

1 Title: Evaluation of FluSight influenza forecasting in the 2021-22 and 2022-23 seasons with a
2 new target laboratory-confirmed influenza hospitalizations

3 Abstract:

4 Accurate forecasts can enable more effective public health responses during seasonal influenza
5 epidemics. Forecasting teams were asked to provide national and jurisdiction-specific
6 probabilistic predictions of weekly confirmed influenza hospital admissions for one through four
7 weeks ahead for the 2021-22 and 2022-23 influenza seasons.

8 Across both seasons, 26 teams submitted forecasts, with the submitting teams varying between
9 seasons. Forecast skill was evaluated using the Weighted Interval Score (WIS), relative WIS,
10 and coverage.

11 Six out of 23 models outperformed the baseline model across forecast weeks and locations in
12 2021-22 and 12 out of 18 models in 2022-23. Averaging across all forecast targets, the FluSight
13 ensemble was the 2nd most accurate model measured by WIS in 2021-22 and the 5th most
14 accurate in the 2022-23 season. Forecast skill and 95% coverage for the FluSight ensemble
15 and most component models degraded over longer forecast horizons and during periods of
16 rapid change.

17 Current influenza forecasting efforts help inform situational awareness, but research is needed
18 to address limitations, including decreased performance during periods of changing epidemic
19 dynamics.

20

21 **Authors:**

22 Sarabeth M. Mathis^{1†}, Alexander E. Webber^{1†}, Tomás M. León², Erin L. Murray², Monica Sun²,
23 Lauren A. White², Logan C. Brooks^{3,5}, Alden Green³, Addison J. Hu³, Daniel J. McDonald⁴, Roni
24 Rosenfeld³, Dmitry Shemetov³, Ryan J. Tibshirani^{3,5}, Sasikiran Kandula⁶, Sen Pei⁷, Jeffrey
25 Shaman^{7,8}, Rami Yaari⁷, Teresa K. Yamana⁷, Pulak Agarwal⁹, Srikar Balusu⁹, Gautham
26 Gururajan⁹, Harshavardhan Kamarthi⁹, B. Aditya Prakash⁹, Rishi Raman⁹, Alexander Rodríguez¹⁰,
27 Zhiyuan Zhao⁹, Akilan Meiyappan¹¹, Shalina Omar¹¹, Prasith Baccam¹², Heidi L. Gurung¹², Steve
28 A. Stage¹³, Brad T. Suchoski¹², Marco Ajelli¹⁴, Allisandra G. Kummer¹⁴, Maria Litvinova¹⁴, Paulo C.
29 Ventura¹⁴, Spencer Wadsworth¹⁵, Jarad Niemi¹⁵, Erica Carcelen¹⁶, Alison L Hill¹⁶, Sung-mok Jung¹⁷,
30 Joseph C. Lemaitre¹⁷, Justin Lessler¹⁷, Sara L Loo¹⁶, Clifton D. McKee¹⁶, Koji Sato¹⁶, Claire Smith¹⁶,
31 Shaun Truelove¹⁶, Thomas McAndrew¹⁸, Wenxuan Ye¹⁸, Nikos Bosse¹⁹, William S. Hlavacek²⁰,
32 Yen Ting Lin²⁰, Abhishek Mallela²⁰, Ye Chen²¹, Shelby M. Lamm²¹, Jaechoul Lee²¹, Richard G.
33 Posner²¹, Amanda C. Perofsky²², Cécile Viboud²², Leonardo Clemente²³, Fred Lu²³, Austin G
34 Meyer²³, Mauricio Santillana²³, Matteo Chinazzi²³, Jessica T. Davis²³, Kunpeng Mu²³, Ana Pastore
35 y Piontti²³, Alessandro Vespignani²³, Xinyue Xiong²³, Michal Ben-Nun²⁴, Pete Riley²⁴, James
36 Turtle²⁴, Chis Hulme-Lowe²⁵, Shakeel Jessa²⁵, V.P. Nagraj²⁶, Stephen D. Turner²⁶, Desiree
37 Williams²⁶, Avranil Basu²⁷, John M. Drake²⁷, Spencer J. Fox²⁸, Graham C. Gibson²⁰, Ehsan Suez²⁸,
38 Edward W. Thommes^{29,30}, Monica G. Cojocar²⁹, Estee Y. Cramer³¹, Aaron Gerding³¹, Ariane
39 Stark³¹, Evan L. Ray³¹, Nicholas G. Reich³¹, Li Shandross³¹, Nutchcha Wattanachit³¹, Yijin Wang³¹,
40 Martha W. Zorn³¹, Majd Al Aawar³², Ajitesh Srivastava³², Lauren A. Meyers³³, Aniruddha Adiga³⁴,
41 Benjamin Hurt³⁴, Gursharn Kaur³⁴, Bryan L. Lewis³⁴, Madhav Marathe³⁴, Srinivasan

42 Venkatramanan³⁴, Patrick Butler³⁵, Andrew Farabow³⁵, Nikhil Muralidhar³⁶, Naren Ramakrishnan³⁵,
43 Carrie Reed¹, Matthew Biggerstaff¹, Rebecca K. Borchering^{1*}

44 ¹Centers for Disease Control and Prevention, Atlanta, Georgia, 30329, USA.

45 ²California Department of Public Health, Richmond, CA, 95899

46 ³Carnegie Mellon University, Pittsburgh, PA, 15213

47 ⁴University of British Columbia, Vancouver, BC, V6T 1Z4, Canada

48 ⁵University of California, Berkeley, Berkeley, CA 94720

49

50 ⁶Norwegian Institute of Public Health, 0213 Oslo, Norway

51 ⁷Columbia University, New York, NY, 10032

52 ⁸Columbia University School of Climate, New York, NY 10025

53 ⁹Georgia Institute of Technology, Atlanta, GA, 30318

54 ¹⁰University of Michigan, Ann Arbor, MI, 48109

55 ¹¹Guidehouse Advisory and Consulting Services, McClean VA, 22102

56 ¹²IEM, Bel Air, MD, 21015

57 ¹³IEM, Baton Rouge, LA, 70809

58 ¹⁴Indiana University School of Public Health, Bloomington, IN, 47405

59 ¹⁵Iowa State University, Ames, IA, 50011

60 ¹⁶Johns Hopkins University, Baltimore, MD, 21205

61 ¹⁷University of North Carolina at Chapel Hill, Chapel Hill, NC,

62 ¹⁸Lehigh University, Bethlehem, PA, 18015

63 ¹⁹London School of Health and Tropical Medicine, London, UK, WC1E 7HT

64 ²⁰Los Alamos National Laboratory, Los Alamos, NM, 87545

65 ²¹Northern Arizona University, Flagstaff, AZ, 86011

66 ²²Fogarty International Center, National Institutes of Health, Bethesda, MD, 20892

67 ²³Northeastern University, Boston, MA, 02115

68 ²⁴Predictive Science Inc, San Diego, CA 92121

69 ²⁵Signature Science, LLC, Austin, TX, 78759

70 ²⁶Signature Science, LLC, Charlottesville, VA, 22911

71 ²⁷University of Georgia, Athens, GA, 30602

72 ²⁸University of Georgia, Athens, GA, 30609

73 ²⁹University of Guelph, Guelph, ON N1G 2W1, Canada

74 ³⁰Sanofi, Toronto, ON, M2R 3T4

75 ³¹University of Massachusetts Amherst, Amherst, MA, 01003

76 ³²University of Southern California, Los Angeles, CA, 90089

77 ³³University of Texas Austin, Austin, TX, 78712

78 ³⁴University of Virginia, Charlottesville, VA, 22903

79 ³⁵Virginia Tech, Arlington, VA 22203

80 ³⁶ Stevens Institute of Technology, Hoboken, NJ, 07030

81 † Indicates co-first authors

82 * Indicates co-corresponding authors

83 **Introduction**

84 Traditional influenza surveillance systems provide a comprehensive picture of influenza activity
85 in the United States [1, 2, 3] and are fundamental for situational awareness and risk
86 communication. However, they measure influenza activity after it has occurred, and do not
87 directly anticipate future trends to inform risk assessment and healthcare preparedness. To
88 address these limitations, the Centers for Disease Control and Prevention (CDC) has supported
89 open influenza forecasting challenges since the 2013–14 season [4]. This collaborative process
90 (named FluSight) has ensured that forecasting targets are relevant to public health. Additionally,
91 forecast data are openly available, which enables transparent evaluation of forecast
92 performance [5, 6].

93 Originally the FluSight collaboration focused on short-term forecasts of outpatient influenza-like-
94 illness (ILI) rates from ILINet [2] and corresponding results have been summarized previously
95 [4, 5, 6]. However, the COVID-19 pandemic resulted in changes in outpatient care-seeking
96 behavior, and the continued co-circulation of SARS-CoV-2 has further complicated the
97 interpretation of ILI data. In the 2021–22 influenza season, the FluSight forecast target shifted to
98 the weekly number of hospital patients admitted with laboratory-confirmed influenza from the
99 Health and Human Services (HHS) Patient Impact and Hospital Capacity Data System [7]. This
100 system was created during the COVID-19 pandemic to gather a complete and unified
101 representation of COVID-19 disease outcomes along with other metrics related to health care
102 capacity. Hospitals registered with Centers for Medicare and Medicaid Services (CMS) are
103 required to report daily COVID-19 and influenza information [8]. Reporting of the influenza data
104 elements, including the previous day's number of admissions with laboratory-confirmed
105 influenza virus infection, became mandatory on February 2, 2022, [8].

106 The COVID-19 pandemic disrupted the typical timing, intensity, and duration of seasonal
107 influenza activity in the United States and many parts of the world [9, 10]. Influenza activity was
108 very low during the 2020–21 season in the U.S., but activity increased during the 2021–22
109 season, with activity peaking later in April, May, and early June 2022 and remaining at higher
110 levels than had been reported during these months in previous seasons [10]. In the 2022-23
111 influenza season, activity began increasing nationally in early October, earlier than previous
112 seasons [2,3,11], and peaked in early December 2022. In this analysis, we summarize the
113 accuracy and reliability of ensemble and component 1- to 4-week ahead forecasts submitted in
114 real-time during the 2021–22 and 2022-23 influenza seasons and identify areas for forecast
115 improvement.

116

117 **Methods**

118 Forecasts of weekly influenza hospital admissions were openly solicited from existing COVID-19
119 and influenza forecasting networks every Monday from January 10, 2022, through June 20,
120 2022, for the 2021-22 season. For the 2022-23 season, forecasts were solicited every Monday
121 from October 17, 2022, through January 9, 2023, then every Tuesday from January 17, 2023,
122 through May 17, 2023. Weeks were defined in terms of MMWR Epiweeks (EW) spanning
123 Sunday to Saturday [12]. Forecasted jurisdictions included the U.S. national level, all fifty states,
124 Washington D.C., and Puerto Rico. Forecasts for the Virgin Islands, while requested, were not
125 included in this evaluation due to low reported hospitalization counts and irregular data
126 submission. Each week, forecasting teams were asked to provide jurisdiction-specific point
127 estimates and probabilistic predictions for 1-, 2-, 3-, and 4-week ahead weekly counts of
128 confirmed influenza hospital admissions. A total of 23 quantiles were requested for the
129 probabilistic forecasts: 0.010, 0.025, 0.050, 0.100, 0.150, ..., 0.950, 0.975, and 0.990. Teams
130 were not required to submit forecasts for all four weeks ahead or for all locations. Additional
131 details of the forecast submission process (e.g., file formatting, submission procedures, and
132 required metadata) are provided in the FluSight-forecast-data GitHub Repository [13].
133

134 The FluSight Ensemble model was generated for all forecasted jurisdictions each week using
135 the unweighted median of each quantile among eligible forecasts. Forecasts were considered
136 eligible for inclusion in the ensemble if they were submitted by 11:59 PM ET on the due date
137 and if all requested quantiles were provided. Modeling teams could further designate whether a
138 particular model's forecasts should be included in the ensemble. If a forecast was designated as
139 "other", it was not included in the FluSight ensemble and not evaluated in this manuscript.
140

141 Baseline forecasts and their prediction intervals were generated each week using the `simplets R`
142 package [14] based on the incident hospitalizations reported in the previous week. The median
143 prediction of the baseline forecasts is the corresponding target value observed in the previous
144 week, and noise around the median prediction is generated using positive and negative 1-week
145 differences (i.e., differences between consecutive reports) for all prior observations, separately
146 for each jurisdiction. Sampling distributions were truncated to prevent negative values. The
147 same median prediction is used for the 1-through 4-week ahead forecasts. Further details on
148 the generation of the baseline model's prediction intervals from a smoothed version of this
149 distribution of differences have been described previously [15,16].
150

151 For inclusion in this analysis, forecasting teams must have submitted greater than or equal to
152 75% of the requested targets between the forecast evaluation period of February 21, 2022, to
153 June 20, 2022 (total of 18 weeks) for 2021-22 or October 17, 2022, to May 15, 2023 (total of 30
154 weeks) for 2022-23. These periods translate to 4-week ahead forecast target end dates of
155 March 19, 2022, to July 16, 2022 for the 2021-22 season and November 11, 2022, to June 10,
156 2023 for the 2022-23 season. The start date of the evaluation period for the 2021-22 season
157 was chosen to be the first forecast date following two weeks of mandatory reporting of
158 confirmed influenza hospitalizations [8] to minimize potential effects of reporting changes on
159 forecasts. For 2021-22 and 2022-23, three and 12 models were excluded from the primary
160 analysis, respectively, for not meeting the inclusion criteria.
161

162 Forecasts were evaluated against the reported number of the previous day's laboratory
163 confirmed influenza admissions (Field #34) from the COVID-19 Reported Patient Impact and
164 Hospital Capacity by State Timeseries [17], with data shifted one day earlier to align with
165 admission date and then aggregated to the weekly scale (from Sunday to Saturday) [13], using
166 data as of September 12, 2022, for 2021-22 and June 13, 2023, for 2022-23. This dataset is

167 subject to revision by submitting facilities; therefore, we analyzed backfill and revision for each
168 season (Supplemental Analysis 1). For each of the contributed forecasts included in the
169 analysis, values were rounded to more closely relate the values of prediction intervals of
170 forecasts to the reported numbers of hospital admissions. In particular, forecast values for
171 quantiles less than 0.5 were rounded down, values for quantiles greater than 0.5 were rounded
172 up, and values for the 0.5 quantile were rounded normally. This rounding procedure ensured
173 that teams were not penalized for missing the prediction interval by less than one hospital
174 admission.

175
176 To evaluate forecast performance across all states, D.C., and Puerto Rico, we primarily used
177 the Weighted Interval Score (WIS). The WIS is a proper score that generates interval scores for
178 probabilistic forecasts provided in the quantile format [15,18]. Briefly, interval scores are used to
179 account for dispersion, underprediction, and overprediction. Forecasts with lower absolute WIS
180 values are considered more accurate than forecasts with higher absolute WIS values. The
181 relative WIS compares forecast WIS values from those of the baseline model. Simple means
182 were calculated for absolute and relative WIS to get a score for each model, location, and
183 season. Mean absolute error (MAE) values are also considered for characterizing differences
184 between forecasted and reported weekly hospitalizations [15]. Unless otherwise specified,
185 forecasts of national hospitalizations were not included in summary metrics for accuracy (e.g.,
186 absolute WIS) since these forecasts can have a disproportionate impact on the overall score. To
187 address concerns related to assessing measures of absolute error on a natural scale when
188 forecasts span multiple orders of magnitude [19], we performed an analogous analysis on log-
189 transformed hospitalization counts after adding one to all counts to account for zero counts
190 (Supplemental Analysis 2). We also performed a separate analysis including only national
191 forecasts (Supplemental Analysis 3).

192
193 In addition, we considered coverage values of the quantile-based prediction intervals to assess
194 each model's ability to appropriately capture uncertainty in forecasts. Coverage values are
195 defined as the percent of observed values that fall within the 50% or 95% prediction intervals for
196 the corresponding date. Ideally, the percent coverage values will be equal to the corresponding
197 prediction interval, e.g., 95% percent prediction intervals should contain the reported value 95%
198 of the time.

199
200 Comparing model forecasts is complicated by the fact that not all models submit forecasts for
201 each of the forecast targets and for each forecast week in the evaluation period. To partially
202 account for this, we consider the percent of forecasts submitted as an indicator of how often and
203 how many different types of forecasts were submitted by each team. Following Cramer et al.
204 2022 [15], we also consider a standardized rank score that uses the number of models
205 forecasting a particular location and target and then ranks these forecasts. Ranks were
206 determined by relative WIS performance, with the best performing model for each observation
207 being assigned a rank of 1 and the worst performing model receiving a rank equal to the
208 number of models submitting a forecast for the observation. These ranks were standardized by
209 rescaling so that 0 corresponds to the worst rank and 1 corresponds to the best rank.

210
211

212 **Results**

213 The 2021-22 influenza season was characterized by two distinct waves of activity. The first
214 occurred between November 2021 and January 2022 and the second between February and
215 June 2022, though reporting of influenza hospitalizations was not mandatory in the HHS system
216 until February 2, 2022 (see observed data in Figure 1a). Reported national weekly influenza
217 hospital admissions exceeded 1000 for 22 out of 25 of the forecast weeks (Figure 1a). Updates
218 to weekly counts from the forecast evaluation period were generally minimal (Figures S2 – S4),
219 with 94% of updates during the 2021-22 season resulting in changes of under 10
220 hospitalizations for subnational jurisdictions.

221 The 2022-23 influenza season was characterized by an early start, reaching 1000 hospital
222 admissions nationally before October 2022. A sharp increase nationally through October and
223 November led to a peak of 26,600 hospital admissions in early December. Hospital admissions
224 decreased rapidly after December, with 3,000 weekly hospital admissions by the end of
225 January, and eventually dropped below 1000 confirmed weekly admissions nationally by May
226 2023. Weekly numbers of admissions exceeded 1000 for 27 out of 34 of the forecast weeks
227 (Figure 1b, Figure S4). In the 2022-23 season, 83% of updates for weekly admissions resulted
228 in changes of under 10 hospitalizations for subnational jurisdictions.

229 **Models Included**

230 For both the 2021-22 and 2022-23 influenza seasons, 26 modeling teams submitted forecasts
231 and 21 and 16 respectively, were eligible for end-of-season evaluation, not including the
232 FluSight baseline and ensemble models. The number and types of models submitted varied
233 across weeks with a range of methodological approaches (see Table S1). For the 2021-22
234 season, a median of 21 models were submitted (range: 15-22), with most having a statistical
235 component, three mechanistic, and six ensembles of component models. In 2022-23 there was
236 a median of 20 models (range: 15 to 26) submitted each week, with many having a statistical
237 component, three mechanistic, and four ensemble models. Modeling teams varied across
238 seasons, with 13 modeling groups having submitted eligible forecasts for both seasons. When
239 only national forecasting targets were considered, no additional teams were included for the
240 2021-22 season, but two teams, NIH-Flu_ARIMA and ISU_NiemiLab-Flu met inclusion criteria
241 for 2022-23 (Supplemental Analysis 3).

242 **Relative WIS**

243 Over the evaluation period more models outperformed the FluSight baseline model in 2022-23
244 (12) than in 2021-22 (6) based on relative WIS (Table 1). Within each season, the models that
245 achieved an overall relative WIS less than or equal to one represent a variety of modeling
246 strategies, including a basic quantile autoregression fit, a mechanistic compartmental model
247 with stochastic simulations, an ensemble of time-series baseline models, a random walk model,
248 a random forest ensemble, and the FluSight hub ensemble (Table S1). Similar results were
249 observed when models were evaluated based on their point forecasts alone (see MAE
250 estimates in Table 1).

251
252 Few teams outperformed the FluSight Ensemble in relative WIS for both seasons. The CMU-
253 TimeSeries model was the only model that outperformed the ensemble for both the 2021-22
254 and 2022-23 seasons while the MOBS-GLEAM_FLUH, PSI-DICE and MIGHTE-Nsemble
255 models outperformed the ensemble only in the 2022-23 season.

256

257 For both seasons, forecasts from the FluSight Ensemble were ranked among the top 50% of all
258 model forecasts for the same location, date, and target, more than three-fourths of the time
259 (79.89% in 2021-22 and 79.02% in 2022-23) (Figure 2). Three models consistently ranked in the
260 top 25% for 2021-22 and 2022-23 seasons, respectively: CMU-TimeSeries (42.49%, 36.32%),
261 PSI-DICE (39.24%, 39.84%), and MOBS-GLEAM_FLUH (38.89%, 50.31%). Several models,
262 ten in 2021-22 and eleven in 2022-23, had bimodal rank distributions, with a combined majority
263 of their forecasts falling in either the bottom 25% or top 25% (Figure 2).

264
265

Log-Transformed Analysis

266 For both seasons, the analysis using log-transformed hospitalization counts resulted in the
267 same top five performing teams in terms of absolute and relative WIS. For the 2021-22 season,
268 all teams were ranked the same for the log-transformed and non-transformed analyses. In
269 2022-23, MIGHTE-Nsemble and PSI-DICE performed better than CMU-TimeSeries for the log-
270 transformed analysis (Table 1 and Supplemental Analysis 2).

271

Relative WIS and Spatial Variation

273 Model performance varied by spatial jurisdiction. For individual states, relative WIS values
274 varied across models ranging from 0.46 to 12.46 in 2021-22 and 0.31 to 12.28 in 2022-23
275 (Figure 3). More models, including the ensemble, performed better at the state-level than the
276 baseline in 2022-23 compared to 2021-22. The relative WIS of the FluSight Ensemble had the
277 smallest range of values across all locations from 0.58 to 1.06 in 2021-22 to 0.63 to 1 in 2022-
278 23 (Figure 3 and Figure S1). To further examine forecast performance across jurisdictions, we
279 considered the percent of jurisdictions that the relative WIS value for a given model and location
280 pair was less than the baseline (i.e., lower than 1). The FluSight Ensemble performed as well as
281 or better than the baseline for all forecast jurisdictions for 2022-23 and 51 out of 52 forecast
282 jurisdictions for 2021-22, a larger number of jurisdictions than any submitted model (Figure 3).
283 In 2022-23, 12 models performed better than the baseline at the jurisdiction-level at least 50%
284 of the time, compared to six models in 2021-22. In general, the models with lower (better)
285 relative WIS values were consistent between the analysis with all spatial jurisdictions and the
286 analysis considering only national forecast targets for both seasons (Supplemental Analysis 3).

287

288

Absolute WIS

290 Across forecasted weeks, the FluSight Ensemble's worst performance in terms of absolute WIS
291 (maximum values) for 1-week ahead targets on March 19, 2022 for 2021-22 and on November
292 26, 2022 for 2022-23 (Figure 4a). For the 4-week ahead horizon, maximum absolute values,
293 indicating the worst performance, for each season occurred on June 04, 2022, and December
294 03, 2022, respectively (Figure 4a). Minimum, or best, absolute WIS values for each season
295 occurred on July 16, 2022, and May 13, 2023, respectively, both during periods of low flu
296 activity.

297

298

Coverage

300 Model performance for the FluSight Ensemble dropped during periods of relatively rapid change
301 (see Figures 1 and 3). The lowest 1-week horizon 95% value occurred for forecasts with target
302 end dates of March 14, 2022, for 2021-22 and on November 21, 2022, for 2022-23 (Figure 4b,
303 c). Across forecasted weeks in the 2021-22 season, the FluSight Ensemble had a minimum
304 95% coverage value at the 1-week horizon of 75%. Lower 95% coverage for the 1-week horizon

305 was observed in the 2022-23 season with a minimum of 29%. The maximum coverage rate
306 achieved by the FluSight Ensemble in any individual week was 100% in both seasons. Minimum
307 FluSight Ensemble 95% coverage values for forecasts at the 4-week horizon in any individual
308 week were 62% for 2021-22 and 15% for 2022-23.

309
310 Model performance, in terms of coverage, tended to decline at longer time horizons for the
311 FluSight Ensemble, baseline, and individual contributed models (see Table 2). Over the forecast
312 weeks, the 2021-22 FluSight ensemble had slightly higher overall 95% coverage values of
313 89.32%, 86.11%, 85.15%, and 83.33% for the 1- to 4-week ahead horizons respectively,
314 compared to the 2022-23 season during which the FluSight Ensemble had 95% coverage
315 values of 85.79%, 81.64%, 78.78%, and 77.85% for the 1- to 4-week ahead horizons
316 respectively. A similar proportion of models had higher overall 95% coverage values at the 1-
317 week ahead horizon than at the 4-week ahead horizon for 2022-23 (14 of 18 models) and 2021-
318 22 (18 out of 23 models) (Table 2). Out of the forecast targets and across forecast weeks, the
319 FluSight Ensemble's 95% prediction interval contained at least 90% of the corresponding
320 observed values only 55.56% and 64.52% of the time, for 2021-22 and 2022-23 respectively
321 (Table 2). Ideally 95% prediction intervals are just wide enough to capture 95% of eventually
322 observed values.

323

324

325 **Discussion**

326 The 2021-22 influenza season marked the return of from very low levels of seasonal influenza
327 activity observed in the U.S. following the first years of the COVID-19 pandemic, and many
328 components of the 2021-22 and 2022-23 FluSight Forecasting Challenges were new. One of
329 the most substantial changes was the shift from the original FluSight forecasting targets of
330 weekly influenza-like-illness (ILI) percentages to weekly counts of confirmed influenza
331 hospitalizations. The COVID-19 pandemic resulted in the availability of a new data source, the
332 unified HHS-Protect dataset [17], which provided information on laboratory confirmed daily
333 influenza hospitalizations from all 50 states, D.C., and Puerto Rico. Confirmed influenza hospital
334 admissions may more directly inform influenza preparedness and response efforts. During the
335 time period that these forecasting results cover, data were reported daily, with mandatory
336 reporting for influenza admissions from most hospitals in each state, U.S. territories, and D.C
337 starting February 2, 2022. Despite challenges accompanying the shift to the new target of
338 influenza hospitalizations, such as limited historic data from this system for model training, these
339 forecasts provided substantial utility and reinforced a number of lessons learned over the course
340 of previous forecasting activities, both during the pre-pandemic influenza seasons and the
341 COVID-19 pandemic.

342

343 **Forecast performance - accuracy**

344 As demonstrated in this analysis, collaborative forecasting hub approaches provide
345 opportunities to systematically evaluate performance across multiple modeling strategies and
346 enable the creation of ensemble models. Since a particular model's performance often varies
347 within and across seasons [20], it is helpful to have a unified representation of model inputs that
348 can be used to quickly assess expected upcoming trends. Additionally, this work indicates that
349 ensemble models may also provide more consistently reliable and well-calibrated forecasts
350 across spatial jurisdictions.

351

352 Across the evaluation period for both seasons and all forecast jurisdictions, the FluSight
353 ensemble was among the top 5 performing models in terms of Absolute WIS and Relative WIS.

354 Additionally, when considering forecast performance by rank (Figure 2), the FluSight ensemble
355 more accurately predicted weekly influenza hospital admissions than most contributed models
356 with the majority of the FluSight ensemble forecasts falling within the top 50% of submitted
357 forecasts (Table 1, Figure 2). While the PSI-DICE, CMU-TimeSeries, and MOBS-
358 GLEAM_FLUH models have more forecasts in the top 25%, they exhibit higher spatial
359 heterogeneity than the FluSight ensemble in forecast performance (Figure 3). The generally
360 high accuracy of the FluSight Ensemble relative to that of individual models is consistent with
361 previous findings that ensemble models, that utilize the outputs from multiple teams, generally
362 outperform individual models on average [15,21,22,23]. Like most models, ensembles may have
363 decreased performance during periods of rapid change when some individual models may have
364 higher accuracy; however, identifying these time frames and corresponding high-performing
365 models has been difficult a priori [5,6].

366
367 One option to better evaluate forecast performance during periods of change and across
368 multiple magnitudes is to evaluate transformed counts [19]. We did not find notable differences
369 in model performance using this approach in either season. We expected that there might be a
370 stronger influence on performance in the 2022-2023 season which saw a sharp increase in
371 hospitalizations in fall 2022, but it is possible that models were not able to capture this initial rise
372 and thus did not accrue additional benefit in the log transform score. The long tail of the season
373 may also have elevated scores across all models.

374
375 Forecast model performance tended to decline over longer time horizons. For both the 2021-22
376 and 2022-23 FluSight seasons, accuracy declined across the 1- to 4-week ahead horizons. This
377 trend has been observed previously in multiple forecast activities. The U.S. COVID-19 Forecast
378 Hub observed declines in accuracy for forecasted deaths over periods of 1- to 4- weeks ahead,
379 and German and Polish COVID-19 forecast efforts also showed declines in performance at the
380 3- and 4-week ahead horizons [18]. Accuracy scores were also shown to decline over longer
381 time horizons for influenza-like-illness forecasts [20].

382
383 Across the forecast weeks, individual models often showed larger increases in absolute WIS,
384 while the FluSight ensemble had the smallest range of absolute WIS for each season,
385 demonstrating one aspect of stability for the FluSight ensemble. In terms of state-level
386 performance, the FluSight ensemble tended to be more robust than individual models, as
387 measured by relative WIS scores (Figure 3). Similarly, the COVID-19 Forecast Hub ensemble
388 [15] performed better across all locations, with the COVID-19 Hub ensemble being the only
389 model to outperform the baseline in each of the forecast locations [15].

390
391
392
393
394

Forecast performance – coverage

395 Our analysis found that, as the forecast horizon moved from 1- to 4-weeks, the FluSight
396 ensemble 95% prediction interval coverage declined from 89.61% to 83.74% in 2021-22 and
397 from 85.69% to 77.85% in 2022-23. These results highlight room for improvement in model
398 calibration, as almost all models (with the exception of the UMass trends ensemble) were
399 overconfident in their predictions (Table 2). The lack of comparable historical data for model
400 fitting may have contributed to poor calibration of 95% prediction intervals.

401
402 Consistent with past forecasting efforts, forecasting remains difficult in periods of rapid change
403 and epidemic turning points (e.g., during initial increases or periods of peaking activity). This

404 analysis highlights declines in forecast accuracy and coverage during periods of rapid change in
405 influenza hospitalizations during both the 2021-22 and 2022-23 seasons. For example, the only
406 model that had 95% coverage greater than 80% from October to January 2023 when
407 hospitalizations were rapidly increasing and then peaking was LUCompUncertLab-
408 humanjudgment, which did not end up meeting inclusion criteria for the full season analysis.
409 Analogous declines were also observed for COVID-19 case forecasts [24] and mortality
410 forecasts across different waves of the COVID-19 pandemic [15], where forecasts
411 systematically underpredicted during periods of increase and overpredicted during periods of
412 decrease.

413
414 Times of changing dynamics are the most important periods for public health response and
415 communication. While forecasting the magnitude at these times may be less tractable, it is
416 possible that we may be able to provide more reliable information during these difficult
417 forecasting periods so that forecasts are better able to inform critical planning. In general, most
418 ensembles tend to predict less activity than observed when trends are steeply increasing and
419 predict more activity than observed when trends are steeply decreasing, especially when there
420 is between- or within-model uncertainty in the timing of peaks in cases, hospitalizations, or
421 deaths. Thus, it may be possible that an ensemble of forecasts for categorical increases or
422 decreases in activity [25] may have additional utility in terms of preserving valuable information
423 while also maintaining the benefits of the use of ensembles over individual models. As such, the
424 FluSight Forecasting Hub added an experimental target in the 2022-23 season for forecasting
425 categorical rate changes in influenza hospitalizations (e.g., probabilities of increase or
426 decrease) [13]. Assessing the utility of this additional forecast target will be an important area of
427 investigation moving forward. Aside from soliciting a separate forecasting target, it may be
428 possible to determine which forecasting models perform better during different phases of
429 epidemics and then use this information to weight models accordingly when their forecasts are
430 aggregated into an ensemble [26].

431

432 **Influenza forecasting in the COVID-19 era: challenges and opportunities**

433 Several challenges for forecasting existed during the 2021-22 and 2022-23 influenza seasons.
434 First, as noted earlier, the change in the forecasting target from outpatient ILI percentages to
435 counts of influenza-associated hospitalizations from a data collection system established during
436 the COVID-19 pandemic meant that there was little data for forecast calibration and training.
437 This shift also required changes in data processing for teams that had produced ILI forecasts
438 previously. While previous data on influenza-associated hospitalizations was available through
439 the FluSurv-NET system, differences in reporting and the spatial resolution, of the FluSurv-NET
440 system may have complicated the process of utilizing this dataset for the purpose of forecasting
441 model calibration. In addition, reporting within the unified HHS-Protect hospitalization dataset
442 changed throughout this forecasting endeavor. For example, the confirmed influenza hospital
443 admissions field only became mandatory for the 2021–22 season on February 2, 2022, leading
444 to an increase in the number of reported hospitalizations and a change in hospital reporting
445 practices during a period of increasing influenza activity.

446

447 In addition to changing reporting patterns, the COVID-19 pandemic brought other challenges for
448 forecasting influenza, including changing human behavior. The quantity and types of
449 interactions between people likely changed in tandem with perceptions of risk of illness with
450 COVID-19. In addition, the use of nonpharmaceutical interventions (NPIs) aimed at preventing
451 SARS-CoV-2 transmission (e.g., stay-at-home orders, mask wearing) reduced transmission of
452 other respiratory pathogens [9], including influenza. These changes in behavior may be related

453 to the minimal influenza activity observed in the U.S. in the 2020–21 season and the low
454 severity but atypically late influenza season observed in the 2021–22 season. Population-level
455 behavior is difficult to predict, especially in the context of changing public health
456 recommendations and emerging SARS-CoV-2 variants, which complicated the process of
457 forecasting. Despite these challenges, FluSight forecasting teams provided forecasts of
458 confirmed influenza hospitalizations throughout each season, which helped public health
459 officials anticipate trends during the unusually prolonged influenza season in 2021-22, with
460 forecasting efforts extending into June, and then again for the atypically early 2022-23 season.

461
462 While the shift to forecasting for a new target presented a modeling challenge, the utility of the
463 corresponding new data source should be recognized. The HHS-Protect dataset [7] provided, in
464 addition to the state-level timeseries, facility-level data, which is at a higher spatial resolution
465 than other indicators of influenza activity. During the forecasting time frame analyzed here, the
466 data were also reported daily with previous day admission data published as soon as the day
467 after their occurrence, providing a timely source of information. As our data update analysis
468 (Figures S2 – S4) shows, these data demonstrated remarkably stable reporting behavior,
469 particularly during the 2021-22 season, with 94% of updates resulting in changes of under 10
470 hospitalizations for subnational jurisdictions. Stability of reporting decreased slightly during the
471 2022-23 season, with 83% of updates resulting in changes of under 10 hospitalizations for
472 subnational jurisdictions. Degraded forecast performance has been associated with large
473 revisions to initially observed values [6], and consistency in reporting is an important component
474 of a reliable forecasting target. Additionally, this dataset provided national and jurisdictional-level
475 data for confirmed influenza hospital admissions. In contrast with ILI, this indicator eliminated
476 the need to model outpatient visits associated with co-circulating non-influenza pathogens that
477 can cause ILI. Continued availability of rapid, disease-specific indicators of hospitalization, such
478 as those provided by these data, will facilitate improved forecasting utility and possibly
479 improvements in accuracy [27], particularly when forecasts are informed by mechanistic
480 transmission models.

481
482
483 The FluSight forecasting collaboration adapted quickly in 2021 to utilize a novel laboratory
484 confirmed influenza hospital admission dataset. Even with limited calibration data and atypical
485 influenza seasonality in the 2021-22 and 2022-23 seasons, the FluSight ensemble forecast
486 provided more robust forecasts than individual component models across spatial jurisdictions
487 and time horizons. This result mirrors those of other forecasting hubs. Collaborative hubs also
488 offer the ability for frequent feedback and interaction between modeling teams, providing
489 opportunities for rapidly sharing observations about underlying data and insights for forecast
490 development [28]. We observed poor coverage and general performance especially at the
491 beginning of the 2022-23 season and during other periods of rapid change. Collective insights
492 from these challenges can also inform when forecasts should be interpreted with extra caution.
493 Ongoing availability of the confirmed influenza hospitalization dataset, which covers all states,
494 could improve model calibration and ultimately contribute to the improvement of influenza
495 forecast performance and utility, as well as continued exploration and improvement of
496 forecasting and ensembling methodologies. These improvements are needed, particularly to
497 more accurately capture trends and appropriate levels of uncertainty during times of rapid
498 change.

499
500
501

502

504 **Tables and Figures**

505 Table 1: Performance metrics for teams submitting at least 75% of weekly FluSight targets.

Model	Absolute WIS	Relative WIS	MAE	50% Coverage (%)	95% Coverage (%)	% of Forecasts Submitted	Log Absolute WIS	Log Relative WIS
2021-22								
CMU-TimeSeries	12.54	0.74	18.92	47	90	100	0.31	0.78
Flusight-ensemble	13.86	0.82	20.79	48	86	100	0.33	0.83
PSI-DICE	14.03	0.83	20.17	43	82	100	0.33	0.84
UMass-trends_ensemble	14.35	0.84	22.24	71	97	100	0.36	0.91
Sgroup-RandomForest	15.45	0.91	23.87	47	95	100	0.38	0.97
CEID-Walk	15.63	0.94	22.19	52	82	89	0.39	0.99
Flusight-baseline	16.99	1.00	24.10	49	83	100	0.40	1.00
GT-FluFNP	17.57	1.02	23.40	39	69	96	0.38	0.98
MOBS-GLEAM_FLUH	17.17	1.03	22.25	32	63	91	0.42	1.08
SigSci-TSENS	17.79	1.03	24.86	38	72	96	0.40	1.01
IEM_Health-FluProject	17.69	1.04	23.98	50	85	100	0.40	1.02
CU-ensemble	18.32	1.08	25.41	44	77	100	0.39	0.98
LucompUncertLab-TEVA	21.02	1.22	29.99	54	86	89	0.41	1.05
UVAFluX-Ensemble	21.65	1.28	25.76	38	64	99	0.45	1.14
LucompUncertLab-VAR2_plusCOVID	22.03	1.30	28.99	42	74	94	0.42	1.08
UT_FluCast-Voltaire	23.64	1.39	35.19	50	91	99	0.45	1.14
LucompUncertLab-VAR2K_plusCOVID	24.44	1.42	32.43	42	74	89	0.47	1.20
LucompUncertLab-VAR2	25.93	1.53	35.05	39	72	94	0.53	1.35
LucompUncertLab-VAR2K	26.81	1.55	39.35	42	83	89	0.61	1.56
LosAlamos_NAU-Cmodel_Flu	28.69	1.69	36.14	26	59	100	0.63	1.60
Sgroup-SikJalpha	28.94	1.70	38.59	18	46	100	0.49	1.24
GH-Flusight	30.93	1.82	31.89	6	13	94	0.74	1.88
SigSci-CREG	27.36	1.93	31.00	19	44	89	0.80	2.03
2022-23								
MOBS-GLEAM_FLUH	42.20	0.61	57.97	42	78	94	0.37	0.65
CMU-TimeSeries	44.48	0.67	65.94	49	87	94	0.41	0.70
PSI-DICE	47.45	0.70	63.17	48	80	100	0.42	0.71
MIGHTE-Nsemble	48.99	0.72	67.50	53	82	96	0.41	0.70
Flusight-ensemble	51.72	0.76	71.04	56	81	100	0.44	0.74
Umass-trends_ensemble	53.86	0.80	79.40	63	89	100	0.49	0.83

Model	Absolute WIS	Relative WIS	MAE	50% Coverage (%)	95% Coverage (%)	% of Forecasts Submitted	Log Absolute WIS	Log Relative WIS
GT-FluFNP	59.75	0.81	72.88	56	75	89	0.53	0.89
CEPH-Rtrend_fluH	54.20	0.83	70.47	44	78	90	0.58	1.05
Sgroup-RandomForest	54.29	0.83	75.98	53	84	97	0.52	0.88
CU-ensemble	62.23	0.85	75.57	51	70	84	0.51	0.85
UGA_flucast-Okeeffe	62.13	0.94	77.33	50	72	95	0.61	1.02
SigSci-TSENS	64.27	0.96	80.02	58	74	93	0.66	1.09
Flusight-baseline	67.69	1.00	80.05	49	74	100	0.59	1.00
VTSanghani-ExogModel	72.30	1.00	92.56	30	61	81	0.63	1.05
UNC_IDD-InfluPaint	61.14	1.01	77.90	40	75	79	0.52	0.94
UVAFluX-Ensemble	78.71	1.11	94.45	22	41	95	0.61	1.03
SigSci-CREG	79.68	1.36	89.29	38	62	91	0.68	1.16
JHU_IDD-CovidSP	129.16	1.91	174.98	48	80	81	0.49	0.82

506

507 The Absolute WIS column refers to the Weighted Interval Score for each model across all fifty states, D.C., and
508 Puerto Rico forecast targets. The Relative WIS compares the WIS value of each model to the Flusight-baseline
509 model. All models with a relative WIS score less than one outperformed the baseline model when evaluated solely
510 upon WIS. 95% and 50% coverage values are provided for the percent of times that reported weekly incidence
511 values were within the 95% or 50% prediction intervals respectively, across all the forecast targets submitted by each
512 team. The percent of forecasts submitted is determined by the number of forecast targets submitted by each team out
513 of all possible forecast targets occurring within the duration of the evaluation period.

514

515

516

517 Table 2: One-to-four-week coverage and one-to-four-week percent of coverage above 90% for
 518 teams meeting inclusion criteria.

Model	Relative WIS	% WIS Below Baseline	Coverage				% Coverage above 90			
			1 Wk	2 Wk	3 Wk	4 Wk	1 Wk	2 Wk	3 Wk	4 Wk
2021-22										
CMU-TimeSeries	0.74	75.00	90.17	91.45	90.60	86.54	50.00	72.22	61.11	27.78
Flusight-ensemble	0.82	92.31	89.32	86.11	85.15	83.33	55.56	33.33	27.78	38.89
PSI-DICE	0.83	76.92	88.89	83.87	78.31	76.50	38.89	27.78	5.56	0.00
Umass-trends_ensemble	0.84	48.08	96.15	97.65	96.90	96.15	100.00	100.00	100.00	100.00
Sgroup-RandomForest	0.91	44.23	95.41	94.87	94.66	94.12	88.89	88.89	83.33	88.89
CEID-Walk	0.94	80.77	82.09	83.77	81.01	81.85	37.50	37.50	31.25	37.50
Flusight-baseline	1.00	0.00	82.26	84.19	82.48	81.62	27.78	22.22	22.22	22.22
GT-FluFNP	1.02	54.00	70.11	68.67	68.22	70.11	5.56	16.67	16.67	22.22
MOBS-GLEAM_FLUH	1.03	60.00	71.11	65.80	59.79	56.49	0.00	0.00	0.00	0.00
SigSci-TSENS	1.03	46.00	74.11	73.44	70.54	69.20	11.11	5.56	5.56	5.56
IEM_Health-FluProject	1.04	48.08	91.45	86.54	82.59	78.21	72.22	38.89	22.22	22.22
CU-ensemble	1.08	32.69	79.59	80.66	76.50	71.90	16.67	11.11	0.00	0.00
LucompUncertLab-TEVA	1.22	32.69	84.86	85.58	86.06	86.18	25.00	18.75	25.00	31.25
UVAFluX-Ensemble	1.28	25.00	66.05	65.51	62.58	60.95	11.11	0.00	0.00	0.00
LucompUncertLab-VAR2_plusCOVID	1.30	36.54	76.70	74.77	73.30	70.14	17.65	5.88	5.88	5.88
UT_FluCast-Voltaire	1.39	5.77	94.73	90.96	89.13	90.42	83.33	72.22	55.56	61.11
LucompUncertLab-VAR2K_plusCOVID	1.42	25.00	75.72	75.24	74.04	72.72	6.25	0.00	0.00	0.00
LucompUncertLab-VAR2	1.53	9.62	73.87	72.29	72.17	70.81	11.76	5.88	11.76	11.76
LucompUncertLab-VAR2K	1.55	9.62	81.97	81.49	83.05	85.46	6.25	18.75	25.00	37.50
LosAlamos_NAU-Cmodel_Flu	1.69	13.46	65.28	59.29	56.52	54.06	5.56	0.00	0.00	0.00
Sgroup-SikJalpha	1.70	1.92	40.28	45.73	48.08	48.29	0.00	0.00	0.00	0.00
GH-Flusight	1.82	5.77	18.33	12.90	11.99	10.63	0.00	0.00	0.00	0.00
SigSci-CREG	1.93	12.00	46.87	43.98	43.86	43.13	0.00	0.00	0.00	0.00
2022-23										
MOBS-GLEAM_FLUH	0.61	94.12	81.34	77.50	76.84	77.67	41.94	29.03	29.03	23.33
CMU-TimeSeries	0.67	86.54	86.27	87.12	87.25	86.31	58.06	64.52	70.97	70.00
PSI-DICE	0.70	92.31	88.03	81.27	77.17	74.87	64.52	67.74	64.52	60.00
MIGHTE-Nsemble	0.72	90.38	86.16	84.22	81.71	76.00	63.33	60.00	66.67	58.62
Flusight-ensemble	0.76	100.00	85.79	81.64	78.78	77.12	64.52	67.74	64.52	60.00

Model	Relative WIS	% WIS Below Baseline	Coverage				% Coverage above 90			
			1 Wk	2 Wk	3 Wk	4 Wk	1 Wk	2 Wk	3 Wk	4 Wk
Umass-trends_ensemble	0.80	92.31	90.88	89.89	87.41	85.83	77.42	74.19	70.97	70.00
GT-FluFNP	0.81	96.00	75.98	72.70	75.00	77.30	55.17	55.17	55.17	65.52
CEPH-Rtrend_fluH	0.83	67.31	75.82	80.22	79.33	77.21	46.43	50.00	57.14	44.44
Sgroup-RandomForest	0.83	92.31	90.06	84.49	81.86	79.71	73.33	70.00	70.00	65.52
CU-ensemble	0.85	75.00	71.60	71.38	69.90	66.85	46.15	53.85	53.85	52.00
UGA_flucast-Okeeffe	0.94	58.82	80.20	73.07	69.02	65.86	50.00	46.67	40.00	37.93
SigSci-TSENS	0.96	42.00	76.31	74.12	72.93	70.35	54.84	54.84	54.84	56.67
Flusight-baseline	1.00	0.00	78.72	74.26	71.34	68.85	58.06	58.06	58.06	56.67
VTSanghani-ExogModel	1.00	51.92	65.62	61.54	58.00	57.29	0.00	0.00	0.00	4.17
UNC_IDD-InfluPaint	1.01	64.71	75.20	74.25	75.12	75.14	52.00	44.00	64.00	54.17
UVAFluX-Ensemble	1.11	11.76	42.81	43.53	39.35	39.42	0.00	0.00	0.00	0.00
SigSci-CREG	1.36	6.00	68.28	62.27	58.85	55.34	48.39	48.39	45.16	43.33
JHU_IDD-CovidSP	1.91	33.33	86.74	81.67	78.18	73.60	65.38	61.54	53.85	48.00

Table 2: % WIS Below Baseline shows the percent of WIS values for each model below the corresponding Flusight-baseline WIS. The ‘% Coverage above 90’ columns show the percent of weekly 95% coverage values that are greater than or equal to 90% for each model by horizon. Modeling teams are ordered within each season by their relative WIS performance.

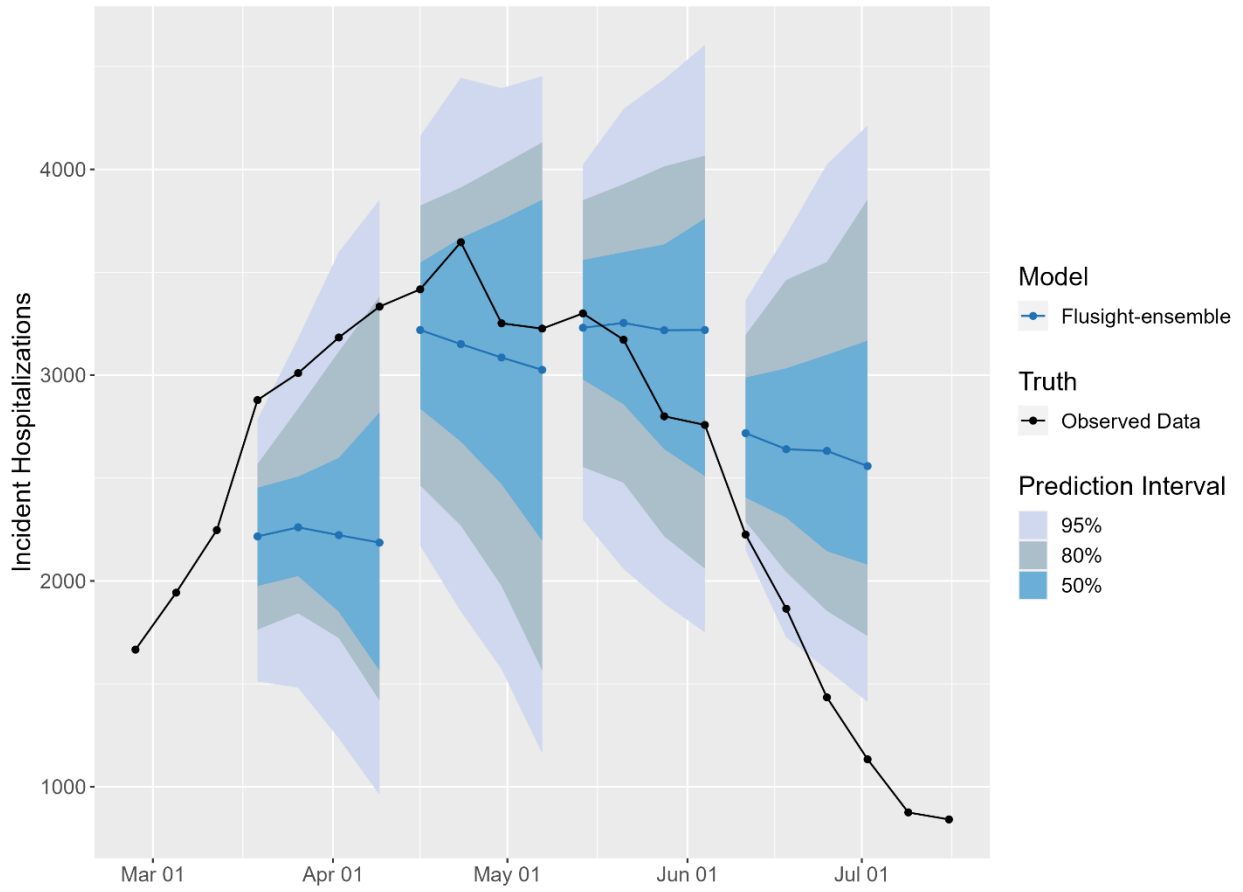
519

520

521

522 Figure 1: National weekly observed hospitalizations (black points) along with FluSight ensemble
 523 forecasts for four weeks of submissions in the 2021-22 season (panel a) and seven weeks of
 524 submissions in the 2022-23 season (panel b). The median FluSight ensemble forecast values
 525 (blue points) are shown with the corresponding 50%, 80%, and 95% prediction intervals (blue
 526 shaded regions).

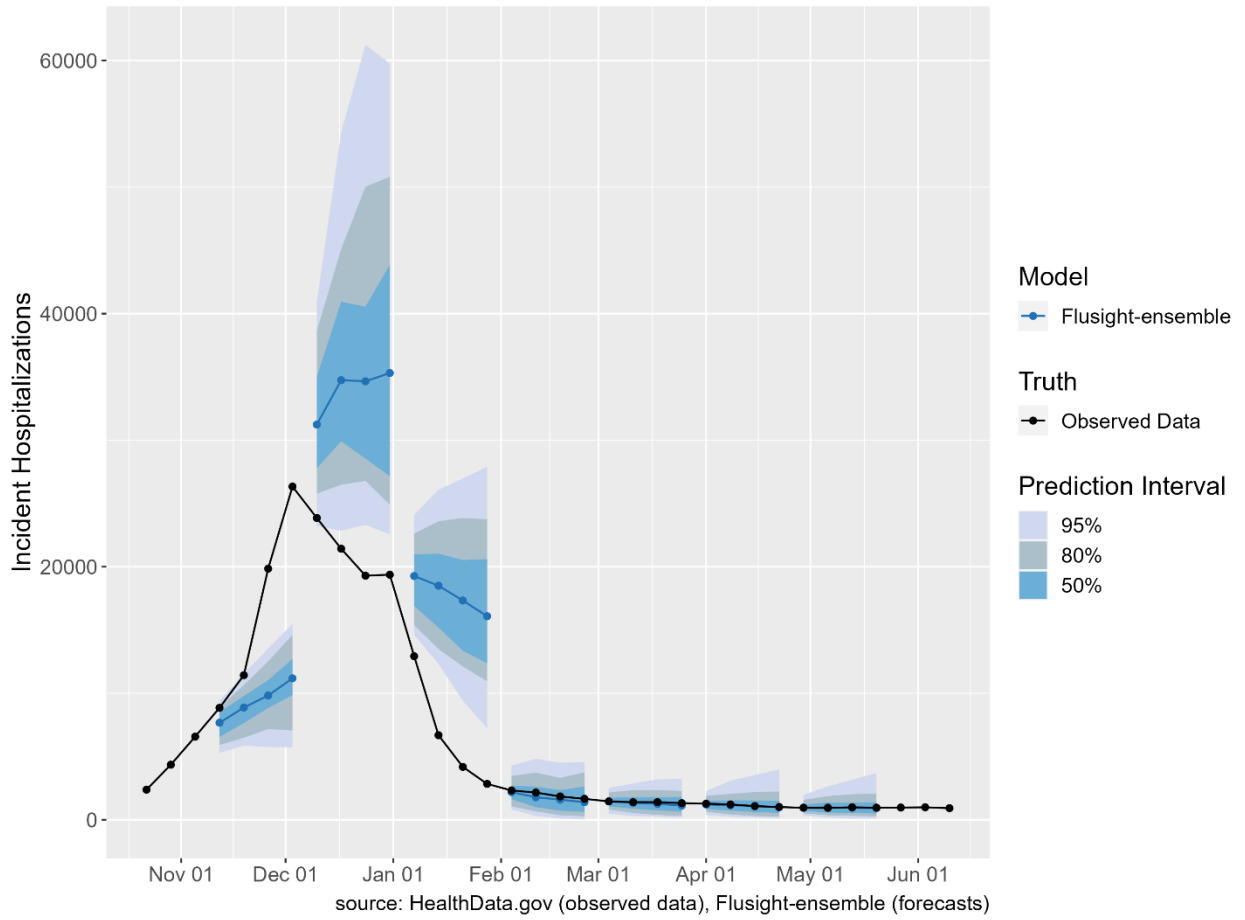
a 2021-22



source: HealthData.gov (observed data), Flusight-ensemble (forecasts)

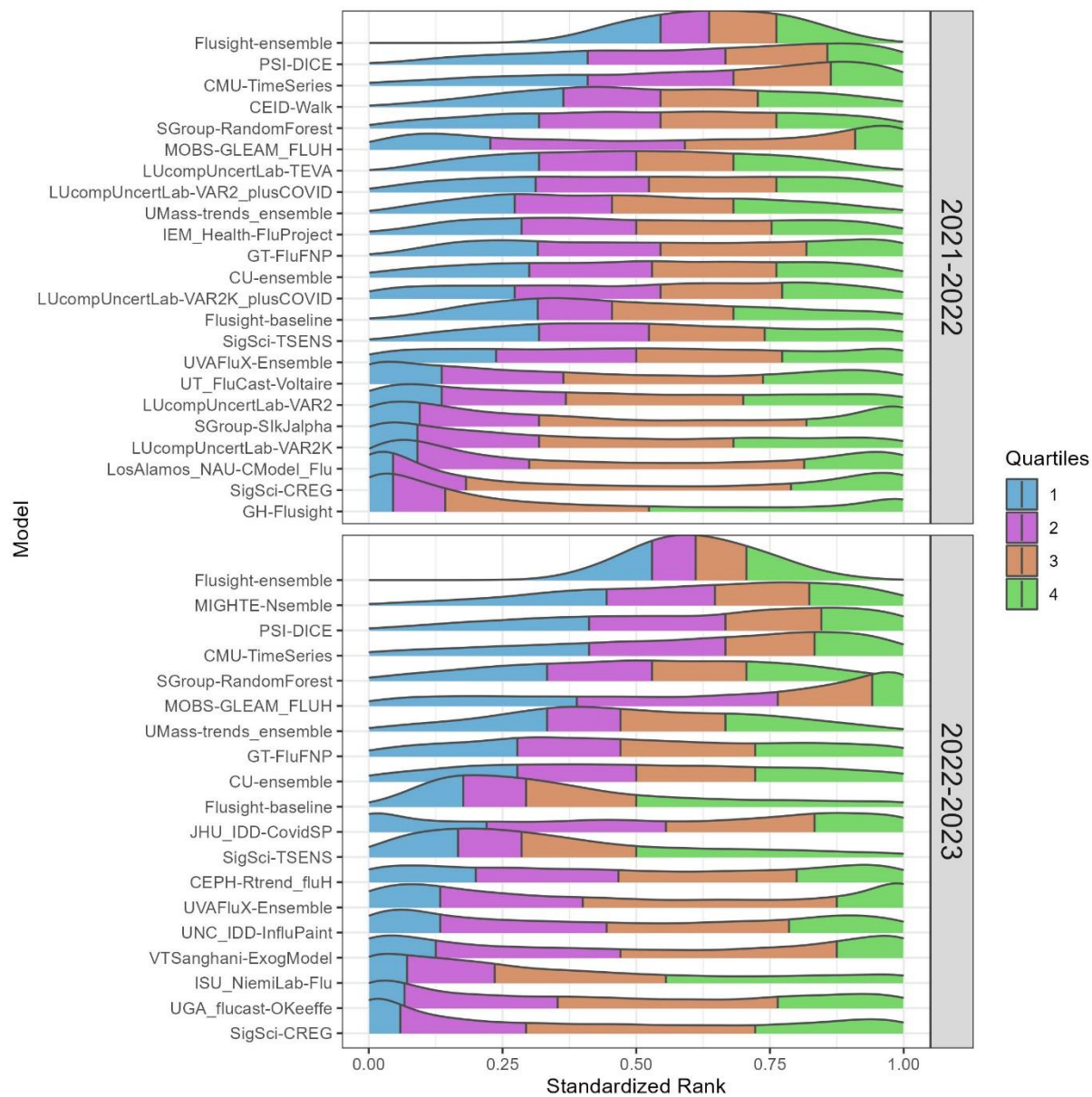
527

b 2022-23



528
529
530
531
532
533
534
535
536
537

538 Figure 2: Standardized rank of weighted interval score (WIS) over all forecast jurisdictions and
 539 horizons (1- to 4-week ahead), for the FluSight ensemble and each team submitting at least
 540 75% of the forecast targets (see Table 1 for qualifying teams and season metrics).



541

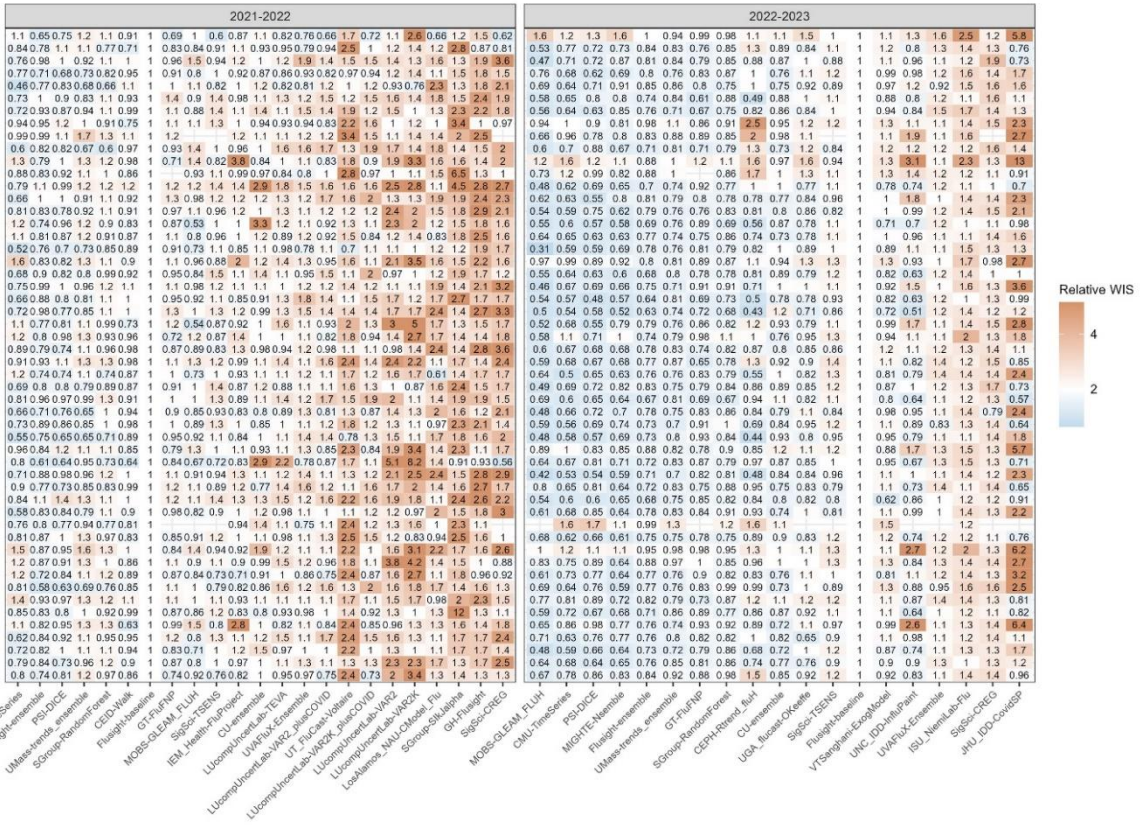
542

543

544

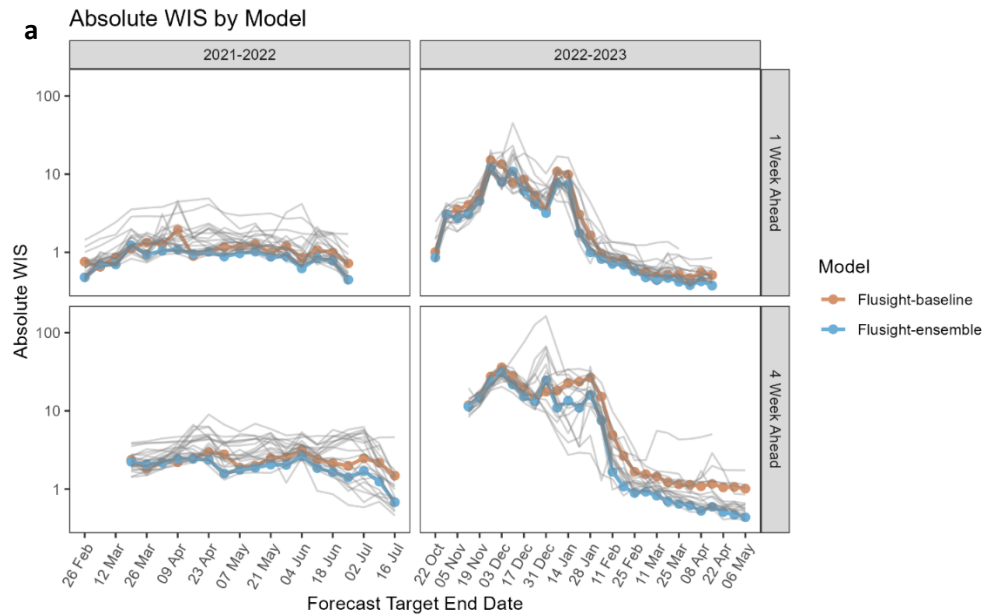
545 Figure 3: State-level WIS values for each team relative to the FluSight baseline model. The
 546 range of Relative WIS values below 1, in blue, indicate better performance than the FluSight
 547 baseline (white). Relative WIS values above 1, in red, indicate poor performance relative to the
 548 FluSight baseline. Teams are ordered on horizontal axis from lowest to highest Relative WIS
 549 values for each season. Analogous jurisdiction-specific relative WIS scores on log transformed
 550 counts are displayed in Figure S7.

551

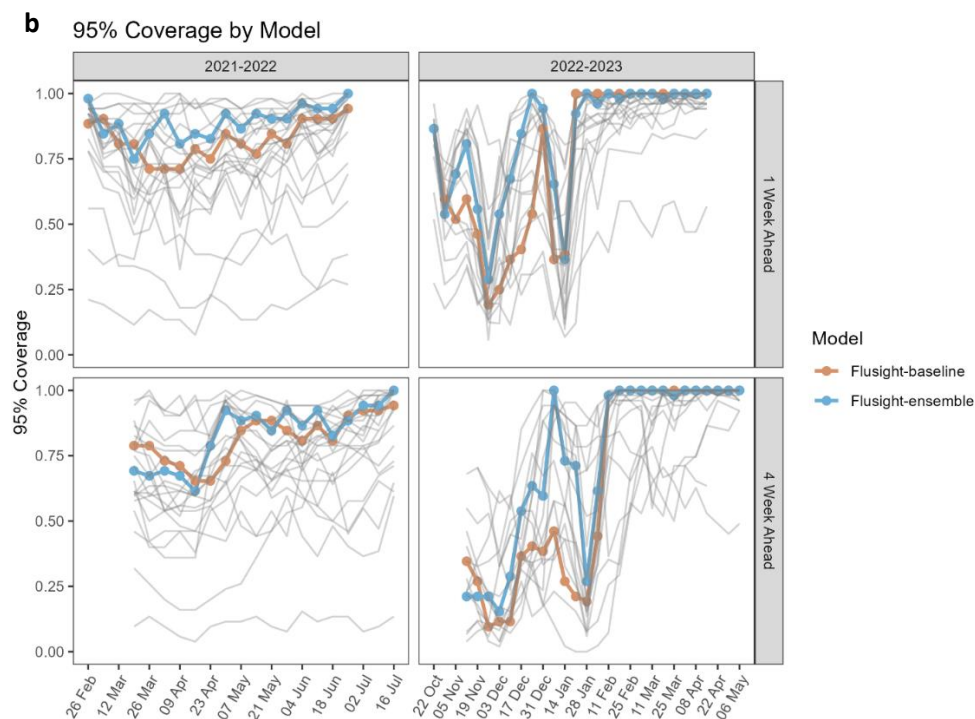


552
 553
 554
 555

556 Figure 4: Time series of log transformed absolute WIS (panel a) and 1- and 4-week ahead 95%
 557 coverage (panel b) for state and territory targets. Note that the forecast evaluation period
 558 translates to 1-week ahead forecast target end dates from February 26 to June 25, 2022,
 559 October 22, 2022, to May 20, 2023, and 4-week ahead forecast target end dates from March 19
 560 to July 16, 2022, and November 5, 2022, to June 10, 2023. Weekly results for the FluSight
 561 baseline and ensemble models are shown in red and blue respectively. Results for individual
 562 contributing models are shown in light gray.



563



564

565

566 **References**

- 567 1. Weekly U.S. Influenza Surveillance Report. <https://www.cdc.gov/flu/weekly/index.htm>. (2023)
- 568 2. U.S. Outpatient Influenza-like Illness Surveillance Network (ILINet).
569 <https://wwwn.cdc.gov/ILINet/>. (2023)
- 570 3. Influenza Hospitalization Surveillance Network (FluSurv-NET).
571 <https://www.cdc.gov/flu/weekly/influenza-hospitalization-surveillance.htm>. (2023)
- 572 4. Lutz C.S., et al. Applying infectious disease forecasting to public health: a path forward using
573 influenza forecasting examples. *BMC Public Health*, **19-1659** (2019).
- 574 5. McGowan C. J., et al.; Influenza forecasting working group, Collaborative efforts to forecast
575 seasonal influenza in the United States, 2015-2016. *Sci. Rep.* **9, 683** (2019).
- 576 6. Reich N. G., et al., A collaborative multiyear, multimodel assessment of seasonal influenza
577 forecasting in the United States. *Proc. Natl. Acad. Sci. U.S.A.* **116, 3146–3154** (2019).
- 578 7. COVID-19 Reported Patient Impact and Hospital Capacity by State Timeseries (RAW).
579 [https://healthdata.gov/Hospital/COVID-19-Reported-Patient-Impact-and-Hospital-Capa/g62h-](https://healthdata.gov/Hospital/COVID-19-Reported-Patient-Impact-and-Hospital-Capa/g62h-syeh)
580 [syeh](https://healthdata.gov/Hospital/COVID-19-Reported-Patient-Impact-and-Hospital-Capa/g62h-syeh). (2023).
- 581 8. COVID-19 Guidance for Hospital Reporting and FAQs For Hospitals, Hospital Laboratory,
582 and Acute Care Facility Data Reporting, 11 June 2023, [www.hhs.gov/sites/default/files/covid-19-](http://www.hhs.gov/sites/default/files/covid-19-faqs-hospitals-hospital-laboratory-acute-care-facility-data-reporting.pdf)
583 [faqs-hospitals-hospital-laboratory-acute-care-facility-data-reporting.pdf](http://www.hhs.gov/sites/default/files/covid-19-faqs-hospitals-hospital-laboratory-acute-care-facility-data-reporting.pdf).
- 584 9. Olsen S.J., et al., Changes in influenza and other respiratory virus activity during the COVID-
585 19 pandemic—United States, 2020–2021. *Morbidity and Mortality Weekly Report*. **Jul 7, 70(29)-**
586 **1013** (2021).
- 587 10. Merced-Morales A., et al. Influenza activity and composition of the 2022–23 influenza
588 vaccine—United States, 2021–22 season. *Morbidity and Mortality Weekly Report*. **Jul 7,**
589 **71(29)-913** (2022).
- 590 11. Centers for Disease Control and Prevention. [https://www.cdc.gov/flu/spotlights/2023-](https://www.cdc.gov/flu/spotlights/2023-2024/22-23-summary-technical-report.htm)
591 [2024/22-23-summary-technical-report.htm](https://www.cdc.gov/flu/spotlights/2023-2024/22-23-summary-technical-report.htm). (2023).
- 592 12. Centers for Disease Control and Prevention - MMWR Weeks definition.
593 https://ndc.services.cdc.gov/wp-content/uploads/MMWR_Week_overview.pdf. (2023).
- 594 13. FluSight Forecast Hub. <https://github.com/cdcepi/Flusight-forecast-data>. (2023).
- 595 14. Simple models for time series forecasting version 0.0.0.1000.
596 <https://rdr.io/github/reichlab/simplets/>. (2023).
- 597 15. Cramer E.Y., et al., Evaluation of individual and ensemble probabilistic forecasts of COVID-
598 19 mortality in the US. *Proc. Natl. Acad. Sci. U.S.A.*, <https://doi.org/10.1073/pnas.2113561119>
599 (2022).
- 600 16. Cramer E.Y., et al. The United States COVID-19 forecast hub dataset. *Scientific Data*. **Aug**
601 **1, 9(1), 1-5** (2022).
- 602 17. HHS Protect Public Data Hub. <https://public-data-hub-dhhs.hub.arcgis.com/>. (2023)

- 603 18. Bracher J., Ray E. L., Gneiting T., Reich N. G.,. Evaluating epidemic forecasts in an interval
604 format. *arXiv preprint arXiv:2005.12881* (2020).
605
- 606 19. Bosse N. I., et al.. Transformation of forecasts for evaluating predictive performance in an
607 epidemiological context. *PLOS Comput Biol*, **19(8):e1011393** (2023).
- 608 20. Reich N.G., et al. Accuracy of real-time multi-model ensemble forecasts for seasonal
609 influenza in the U.S. *PloS Comput Biol*, **15:e1007486** (2019).
- 610 21. Biggerstaff M., Slayton R. B., Johansson M. A., Butler J. C., Improving pandemic response:
611 employing mathematical modeling to confront COVID-19. *Clin. Infect. Dis.*, **10.1093/cid/ciab673**
612 (2021).
- 613 22. Reich N.G., et al. Collaborative hubs: making the most of predictive epidemic modeling.
614 *American Journal of Public Health*. **Apr 13(0):e1-4** (2022).
- 615 23. Howerton E., et al. Informing pandemic response in the face of uncertainty. An evaluation of
616 the U.S. COVID-19 scenario modeling hub. *medRxiv*. (2023).
- 617
- 618 24. Lopez V. K., et al. Challenges of COVID-19 case forecasting in the US, 2020-2021.
619 *medRxiv, Cold Spring Harbor Laboratory Press*, **1 Jan.** (2023).
620
- 621 25. Srivastava A., Singh S., Lee F. Shape-based evaluation of epidemic forecasts. *arXiv*
622 *preprint arXiv:2209.04035*. (2022)
- 623 26. Adiga A., et al. Phase-Informed Bayesian ensemble models improve performance of
624 COVID-19 forecasts. *Proceedings of the AAAI Conference on Artificial Intelligence*, **37(13)**,
625 **15647-15653**. <https://doi.org/10.1609/aaai.v37i13.26855>. (2023).
- 626
- 627 27. Fox J., et al., Real-time pandemic surveillance using hospital admissions and mobility data.
628 *Proc. Natl. Acad. Sci. U.S.A.* **119**, **10.1073/pnas.2111870119** (2022).
- 629 28. Johansson M.A., et al. An open challenge to advance probabilistic forecasting for dengue
630 epidemics. *Proc. Natl. Acad. Sci. U S A*. **116-24**, **268–74** (2019).
- 631 29. Zoltar Forecast Archive. <https://zoltardata.com/project/299>. (2023).
- 632 30. R Core Team. R: A language and environment for statistical computing. R foundation for
633 statistical computing, Vienna, Austria. <<https://www.R-project.org/>>. (2023).

634
635
636

637 **Disclaimers**

638 Any use of trade, firm, or product names is for descriptive purposes only and does not imply
639 endorsement by the U.S. Government. The findings and conclusions in this report are those of
640 the authors and do not necessarily represent the views of the Centers for Disease Control and
641 Prevention or the National Institutes of Health.

642 **Data Availability**

643 The forecasts from models used in this paper are available from the FluSight Forecast Hub
644 GitHub repository (<https://github.com/cdcepi/Flusight-forecast-data>) [13] and the Zoltar forecast
645 archive (<https://zoltardata.com/project/299>) [29]. These are both publicly accessible. The code
646 used to generate all figures and tables in the manuscript will be available in a public
647 repository (<https://github.com/cdcepi/FluSight-manuscripts>) at the time of publication. All
648 analyses were conducted using the R language for statistical computing (version 4.0.3) [30].
649

650 **Acknowledgments**

651 The authors would like to acknowledge Michael A. Johansson and Nicole Samay for their
652 contributions to this work.
653

654 M.B.N., P.R., J.T., S.V., A.A., G.K., B.H., B.L.L., M.V.M., M.A.A, A.S. disclose support for the
655 research of this work from the Centers for Disease Control and Prevention (CDC) and Council
656 of State and Territorial Epidemiologists (CSTE), [Cooperative Agreement number
657 NU380T000297]
658

659 M.C., J.D, K.M., X.X, APP, AV, P.C.V., A.G.K. M.L., M.A. disclose support for the research of
660 this work from the HHS/CDC 6U01IP001137 and HHS/CDC 5U01IP0001137.
661

662 B.A.P., A.R., H.P.K., Z.Z., G.G., P.A., S.S.B., R.V.V., disclose support for the research of this
663 work from NSF (Expeditions CCF-1918770, CAREER IIS-2028586, RAPID IIS-2027862,
664 Medium IIS-1955883, Medium IIS-2106961, PIPP CCF-2200269), CDC MInD program, faculty
665 gifts from Facebook/Meta, and funds/computing resources from Georgia Tech and GTRI.
666

667 M.S., L.C. F.L and A.G. M. disclose support for the research of this work from the Centers for
668 Disease Control and Prevention (CDC) and Council of State and Territorial Epidemiologists
669 (CSTE), [Cooperative Agreement number NU380T000297]
670

671 M.S. discloses support for the research of this work from the National Institutes of Health (grant
672 number R01GM130668) and (in part) by contract 200-2016-91779 with the Centers for Disease
673 Control and Prevention.
674

675 S.V., A.A., G.K., B.H., B.L.L., M.V.M., also disclose support for the research of this work from
676 NSF Expeditions CCF-1918656, VDH Grant VDH-21-501-0135, University of Virginia Strategic
677 Investment Fund Award SIF160.
678

679 L.C.B., A.G., A.J.H., D.J.M., R.R., D.S., R.J.T. disclose support for the research of this work
680 from the Centers for Disease Control and Prevention U011P001121 and Centers for Disease
681 Control and Prevention 75D30123C15907.
682

683 B.T.S., S.A.S., H.L.G., and P.B. disclose support for the research of this work from the Centers
684 for Disease Control and Prevention (CDC) and Council of State and Territorial Epidemiologists
685 (CSTE), [Cooperative Agreement number NU380T000297]
686

687 N.R. discloses support from National Science Foundation grants CCF-1918770, NRT DGE-
688 1545362, and OAC-1835660.
689

690 S.T., C.P.S., A.H. disclose support for the research of this work from the National Science
691 Foundation [2127976]. S.T., C.P.S., A.H., J.L., J.C.L., S.L.L., C.D.M., K.S., S-m.J. disclose
692 support from the Centers for Disease Control and Prevention [200-2016-91781].
693

694 J.L. and J.C.L. disclose support from the National Institutes of Health (NIH 5R01AI102939)
695

696 A.M., Y.T.L., and W.S.H. disclose support for the research of this work from Laboratory Directed
697 Research and Development Program at Los Alamos National Laboratory [20220268ER].

698 W.S.H., R. G. P., S. L., Y. C. disclose support for the research of this work from National
699 Institute of Health [R01GM111510].

700

701 **Author contributions:** S.M.M., A.E. W, M.B, and R.K.B. contributed to conceptualization.
702 S.M.M. and R.K.B. wrote the original draft of the manuscript. S.M.M. and A.E.W. performed the
703 formal analysis. M.B. and C.R. performed supervision and project administration. All authors
704 contributed modeling data and T.M.L., E.L.M., M.S., L.A.W., L.C.B., A.G., A.J.H., D.J.M., R.R.,
705 D.S., R.J.T., S.K., S.P., J.S., R.Y., T.K.Y., Y.A., S.B., G.G., H.K., B.A.P., R.R., A.R., Z.Z., A.M.,
706 S.O., P.B., H.L.G., S.A.S., B.T.S., M.A., A.G.K., M.L., P.C.V., S.W., J.N., E.C., A.L.H., S.J.,
707 J.C.L., J.L. S.L.L., C.D.M., K.S., C.S., S.T., T.M., W.Y., N.B., W.S.H., Y.T.L., A.M., Y.C., S.M.L.,
708 J.L., R.G.P., A.C.P., C.V., L.C., F.L., A.G.M., M.S., M.C., J.T.D., K.M., A.P.P., A.V., X.X.,
709 M.B.N., P.R., J.T., C.H.L., S.J., V.P.N., S.D.T., D.W., A.V., J.M.D., S.J.F., G.C.G., E.S., E.W.T.,
710 M.G.C., W.Y.C., A.G., A.S., E.L.R., N.G.R., L.S., N.W., Y.W., M.W.Z., M.A.A., A.S., L.A.M.,
711 A.A., B.H., G.K., B.L.L., M.M., S.V., P.B., A.F., N.M., AND N.R. submitted forecast data for the
712 analysis. All authors contributed to the review and editing of the manuscript.

713

714 **Competing interests:** E.W.T. is an employee of Sanofi, which manufactures influenza
715 vaccines. J.S. and Columbia University disclose partial ownership of SK Analytics. J.S.
716 discloses consulting for BNI.

717

718

719

720

721

722

