Confirmed and Forecasted Case Data.
398     https://covid-19.bsvgateway.org/ (accessed June 23, 2020).

399  17 Imperial College COVID-19 LMIC Reports. https://mrc-ide.github.io/global-lmic-reports/
400     (accessed June 23, 2020).

401  18 Srivastava A, Xu T, Prasanna VK. Fast and Accurate Forecasting of COVID-19 Deaths Using the
402     SIkJ$\alpha$ Model. *arXiv:200705180 [physics, q-bio]* 2020; published online July 12.
403     http://arxiv.org/abs/2007.05180 (accessed Aug 23, 2020).

404  19 COVID-19 estimation updates. Institute for Health Metrics and Evaluation. 2020; published
405     online March 24. http://www.healthdata.org/covid/updates (accessed June 23, 2020).

406  20 Rivers C, George D. How to Forecast Outbreaks and Pandemics. 2020; published online July
407     5. https://www.foreignaffairs.com/articles/united-states/2020-06-29/how-forecast-
408     outbreaks-and-pandemics (accessed July 8, 2020).

409  21 IHME COVID-19 Estimation Update: May 4th, 2020.
410     http://www.healthdata.org/sites/default/files/files/Projects/COVID/Estimation_update_050
411     420.pdf (accessed July 6, 2020).

412  22 IHME COVID-19 Estimation Update: May 29th, 2020.
413     http://www.healthdata.org/sites/default/files/files/Projects/COVID/Estimation_update_05.3
414     0.2020.pdf (accessed July 6, 2020).

415  23 Reiner RC, Barber RM, Collins JK, *et al.* Modeling COVID-19 scenarios for the United States.
416     *Nature Medicine* 2020; : 1–12.

417  24 Coronavirus Pandemic (COVID-19) - Statistics and Research - Our World in Data.
418     https://ourworldindata.org/coronavirus (accessed June 28, 2020).

419  25 nytimes/covid-19-data. The New York Times, 2020 https://github.com/nytimes/covid-19-
420     data (accessed June 28, 2020).

421    26 COVID-19 Map. Johns Hopkins Coronavirus Resource Center.
422        https://coronavirus.jhu.edu/map.html (accessed June 23, 2020).

423    27 Covid-19 data - Tracking covid-19 excess deaths across countries | Graphic detail | The
424        Economist. https://www.economist.com/graphic-detail/2020/07/15/tracking-covid-19-
425        excess-deaths-across-countries (accessed July 26, 2020).

426    28 A greater tragedy than we know: Excess mortality rates suggest that COVID-19 death toll is
427        vastly underestimated in LAC. UNDP.
428        https://www.latinamerica.undp.org/content/rblac/en/home/presscenter/director-s-graph-
429        for-thought/a-greater-tragedy-than-we-know--excess-mortality-rates-suggest-t.html
430        (accessed July 20, 2020).

431    29 CDC. Coronavirus Disease 2019 (COVID-19). Centers for Disease Control and Prevention.
432        2020; published online Feb 11. http://www.cdc.gov/coronavirus/2019-ncov/covid-
433        data/forecasting-us.html (accessed June 23, 2020).

434    30 Dicker D, Nguyen G, Abate D, *et al.* Global, regional, and national age-sex-specific mortality
435        and life expectancy, 1950–2017: a systematic analysis for the Global Burden of Disease Study
436        2017. *The Lancet* 2018; **392**: 1684–735.

437

## Acknowledgements

441

## Competing Interests

443    The authors declare they have no competing interests as defined by Nature Research that might be
444    perceived to influence the results and/or discussion reported in this manuscript.

445

## Author Contributions

447    JF, PL, TV, SIH, CJLM, and EG conceptualized and designed the study, with substantial input from RR, RB,
448    JC, SL, and DP. JF and PL acquired the data, and JF, PL, CT, and AC wrote the analytical code to conduct
449    the analysis. JF, PL, and CJLM drafted the first draft of the article and all authors meaningfully revised.
450    SIH, CJLM, and EG supervised the work.

451

452 **Tables and Figures**

| Model | Data Access | Model Type | Mortality Input Data | Model Outputs | Geographies | Range | Model Structure and Assumptions |
|---|---|---|---|---|---|---|---|
| IHME - CurveFit | http://www.healthdata.org/covid/data-downloads | Statistical Curve Fit | JHU+ local and national governments | Hospital and ICU Admissions, Ventilator, Hospital Beds Utilization; Confirmed Daily and Cumulative Cases; Daily and Cumulative Mortality | 34 Countries* | Aug 4th** | Statistical curve fit model aimed at predicting peak of hospital resource use as a function of social distancing. |
| IHME - CF SEIR | http://www.healthdata.org/covid/data-downloads | Curve Fit + SEIR | JHU+ local and national governments | Hospital and ICU Admissions, Ventilator, Hospital Beds Utilization; Confirmed Daily and Cumulative Cases; Daily and Cumulative Mortality | 52 Countries* | Aug 4th** | Hybrid curve fit (next 8 days) and SEIR model (after 8 days) with additional parameters including mobility, testing, temperature, and population density. |
| IHME – MS SEIR | http://www.healthdata.org/covid/data-downloads | Mortality Spline + SEIR | JHU+ local and national governments | Hospital and ICU Admissions, Ventilator, Hospital Beds Utilization; Confirmed Daily and Cumulative Cases; Daily and Cumulative Mortality | 163 Countries* | Feb 1st | Covariate-adjusted (population, testing, mandates, flu/pneumonia seasonality, mask use, etc.) SEIR model based on daily deaths estimates harmonized with testing, hospitalization via a random knot spline. |
| Youyang Gu | https://github.com/youyanggu/covid19_projections | SEIR | JHU | Daily and Cumulative Mortality; Daily, Active, and Cumulative Cases | 73 Countries* | Nov 1st** | SEIR model with three R0 values corresponding to: 1) Pre-mitigation 2) Post mitigation 3) Post reopening. Performs grid search to optimize parameter selection. |
| MIT - DELPHI | https://github.com/COVIDAnalytics/DELPHI | SEIR | JHU | Cumulative Mortality; Active Cases, Cumulative Detected Cases, Active Hospitalized Cases; Cumulative Hospitalized Cases | 159 Countries* | Jan 15th | Standard SEIR model adjusted for effective meta-analysis driven parameters of contact rate, under-detection, hospitalization, and societal-governmental response measures (4 phased non-linear parametric model). |
| Imperial-LMIC | https://github.com/mrc-ide/global-lmic-reports | SEIR | Euro-CDC | Daily and Cumulative Cases; Daily and Cumulative Mortality; ICU incidence, ICU Demand, Hospital Incidence, Hospital Demand | 176 Countries | Jan 22nd | Modeled using SQUIRE, an age standardized SEIR model with parameters for healthcare capacity and disease severity. Incorporates mobility dependent R0 based on Google mobility data. Baseline scenario assumes current levels of mobility and interventions persist. |
| LANL -GR | https://covid-19.bsvgateway.org/ | Dynamic Growth | JHU | Confirmed Daily and Cumulative Cases; Daily and Cumulative Mortality | 153 Countries* | Dec 15th | Estimates cases driven by a dynamic growth parameter, adjusted based on trends in observed cases. Mortality driven by estimated CFR, assumed to be consistent over the forecast period. |
| USC SIKJalpha | https://github.com/scc-usc/ReCOVER-COVID-19 | SIKJalpha | JHU | Confirmed Daily and Cumulative Cases; Daily and Cumulative Mortality | 177 Countries* | Feb 10th | Application of SIKJalpha epidemiological model which models temporally varying infection rates and human mobility. Models CFR as a function of cases with different infection times. |

**Table 1. Models Included in the Study**
All eight models included in the study are shown. The full list of models assessed for inclusion is shown in the supplemental review file.
Range indicates the last date upon which forecasts are available in the most current version of each model.
*Includes state-level estimates for the United States.
**No longer actively producing forecasts at the time of publication.

**Figure 1. Cumulative Mortality Forecasts and Prediction Errors by Model – Example for United States**

The most recent version of each model is shown on the top left. The middle row shows all iterations of each model as separate lines, with the intensity of color indicating model date (darker models are more recent). The vertical dashed lines indicate the first and last model release date for each model. The bottom row shows all errors calculated at weekly intervals. The top right panel summarizes all observed errors, using median error and median absolute error, by weeks of forecasting, and month of model estimation. Errors incorporate an intercept shift to account for differences in each model's input data. This figure represents an example for the United States of country-specific plots made for all locations examined in this study. Graphs for all geographies can be found in the supplement. Note that while certain model uses different input data source than the other modelling groups causing apparently discordant past trends in the top left panel. We plot raw estimates on the top left panel, however we implement an intercept shift to account for this issue in the calculation of errors.

456



**Figure 2. Illustration of Analytical Framework**

This figure highlights the analytical framework presented in the main text. Part A highlights the "most current" approach, which is used to select the data shown in Figure 3. Part B highlights the "month stratified" approach used for Figures 4 and 5. The Y axis shows the number of weeks of extrapolation for each scenario, while the x axis shows a range of model date—the date on which a model was released. The thick band in each plot highlights the 4-week window of model dates used for each extrapolation week value. The thin line shows the period for which each set of models is extrapolating before errors are calculated. In the top panel, the most recent four weeks of model dates are used for each extrapolation length. Therefore, for 1-week errors models from October were used, whereas for 12-week errors, models from July and August were used. In the bottom panel, models from July are used in all cases. The analytic strategy highlighted in the top panel provides the most recent evidence possible for each extrapolation length. The strategy in the bottom allows for more reliable assessment of how errors grow with increased extrapolation time.

**Figure 3. Most Current - Cumulative Mortality Accuracy – Median Absolute Percent Error**

Median absolute percent error values, a measure of accuracy, were calculated across all observed errors at weekly intervals, for each model by weeks of forecasting and geographic region. Values that represent fewer than five locations are masked due to small sample size. Models were included in the global average when they included at least five locations in each region. Pooled summary statistics reflect values calculated across all errors from all models, in order to comment on aggregate trends by time or geography. Results are shown here for the most recent four week window allowing for the calculation of errors at each point of extrapolation (see Figure 2 and methods). Results from other months are shown in the supplement.

Color scale: −100%  −80%  −60%  −40%  −20%  0%  20%  40%  60%  80%  100%

458

**Global** (Forecasting Weeks vs Model)

| Week | Delphi | LANL | Imperial | SIKJalpha | IHME−MS−SEIR | Pooled |
|---|---|---|---|---|---|---|
| 12 | −18% | | 94% | −5% | −4% | −6% |
| 11 | −15% | | 96% | −3% | −4% | −5% |
| 10 | −12% | | 84% | −2% | −4% | −4% |
| 9 | −11% | | 76% | −2% | −4% | −3% |
| 8 | −9% | | 68% | −1% | −3% | −2% |
| 7 | −8% | | 55% | −1% | −3% | −2% |
| 6 | −6% | −8% | 48% | 0% | −2% | −2% |
| 5 | −4% | −6% | 39% | 0% | −2% | −1% |
| 4 | −2% | −4% | 28% | 0% | −2% | −1% |
| 3 | −1% | −3% | 20% | 0% | −1% | 0% |
| 2 | 0% | −1% | 12% | 0% | −1% | 0% |
| 1 | 0% | 0% | 5% | 0% | 0% | 0% |

**Eastern Europe, Central Asia**

| Week | Delphi | LANL | Imperial | SIKJalpha | IHME−MS−SEIR | Pooled |
|---|---|---|---|---|---|---|
| 12 | −42% | | 51% | −4% | −14% | −14% |
| 11 | −36% | | 50% | 0% | −13% | −12% |
| 10 | −29% | | 47% | 2% | −12% | −11% |
| 9 | −25% | | 45% | 4% | −11% | −7% |
| 8 | −21% | | 41% | 4% | −9% | −5% |
| 7 | −18% | | 36% | 4% | −8% | −3% |
| 6 | −15% | −10% | 27% | 3% | −6% | −3% |
| 5 | −13% | −7% | 23% | 2% | −5% | −2% |
| 4 | −10% | −6% | 16% | 2% | −4% | −2% |
| 3 | −7% | −4% | 12% | 1% | −3% | −1% |
| 2 | −4% | −2% | 8% | 0% | −2% | −1% |
| 1 | −1% | −1% | 4% | 0% | −1% | 0% |

**High−income**

| Week | Delphi | LANL | Imperial | SIKJalpha | IHME−MS−SEIR | Pooled |
|---|---|---|---|---|---|---|
| 12 | −11% | | −5% | −4% | −6% | −6% |
| 11 | −8% | | −4% | −3% | −5% | −5% |
| 10 | −6% | | −2% | −2% | −4% | −4% |
| 9 | −5% | | −1% | −1% | −4% | −3% |
| 8 | −4% | | −1% | −1% | −3% | −2% |
| 7 | −3% | | −1% | 0% | −2% | −2% |
| 6 | −2% | −6% | −1% | 0% | −2% | −1% |
| 5 | −1% | −4% | 0% | 0% | −2% | −1% |
| 4 | −1% | −2% | 0% | 0% | −1% | −1% |
| 3 | 0% | −1% | 0% | 0% | −1% | 0% |
| 2 | 0% | 0% | 0% | 0% | 0% | 0% |
| 1 | 0% | 0% | 0% | 0% | 0% | 0% |

**Latin America and Caribbean**

| Week | Delphi | LANL | Imperial | SIKJalpha | IHME−MS−SEIR | Pooled |
|---|---|---|---|---|---|---|
| 12 | −22% | | 113% | −10% | 3% | −5% |
| 11 | −18% | | 118% | −10% | 0% | −5% |
| 10 | −16% | | 112% | −10% | −1% | −6% |
| 9 | −15% | | 100% | −10% | −4% | −6% |
| 8 | −14% | | 88% | −8% | −4% | −6% |
| 7 | −12% | | 75% | −4% | −5% | −4% |
| 6 | −10% | −19% | 68% | −2% | −7% | −5% |
| 5 | −7% | −15% | 58% | 0% | −4% | −3% |
| 4 | −6% | −12% | 47% | 0% | −3% | −2% |
| 3 | −5% | −8% | 36% | 0% | −2% | −1% |
| 2 | −2% | −5% | 22% | 0% | −2% | 0% |
| 1 | −1% | −1% | 11% | 0% | −1% | 0% |

**North Africa and Middle East**

| Week | Delphi | LANL | Imperial | SIKJalpha | IHME−MS−SEIR | Pooled |
|---|---|---|---|---|---|---|
| 12 | −27% | | 418% | −15% | 17% | −5% |
| 11 | −24% | | 384% | −13% | 13% | −6% |
| 10 | −22% | | 301% | −11% | 10% | −6% |
| 9 | −19% | | 230% | −8% | 9% | −6% |
| 8 | −15% | | 180% | −6% | 8% | −4% |
| 7 | −13% | | 150% | −5% | 6% | −3% |
| 6 | −11% | −7% | 120% | −4% | 3% | −4% |
| 5 | −8% | −5% | 92% | −2% | 1% | −3% |
| 4 | −6% | −4% | 61% | −1% | 0% | −2% |
| 3 | −4% | −3% | 42% | −1% | 0% | −1% |
| 2 | −2% | −2% | 26% | 0% | 0% | 0% |
| 1 | −1% | −1% | 11% | 0% | 0% | 0% |

**Southeast, East Asia, Oceania**

| Week | Delphi | LANL | Imperial | SIKJalpha | IHME−MS−SEIR | Pooled |
|---|---|---|---|---|---|---|
| 12 | 16% | | 33% | −11% | −38% | −9% |
| 11 | 13% | | 33% | −7% | −27% | −3% |
| 10 | 8% | | 33% | −6% | −14% | −2% |
| 9 | 7% | | 22% | −3% | −14% | −2% |
| 8 | 4% | | 56% | −3% | −14% | −2% |
| 7 | 7% | | 72% | −2% | −8% | −1% |
| 6 | 6% | | 67% | −1% | −5% | −1% |
| 5 | 3% | | 52% | −1% | −3% | −1% |
| 4 | 4% | | 39% | −1% | −3% | −1% |
| 3 | 2% | | 24% | 0% | −2% | 0% |
| 2 | 1% | | 12% | 0% | −1% | 0% |
| 1 | 0% | | 3% | 0% | 0% | 0% |

**Sub−Saharan Africa**

| Week | Delphi | LANL | Imperial | SIKJalpha | IHME−MS−SEIR | Pooled |
|---|---|---|---|---|---|---|
| 12 | −30% | | 645% | 5% | 6% | 12% |
| 11 | −26% | | 558% | 6% | 7% | 8% |
| 10 | −25% | | 425% | 5% | 6% | 4% |
| 9 | −27% | | 312% | 4% | 1% | 3% |
| 8 | −28% | | 252% | 2% | −2% | 3% |
| 7 | −22% | | 189% | 1% | −2% | 3% |
| 6 | −20% | −4% | 141% | 1% | −2% | 0% |
| 5 | −17% | −4% | 108% | 0% | −1% | 0% |
| 4 | −14% | −4% | 77% | −1% | −1% | 0% |
| 3 | −11% | −2% | 50% | −1% | −1% | 0% |
| 2 | −3% | −1% | 29% | 0% | 0% | 0% |
| 1 | 0% | 0% | 11% | 0% | 0% | 0% |

**Figure 4. Month Stratified July Models - Cumulative Mortality Bias - Median Percent Error**

Median percent error values, a measure of bias, were calculated across all observed errors at weekly intervals, for each model, by weeks of forecasting and geographic region. Values that represent fewer than five locations are masked due to small sample size. Models were included in the global average when they included at least five locations in each region. Pooled summary statistics reflect values calculated across all errors from all models, in order to comment on aggregate trends by time or geography. Results are shown here for models released in July, and results from other months are shown in the appendix.

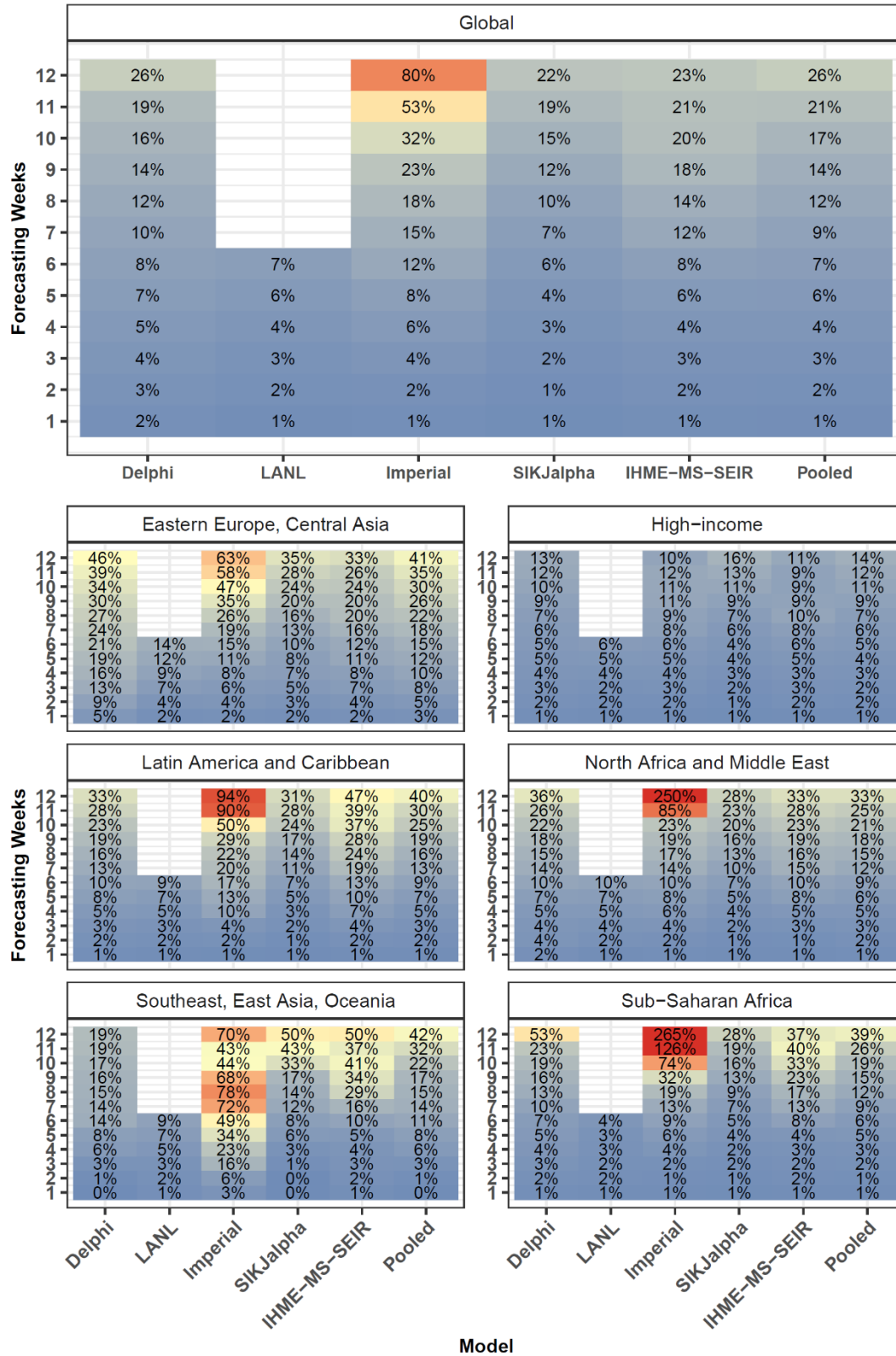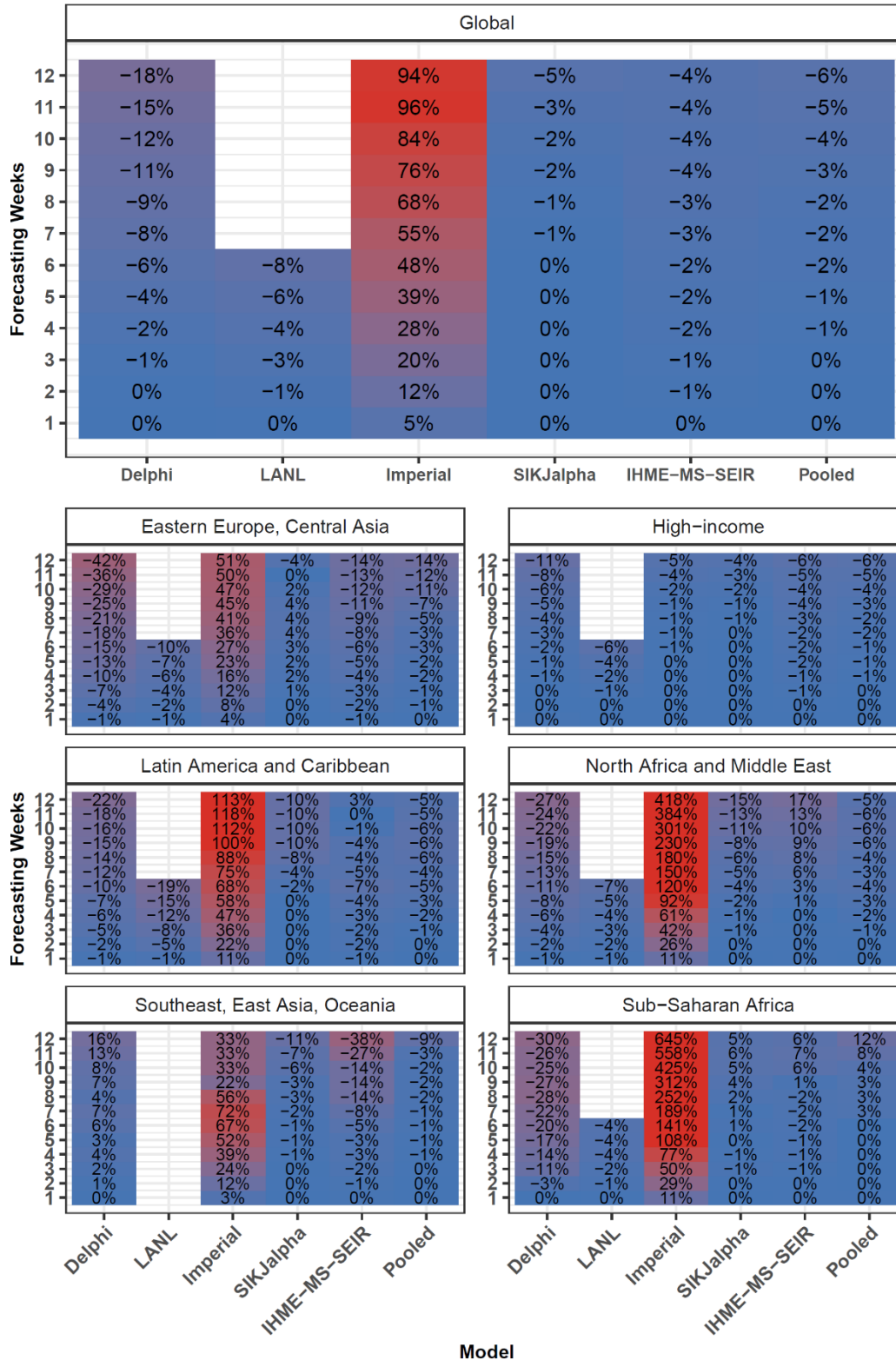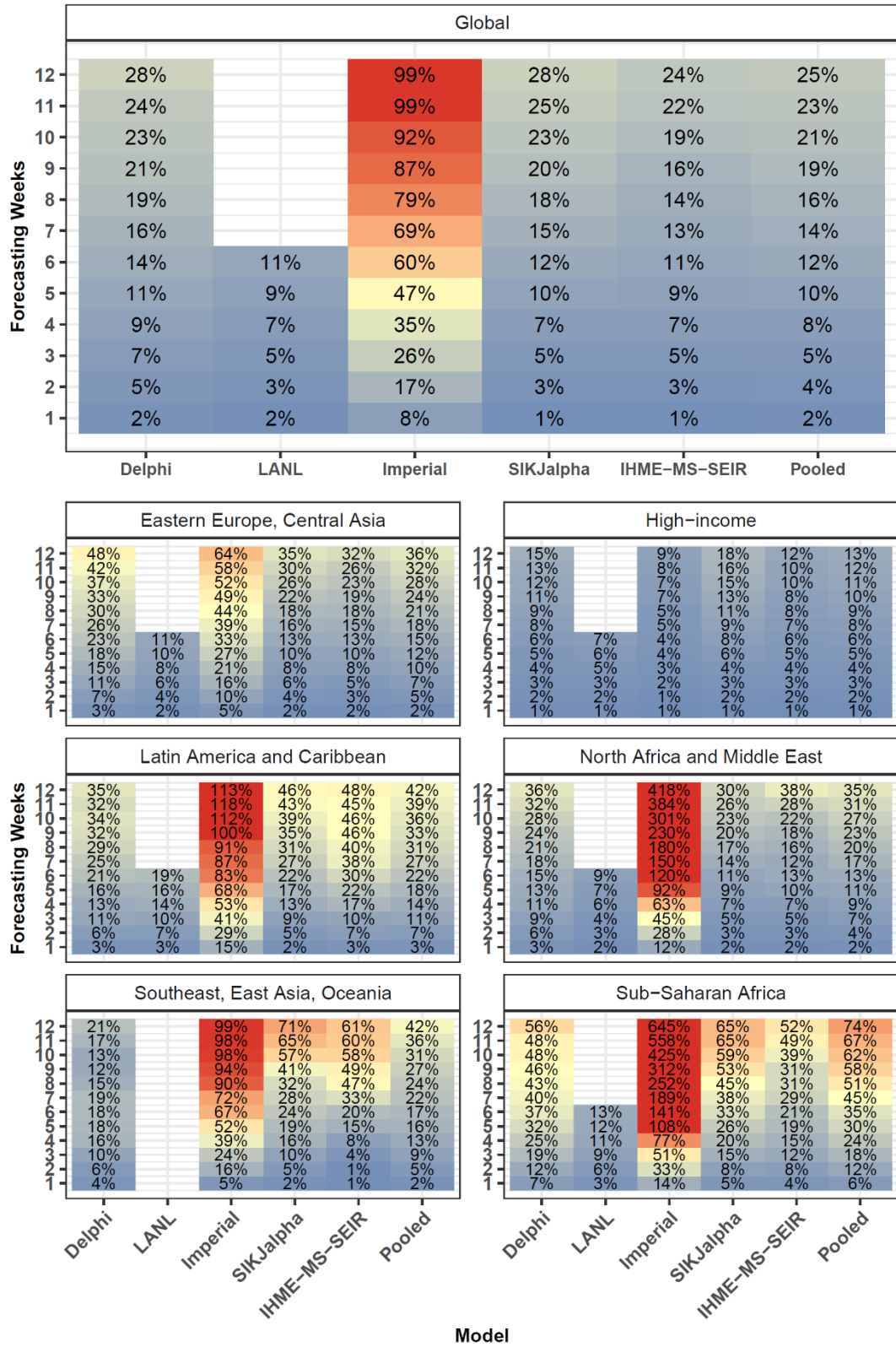**Figure 5. Month Stratified July Models - Cumulative Mortality Accuracy – Median Absolute Percent Error**
Median absolute percent error values, a measure of accuracy, were calculated across all observed errors at weekly intervals, for each model by weeks of forecasting and geographic region. Values that represent fewer than five locations are masked due to small sample size. Models were included in the global average when they included at least five locations in each region. Pooled summary statistics reflect values calculated across all errors from all models, in order to comment on aggregate trends by time or geography. Results are shown here for models released in July, and results from other months are shown in the supplement.
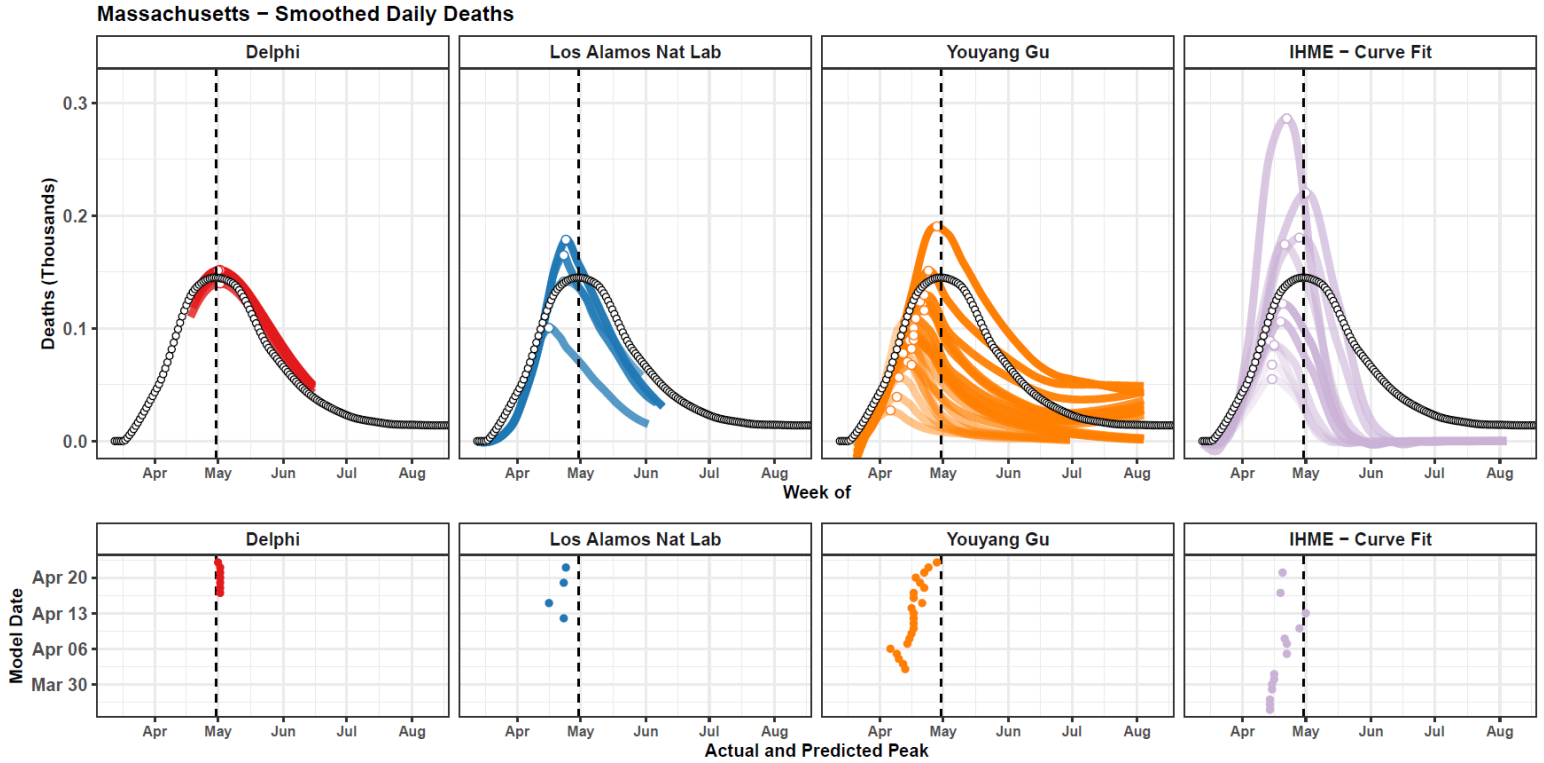
**Figure 6. Observed vs Predicted Peak in Daily Deaths– Example for Massachusetts**

Observed daily deaths, smoothed using a loess smoother, are shown as black-outlined dots (top). The observed peak in daily deaths is shown with a vertical black line (bottom). Each model version that was released at least one week prior to the observed peak is plotted (top) and its estimated peak is shown with a point (top and bottom). Estimated peaks are shown in the bottom panel with respect to their predicted peak date (x-axis) and model date (y-axis). Values are shown for the Massachusetts, and similar graphs for all other locations are available in the appendix. Massachusetts was chosen as the example location as the United States (used as the example for Figure 1) peaked earlier, only allowing for two models to provide peak timing errors, whereas Massachusetts peaked later, allowing for four models, making for a more illustrative example.
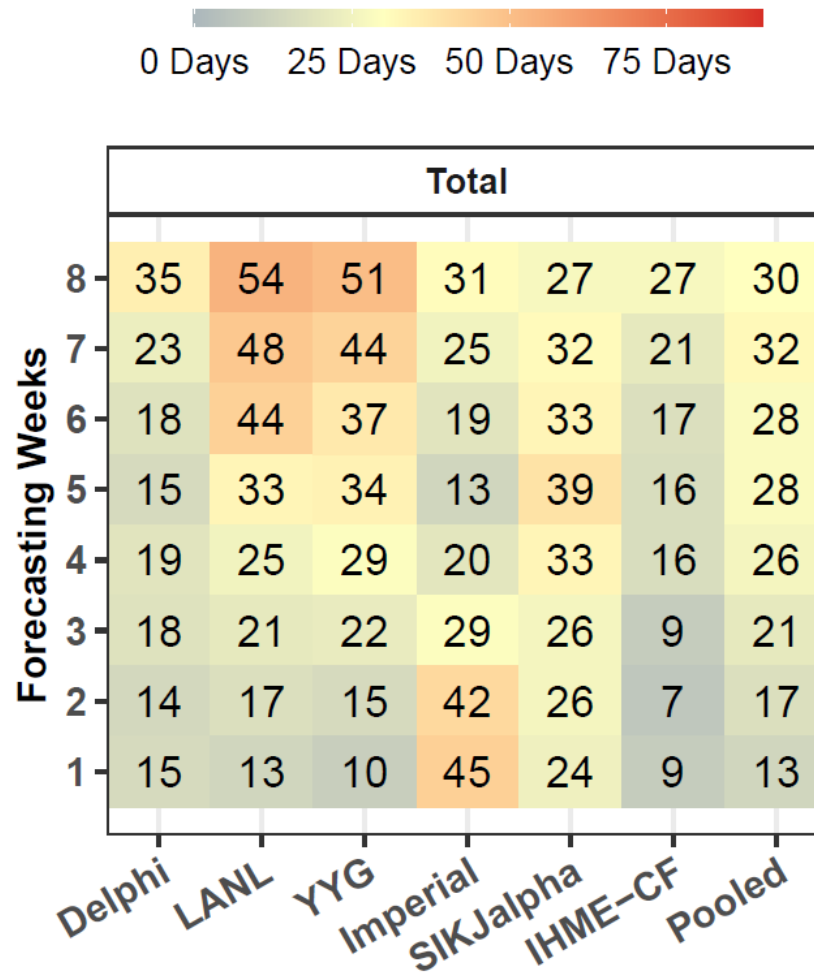
**Figure 7. Peak Timing Accuracy – Median Absolute Error in Days**

Median absolute error in days is shown by model and number of weeks of forecasting. Models that are not available for at least 40 peak timing predictions are not shown. Errors only reflect models released at least seven days before the observed peak in daily mortality. One week of forecasting refers to errors occurring from seven to 13 days in advance of the observed peak, while two weeks refers to those occurring from 14 to 20 days prior, and so on, up to six weeks, which refers to 42-48 days prior. Errors are pooled across month of estimation, as we found little evidence of change in peak timing performance by month (see appendix).