

# A Collaborative Ensemble Approach to Real-Time Influenza Forecasting in the U.S.: Results from the 2017/2018 Season

(working author list subject to re-ordering pending final contributions)

Nicholas G Reich<sup>1</sup>, Craig McGowan<sup>2</sup>, Logan Brooks<sup>3</sup>, Sasikiran Kandula<sup>4</sup>,  
Evan Moore<sup>1</sup>, Dave Osthus<sup>5</sup>, Evan Ray<sup>6</sup>, Abhinav Tushar<sup>1</sup>, Teresa Yamana<sup>4</sup>,  
Willow Crawford-Crudell<sup>7</sup>, Graham Casey Gibson<sup>1</sup>, Rebecca Silva<sup>8</sup>  
Matthew Biggerstaff<sup>2</sup>, Michael A Johansson<sup>9</sup>, Roni Rosenfeld<sup>3</sup>, Jeffrey Shaman<sup>4</sup>

<sup>1</sup>University of Massachusetts-Amherst, Amherst, USA

<sup>2</sup>Influenza Division, Centers for Disease Control and Prevention, Atlanta, USA

<sup>3</sup>Carnegie Mellon University, Pittsburgh, USA

<sup>4</sup>Columbia University, New York, USA

<sup>5</sup>Los Alamos National Laboratory, Los Alamos, USA

<sup>6</sup>Mount Holyoke College, South Hadley, USA

<sup>7</sup>Smith College, Northampton, USA

<sup>8</sup>Amherst College, Amherst, USA

<sup>9</sup>Division of Vector-Borne Diseases, Centers for Disease Control and Prevention, Atlanta, USA

## Abstract

TBD.

# 1 Introduction

Seasonal influenza results in a substantial, annual public health burden in the United States and worldwide. In the influenza season running from October 2017 through May 2018, one of the largest seasonal outbreaks on record, the United States Centers for Disease Control and Prevention estimates there were an estimated XX million cases and XXX,000 hospitalizations.[1] The CDC utilizes a variety of surveillance methods to assess the severity of an influenza season, including monitoring outpatient visits for influenza-like illness (ILI), influenza-related hospitalizations, and virologic characteristics.[2] However, like all surveillance systems, these are constrained to describing events that have already taken place, and the total burden, along with the timing of the epidemic, can vary substantially from season to season.[1] Forecasts of an influenza season offer the possibility of providing actionable information to improve public health responses, and recent years have seen a large amount of peer-reviewed research describing efforts to predict seasonal influenza.[3, 4, 5, 6, 7, 8, 9, 10]

Ensemble models, i.e. methods that bring together predictions from multiple different component models, have long been seen as a valuable method for improving predictions over any single model. This "wisdom of the crowd" approach has both theoretical and practical advantages. First, it allows for an ensemble forecast to incorporate signals from different data sources and models that may highlight different features of a system. Second, combining signals from models with different biases may allow those biases to offset and result in an ensemble that is more accurate than the individual ensemble components. Weather and climate models have utilized ensemble systems for these very purposes, and recent work has extended ensemble forecasting to infectious diseases, including influenza, dengue fever, lymphatic filariasis, and Ebola hemorrhagic fever.[11, 12, 13, 14]

Since the 2013/2014 influenza season, the CDC has run an annual prospective influenza forecasting competition, known as the FluSight challenge, in collaboration with outside researchers. Participating teams submit probabilistic forecasts for various influenza-related targets of interest to public health officials weekly from early November through mid May. Among other government-sponsored infectious disease forecasting competitions in recent years,[15, 16] this challenge been unique in its prospective orientation over multiple outbreak seasons Also, it has provided a venue for close interaction and collaboration between government public health officials and academic and private-sector researchers.

The FluSight challenge has been designed and retooled over the years with an eye towards maximizing the public health utility and integration of forecasts with real-time public health decision making. All forecast targets are derived from the trajectories of U.S. region-level weighted influenza-like illness (wILI), an estimate of the percentage of outpatient visits due to ILI weighted by state populations. ILI is perhaps the most frequently used measure of the burden of influenza-like respiratory illness in epidemiological surveillance. Weekly submissions to the FluSight challenge contain probabilistic and point forecasts for seven targets in each of 11 regions in the U.S. (national-level plus the 10 Health and Human Services (HHS) regions, Figure 1A). There are two classes of targets: "week-ahead" and "seasonal". "Week ahead" targets refer to the four weekly targets (incidence 1, 2, 3 and 4 weeks in the future) that are different for each week of the season. "Seasonal" targets refer to quantities (outbreak onset, outbreak peak week, and outbreak peak intensity) that do not change for a region within a season (see Figure 1B and Methods).

In March 2017, a group of influenza forecasters from different institutions who have worked with the CDC in the past established the FluSight Network. This research consortium worked collaboratively throughout 2017 and 2018 to build and implement in real-time a multi-institution ensemble with performance-based model weights. During the 2015/2016 and 2016/2017 FluSight challenges, analysts at the CDC built a simple ensemble model

for all targets by taking the arithmetic mean of all submitted models. This model was one of the top performing models each season (McGowan et al., under revision).

A central goal of the FluSight Network was to demonstrate the benefit of performance-based weights in a real-time, multi-team ensemble setting by outperforming the “simple average” ensemble that CDC uses to inform decision making and situational awareness during the annual influenza season. In this paper, we describe the development of this collaborative ensemble model and present results from both retrospective (2010 - 2017) and prospective (2017-2018) forecast evaluations. The FluSight Network assembled 21 ensemble components to build ensemble models for seasonal influenza outbreaks (Table 1). These components encompassed a variety of different modeling philosophies, including Bayesian hierarchical models, mechanistic models of infectious disease transmission, statistical learning methodologies, and classical statistical models for time-series data. We show that using ensemble models informed by past component performance consistently improved forecast accuracy. Given the fortuitous timing of this experiment, during the most severe seasonal influenza season on record, we provide the first evidence from a real-time forecasting study that performance-based weights can improve ensemble forecast performance high severity infectious disease outbreaks. This research is an important example of a collaboration between government and academic public health experts, setting an important precedent and prototype for real-time collaboration in more severe outbreaks, such as a global influenza pandemic.

## 2 Results

### 2.1 Summary of ensemble components

Individual component model forecast performance in the seven training seasons (2010/2011 - 2016/2017) varied widely across region, season, and target. A detailed comparative analysis of component forecast performance can be found elsewhere[22], however we summarize a few key insights from model performance here. A seasonal baseline model, whose forecasts for a particular target are based on data from previous seasons and do not update based on data from the current season, was used as a reference point for other models. Over 50% of the ensemble components out-performed the seasonal baseline model in forecasting 1-, 2-, and 3-week ahead incidence as well as season peak percentage and season peak week. However, season-to-season variability in forecast performance was large, as 10 models had, in at least one season, better overall average accuracy than the model with the best average performance across all seasons. To evaluate model accuracy, we followed CDC convention and used a metric that takes the geometric average of the probabilities assigned to the eventually observed value. This measure, which we refer to as “forecast score”, can be interpreted as the average probability a given forecast model assigned to the eventually observed value. As such, higher values, on a scale of 0 to 1, indicate more accurate models.

### 2.2 Choice of ensemble model based on cross-validation

The FSNetwork Target-Type Weights (FSNetwork-TTW) ensemble model outperformed all other ensemble and ensemble components in the training phase by a slim margin. This model was one of five pre-specified ensemble approaches defined prior to any systematic evaluation of ensemble component performance in previous seasons and prior to the 2017/2018 season. Using 42 estimated weights, one for each model and target-type (week-ahead

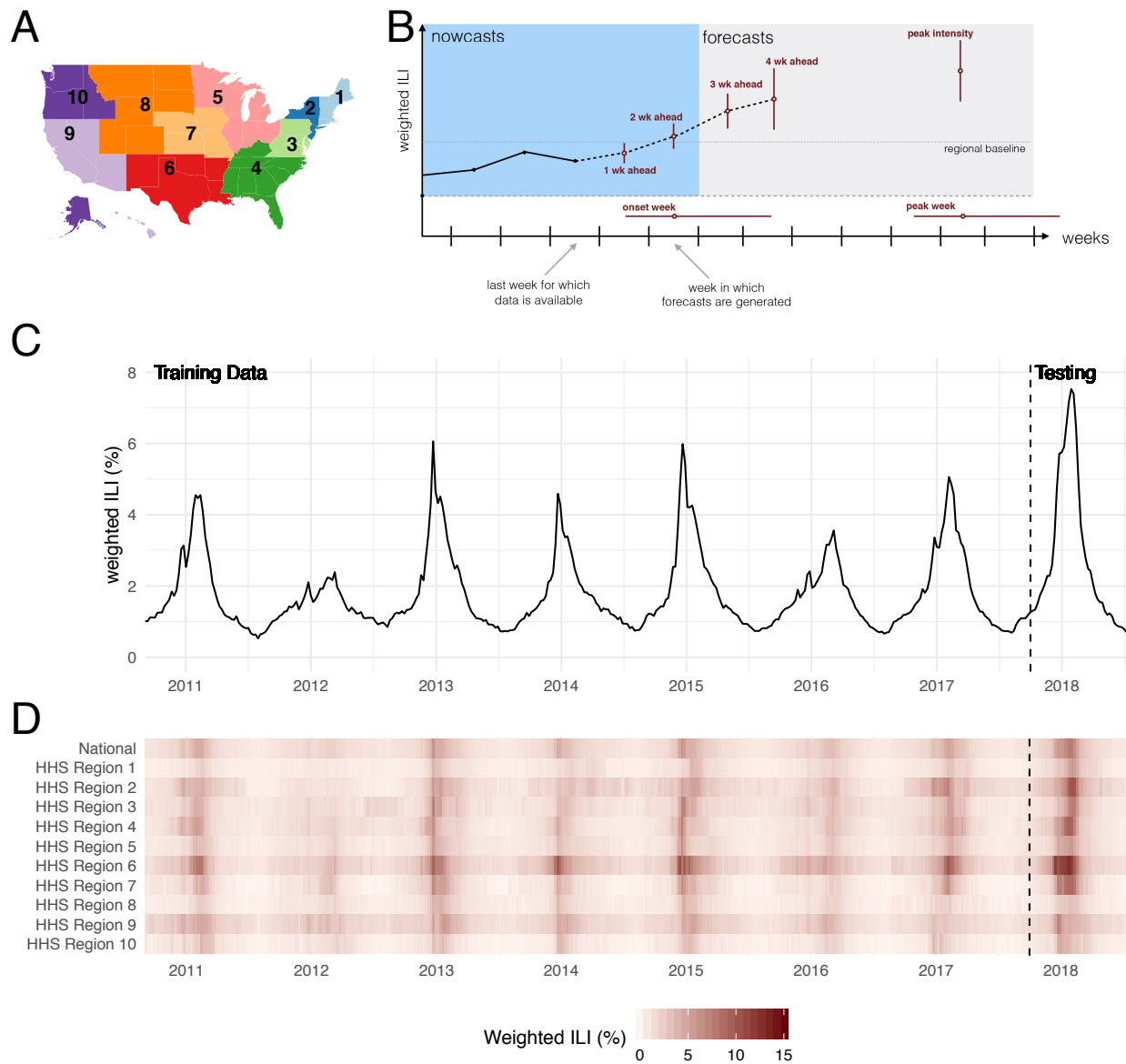


Figure 1: (A) Map of the 10 U.S. Health and Human Services regions. Influenza forecasts are made at this scale. (B) A figure showing the structure of a single forecast. Seven forecasting targets are illustrated with a point estimate (dot) and interval (uncertainty bars). The five targets on the wILI scale are shown with uncertainty bars spanning the vertical wILI axis, while the two targets for a time-of-year outcome are illustrated with horizontal uncertainty bars along the temporal axis. The onset is defined relative to a region- and season-specific baseline wILI percentage defined by the CDC.[17] Arrows illustrate the timeline for a typical forecast for the CDC FluSight challenge, assuming that forecasts are generated or submitted to the CDC using the most recent reported data. These data include the first reported observations of wILI% from two weeks prior. Therefore, 1 and 2 week-ahead forecasts are referred to as nowcasts, i.e., at or before the current time. Similarly, 3 and 4 week-ahead forecasts are forecasts, or estimates about events in the future. This figure has appeared previously in (Reich et al, under review). (C) Publicly available weighted influenza-like illness (wILI) data from the CDC website for the national level. The y-axis shows the weighted percentage of doctor's office visits in which a patient presents with influenza-like illness for each week from September 2010 through July 2018, which is the time period for which the models presented in this paper made seasonal forecasts. (D) Publicly available wILI data for each of the 10 HHS regions. Darker colors indicate higher wILI.

Team	Model Abbr	Model Description	Ens. Model
FSNetwork	EW	Equal Weights (number of unique weights = 1)	x
	CW	Constant Weights (21)	x
	TTW	Target-Type Weights (42)	x
	TW	Target Weights (147)	x
	TRW	Target-Region Weights (1,617)	x
CU	EAKFC_SEIRS	Ensemble Adjustment Kalman Filter SEIRS	[18]
	EAKFC_SIRS	Ensemble Adjustment Kalman Filter SIRS	[18]
	EKF_SEIRS	Ensemble Kalman Filter SEIRS	[4]
	EKF_SIRS	Ensemble Kalman Filter SIRS	[4]
	RHF_SEIRS	Rank Histogram Filter SEIRS	[4]
	RHF_SIRS	Rank Histogram Filter SIRS	[4]
	BMA	Bayesian Model Averaging	[11]
Delphi	BasisRegression	Basis Regression (epiforecast defaults)	[19]
	DeltaDensity1	Delta Density (epiforecast defaults)	[9]
	EmpiricalBayes1	Empirical Bayes (conditioning on past four weeks)	[20, 19]
	EmpiricalBayes2	Empirical Bayes (epiforecast defaults)	[20, 19]
	EmpiricalFuture	Empirical Futures (epiforecast defaults)	[19]
	EmpiricalTraj	Empirical Trajectories (epiforecast defaults)	[19]
	DeltaDensity2	Markovian Delta Density (epiforecast defaults)	[9]
	Uniform	Uniform Distribution	
	Stat	Ensemble (combination of 8 Delphi models)	x [9]
LANL	DBM	Dynamic Bayesian SIR Model with discrepancy	[8]
ReichLab	KCDE	Kernel Conditional Density Estimation	[21]
	KDE	Kernel Density Estimation and penalized splines	[12]
	SARIMA1	SARIMA model without seasonal differencing	[12]
	SARIMA2	SARIMA model with seasonal differencing	[12]

Table 1: List of models, with key characteristics. Team abbreviations are translated as: FSNetwork = FluSight Network, CU = Columbia University, Delphi = Carnegie Mellon, LANL = Los Alamos National Laboratories, ReichLab = University of Massachusetts Amherst. The ‘Ext data’ column notes models that use data external to the ILINet data from CDC. The ‘Mech. model’ column notes models that rely to some extent on an mechanistic or compartmental model formulation. The ‘Ens. model’ column notes models that are ensemble models.

and seasonal) combination, the FSNetwork-TTW model built a weighted model average using a predictive density stacking approach (see Figure 3 and Methods). In the training period consisting of the seven influenza seasons prior to 2017/2018, this model achieved a leave-one-season-out cross-validated average forecast score of 0.41, compared with the FSNetwork Target Weights (FSNetwork-TW) model with a score of 0.40, the FSNetwork Constant Weights (FSNetwork-CW) model with a score of 0.40, and the FSNetwork Target-Region Weights (FSNetwork-TRW) model with a score of 0.40 (Figure 4). We chose the target-type weights model as the model that would be submitted in real-time to the CDC during the 2017/2018 season, based on the pre-specified criteria of it having the highest score of any approach in the cross-validated training phase.

Using out-of-sample cross-validated performance of all ensemble components across the seven training seasons, we estimated weights for the chosen FSNetwork-TTW ensemble model that would be used for the 2017/2018 real-time forecasting. The FSNetwork-TTW model assigned non-negligible weight (greater than 0.001) to 8 models for week-ahead targets and 6 models for seasonal targets (Figure 3). For week-ahead targets, the highest non-zero weight (0.42) was given to the Delphi-DeltaDensity1 model. For seasonal targets, the highest weight (0.26) was given to the LANL-DBM model. In the weights for the seasonal targets, six models shared the weight, with none of the six having less than 0.11 weight. All four research teams had at least one model with non-negligible weight in the chosen model. It must be noted that ensemble weights themselves are not a measure of a component's standalone accuracy nor its contribution to the overall accuracy of the ensemble.

## 2.3 Summary of ensemble real-time performance in 2017/2018 season

The 2017/2018 influenza season in the U.S. exhibited features that were unlike that of any season in the past 15 years. As measured by wILI percent at the national level, the 2017/2018 season was on par with the other two highest peaks on record (since 1997): the 2003/2004 season and the 2009 H1N1 pandemic. In some regions, for example HHS Region 2 (New York and New Jersey) and HHS Region 4 (southeastern states), the highest reported wILI percent for the 2017/2018 season was more than 20% above previously observed peaks. Because all forecasting models rely, to some extent, on future trends mimicking observed patterns in the past, the anomalous dynamics in 2017/2018 posed a challenging "test season" for all models, including the new ensembles. Indeed, some of the models that saw the largest drop in performance (e.g., Delphi-DeltaDensity1, ReichLab-KCDE, and Delphi-DeltaDensity2) are ones that explicitly rely on using kernel conditional density estimation to find previously observed regions of the time-series that bear a similarity to current trends.

In spite of these unusual dynamics negatively impacting the forecast accuracy of the top-performing ensemble components, the FSNetwork-TTW ensemble model that varied weights by model and target-type showed the best performance among all selected models in the 2017/2018 season (Figure 4).

- compare all ensembles with single best model from each team in CV stage and with unweighted average from FSN from 2015/2016 - 2017/2018 -

Despite being optimized for high log-score values, the FSNetwork-TTW showed robust performance across a variety of different evaluation metrics in the 2017/2018 season. Its root mean squared error was XX XX [add rank raster?] (See Appendix). According to the probability integral transform metric[], the FSNetwork-TTW model was well-calibrated for all four week-ahead targets showing no significant deviations. ADD BIAS!! It was slightly less well-calibrated for peak performance, and showed indications of having too narrow predictive distributions over the 2017/2018 season.

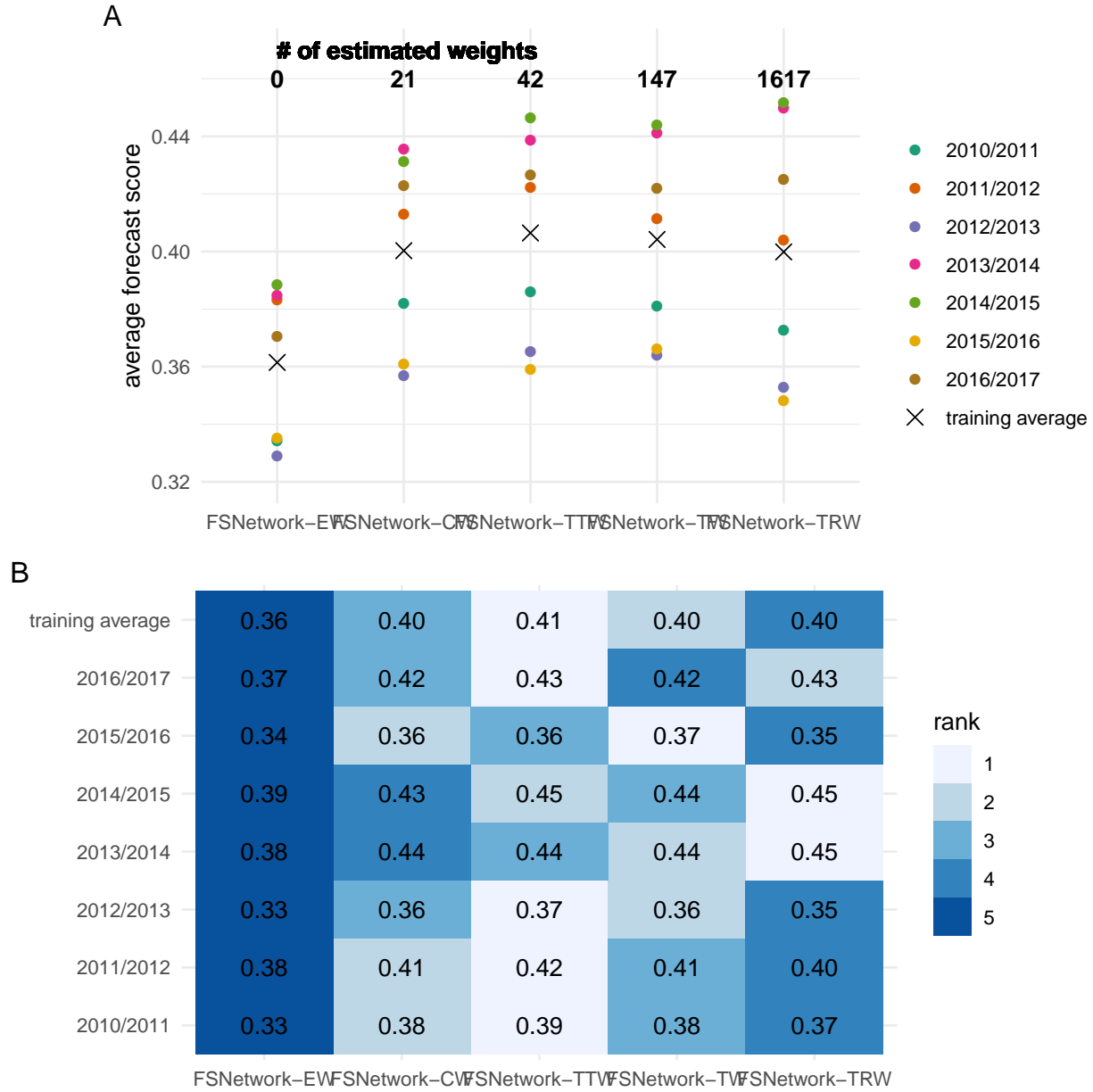


Figure 2: Comparison of five ensemble models during the training phase. In both panels, the models are sorted from simplest (left) to most complex (right), with the number of estimated weights for each model shown at the top of Panel A. (A) forecast scores by season, with overall average given by the X. (B) Model Rank within each season and average forecast score (across weeks, targets, and regions).

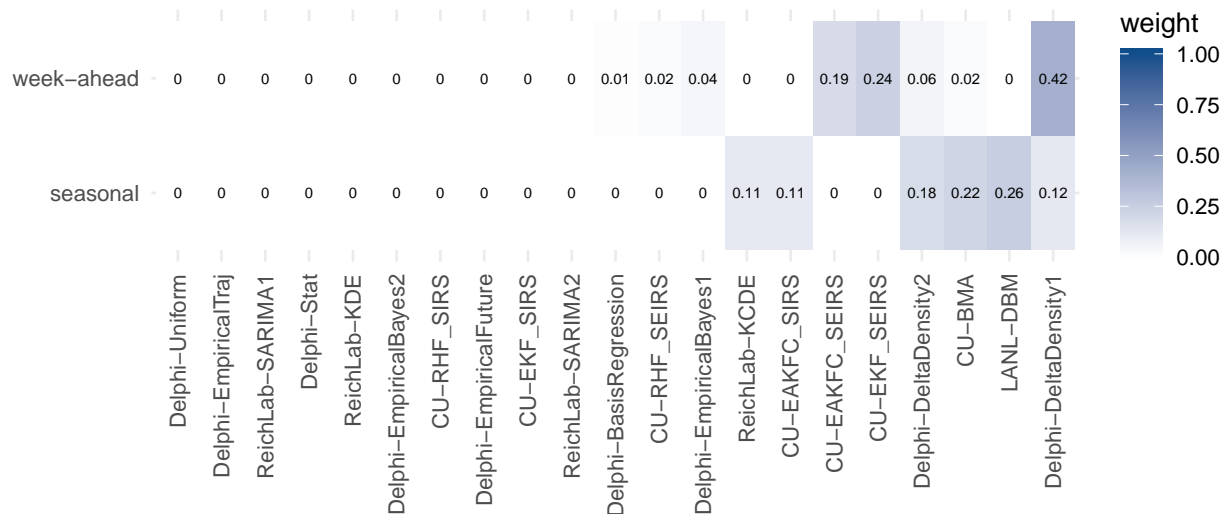


Figure 3: Model weights for the FluSight Network Target-Type Weights (FSNetwork-TTW) model. Weights were estimated using cross-validated forecast performance in the 2010/2011 through the 2016/2017 seasons.

The model accuracy patterns observed in the cross-validation phase was largely preserved in the real-time phase. In the cross-validation phase, ensemble models had slightly higher average performance than all ensemble components, and ranked as the top model in each season].

When taking into account the performance of ensemble components in the 2017/2018 season, the weights for a subsequent hypothetical ensemble using the same components would be different. Components that received lots of weight in the original ensemble but did particularly poorly in the 2017/2018 season saw the largest drop in weight. Overall, three components were added to the list of six existing components that received more than 0.001 weight for seasonal targets: CU-EAKFC\_SEIRS, CU-EKF\_SEIRS, and ReichLab-SARIMA2. One component (ReichLab-SARIMA2) was added to the list of eight existing components that received more than 0.001 weight for week-ahead targets.

### 3 Discussion

Ensembles hold promise for giving decision makers the ability to use “one answer” that combines the strengths of many different modeling approaches while mitigating their weaknesses. This work presents the first attempt to systematically combine infectious disease forecasts from multiple research groups in real-time using an approach that factors in past performance of each component method. Of the 29 models submitted to the CDC in 2017/2018 as part of their annual FluSight forecasting challenge, this ensemble was the second-highest scoring model overall. (The top scoring model was an ensemble of human judgement forecasts.[23]) In the 2018/2019 influenza season (forthcoming, at the time of writing), based on results from this study, the CDC used forecasts from the FluSight Network ensemble model in internal and external communication and planning reports.

Even in a very unusual influenza season, the ensemble approach was a steady contributor and did not see a large reduction in overall performance compared to performance during the training seasons. This bodes well for the



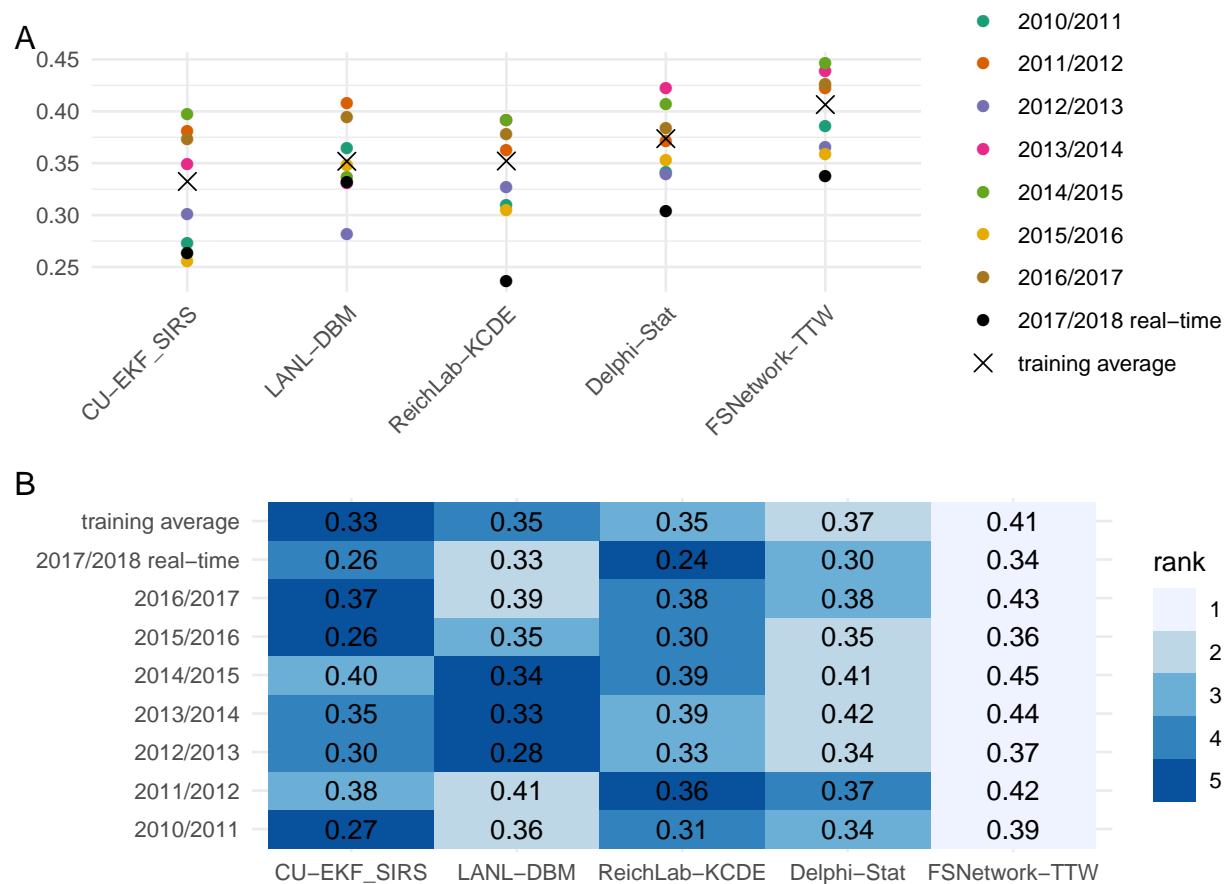


Figure 4: Average forecast score, aggregated across targets, regions, and weeks, plotted separately for selected models and each season. Models shown include the FSNetwork-TTW model, the top performing model from each team during the training phase and, for the last three seasons, the unweighted average of all FluSight models received by CDC. (A) Models are sorted from lowest average scores (left) to highest scores (right). Higher scores indicate better performance. Dots show average scores across all targets, regions, and weeks within a given season. The 'X' marks the geometric mean of the seven seasons. The black dot represents the average score for the prospective, real-time forecasts from the 2017/2018 season. (B) Model ranks for each season and for all training seasons combined. The color of each cell indicates the rank and the forecast score is shown. Note that a component's standalone accuracy does not necessarily correlate to its contribution to the overall ensemble accuracy. See discussion at end of Section efsupsec:comp-models.

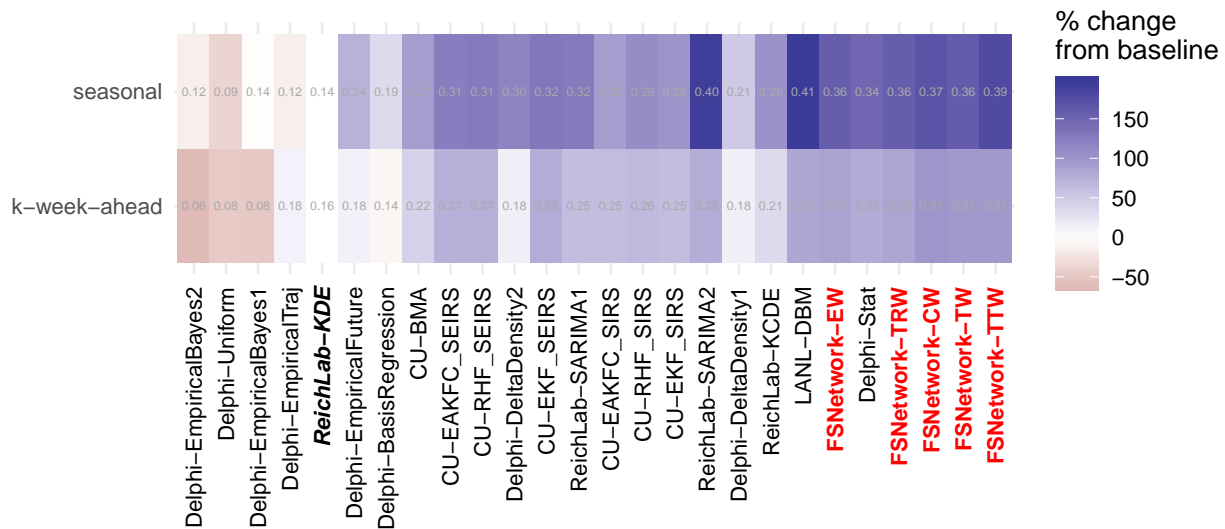


Figure 5: Average forecast score in 2017/2018 by model target-type. The text within the grid shows the score itself. The white midpoint of the color scale is set to be the target-specific average of the historical baseline model, ReichLab-KDE, with darker blue colors representing models that have better scores than the baseline and darker red scores representing models that have worse scores than the baseline. The models are sorted in descending order from most accurate (right) to least accurate (left).

long-term robustness of models such as this one, compared to single models. During the training and test phases, the weighted ensemble approaches considered generally outperformed the FSNetwork Equal Weight ensemble, illustrating the value in incorporating information on prior component performance. As shown by the FSNetwork Target-Type Weights component weighting structure presented above (Figure 3), no one model was ever used by the ensemble as the best answer and it instead relied on a combination of components to optimize performance. Overall, the ensemble methods outperformed the components and while there were some specific situations where the ensemble components did better in our prospective testing seasons, that is unsurprising given the number of models used. An important consideration in this study was that we only passed along one ensemble model into the testing phase and did not pre-specify any comparisons with ensemble components.

A critical limitation to our approach is that, as currently implemented, it relies on having multiple years of past performance for each component to be able to construct a reliable ensemble. This is not a viable method for use in emerging pandemics, where there may not be any historical data on how models have performed nor reliable real-time data to train on. However, preliminary work on “follow-the-leader”-type approaches to dynamically updating the weights, which could remove the requirement that all models have a substantial track-record of performance, has shown some promise, though such approaches still rely on accurately reported real-time data. Furthermore, a simple average of forecasts remains available in such situations and, as illustrated by the relatively strong performance of the FluSight Network Equal Weights model, can still offer advantages over individual models.

One risk of complex ensemble approaches is that they may be “overfit” to the data, resulting in models that place too much emphasis on one approach in a particular scenario or setting. This is a particular concern in applications such as this one, where the number of observations is fairly limited (hundreds to thousands of observations instead of hundreds of thousands). Against this backdrop, the relative simplicity of the FSNetwork Target-Type weights model is a strength, as there is less danger of these models being overfit to the data. Additionally, approaches that use regularization or penalization to reduce the number of effective parameters estimated by a particular

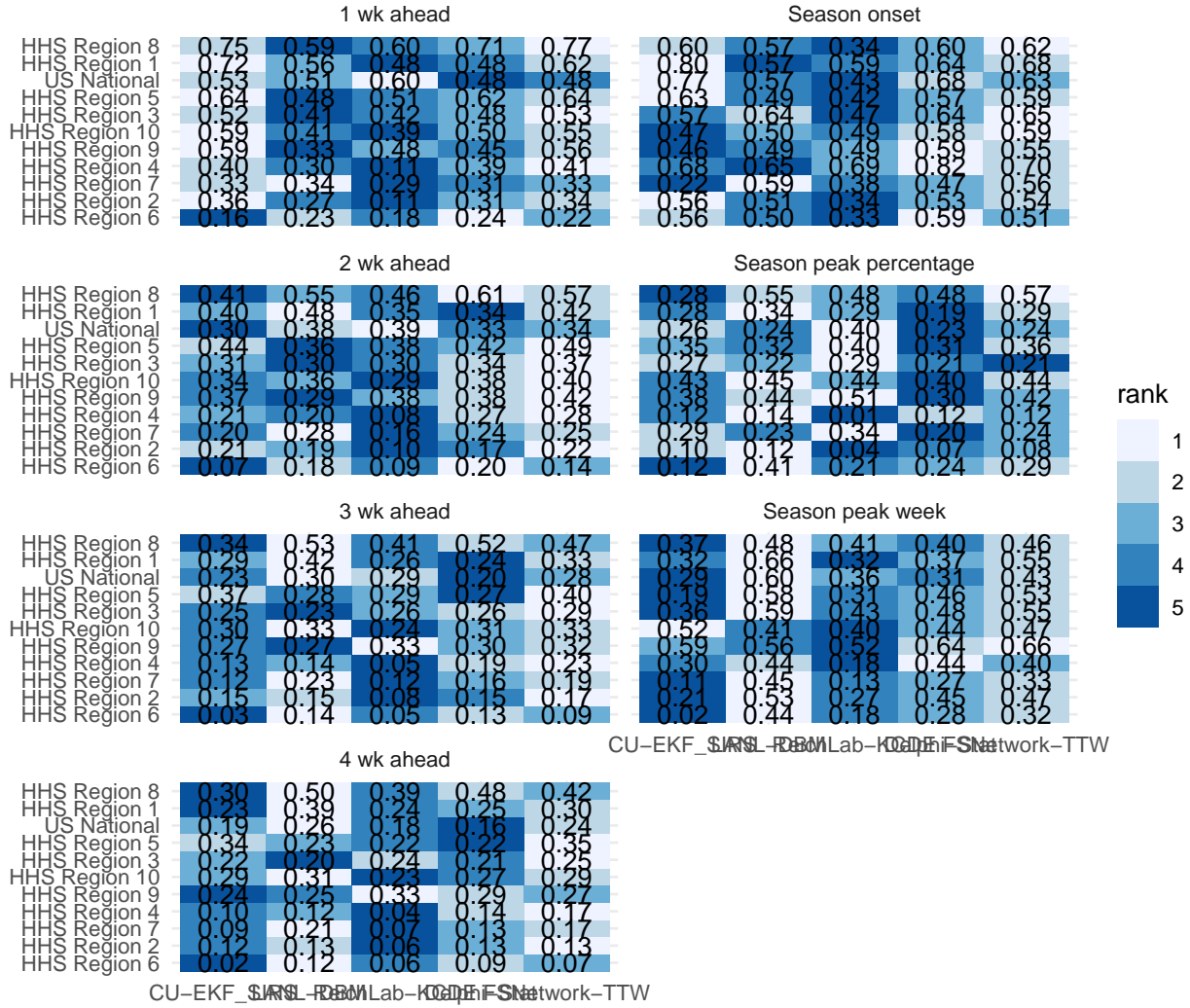


Figure 6: Model scores and ranks by target and region for 2017/2018. Only selected models are shown. Regions are sorted with the most predictable region overall (i.e. highest forecast scores) at the top. Color indicates model rank in the 2017/2018 season. The average forecast score is printed.

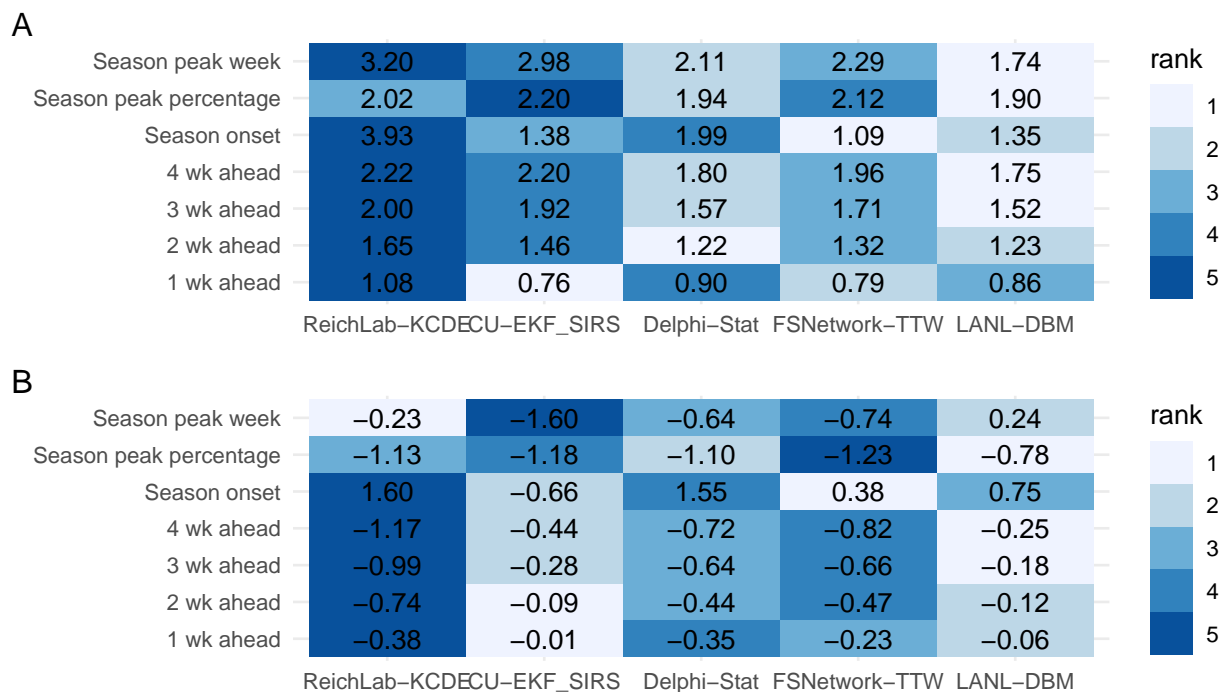


Figure 7: Root mean squared error (RMSE) and bias by target for selected models (with rank) in the 2017/2018 season. Evaluations are for all weeks in the 2017/2018 season. Models are sorted with lowest RMSE on right.

model have been shown to have some practical utility in similar settings and may also have a role to play in future ensembles for infectious disease forecasting.[12]

As this success of this collaborative effort shows, there are significant gains to be made by working across disciplines and research groups, incorporating experts from government, academia and industry. However, at this point, collaborative efforts that consist of pooling results together (in many scientific applications, not just infectious disease forecasting) largely rely on bespoke technological solutions. We built a highly customized solution that relied on GitHub, Travis Continuous Integration server, and model code in R, python, and MatLab. In all, the seven seasons of training data consisted of about 95MB and over 1.5m rows of data per model and about 2GB of forecast data for all models combined. The real-time forecasts for the 2017/2018 season added about 300MB of data. To move ensemble infectious disease forecasting into a more generalizable, operational phase, technological advancements are necessary to both standardize data sources, model structures, and forecast formats as well as develop modeling tools that can facilitate the development and implementation of component and ensemble models.

Public health officials are still learning how to best integrate infectious disease forecasts into real-time decision making. Real-time implementation and testing of forecasting methods plays a central role in planning and assessing what targets should be forecasted for maximum public health impact. Close collaboration between public health policy-makers and quantitative modelers is necessary to ensure that forecasts have maximum impact and are appropriately communicated to the public and the broader public health community. By using the forecasting structure of the CDC's FluSight influenza forecasting challenge, the results of our collaborative ensemble effort were easily shared with officials at CDC. This allowed officials to utilize the ensemble results and develop an understanding of the value of ensemble forecasting approaches. Continuing to work closely with CDC in further

collaborative ensemble efforts will help move the field of infectious disease forecasting forwards and further illustrate the public health utility of ensemble approaches.

## 4 Methods

### 4.1 Influenza Data

Forecasting targets for the CDC FluSight challenge are based on the US Outpatient Influenza-like Illness Surveillance Network (ILINet). ILINet is a syndromic surveillance system that measures the weekly percentage of outpatient visits due to influenza-like illness (ILI) from a network of more than 2,800 providers, and publishes a weighted estimate of ILI (wILI) based on state populations. Estimates of wILI are reported weekly by the CDC's Influenza Division for the United States as a whole as well as for each of the 10 Health and Human Services (HHS) regions. Reporting of 'current' wILI is typically delayed by approximately one to two weeks from the calendar date of a doctor's office visit as data are collected and processed, and each weekly publication can also include revisions of prior reported values if new data become available. Larger revisions have been shown to be associated with decreased forecast accuracy.[\[22\]](#) For the U.S. and each HHS Region, CDC publishes an annual baseline level of ILI activity based on off-season ILI levels.

### 4.2 Forecast Targets and Structure

As the goal was to submit our ensemble forecast in real-time to the CDC FluSight forecasting challenge, we adhered to guidelines and formats set forth by the challenge in determining forecast format. A season typically consists of forecast files generated weekly for 33 weeks, starting with epidemic week 43 (EW43) of one calendar year and ending with EW18 of the following year. Forecasts for the CDC FluSight challenge consist of seven targets: three seasonal targets and four short-term or 'week-ahead' targets (Figure 1B). The seasonal targets consist of season onset, defined as the first MMWR week where wILI is at or above baseline and remains above for three consecutive weeks, season peak week, defined as the MMWR week of maximum wILI, and season peak percentage, defined as the maximum wILI value for the season. The short-term targets consist of forecasts for wILI values 1, 2, 3, and 4 weeks ahead of the most recently published data. With the two-week reporting delay in the publication of ILINet, these forecasts are for the level of wILI occurring 1 week prior to the week the forecast is made, the current week, and the two weeks after the forecast is made (Figure 1B). Forecasts are created for all targets for the US as a whole and for each of the 10 HHS Regions (Figure 1A,C,D).

For all targets, forecasts consist of probability distributions within bins of possible values for the target. For season onset and peak week, forecast bins consist of individual weeks within the influenza season, with an additional bin for onset week corresponding to a forecast of no onset. For short-term targets and peak intensity, forecast bins consist of levels of observed wILI rounded to the nearest 0.1% up to 13%, which is the level of resolution publicly for ILINet reported by the CDC. Formally, the bins are defined as  $[0.00, 0.05)$ ,  $[0.05, 0.15)$ ,  $\dots$ ,  $[12.85, 12.95)$ ,  $[12.95, 100]$ .

### 4.3 Forecast Evaluation

Submitted forecasts were evaluated using the modified log score used by the CDC in their forecasting challenge, which provides a simultaneous measure of forecast accuracy and precision. The log score for a probabilistic forecast  $m$  is defined as  $\log f_m(z^*|\mathbf{x})$ , where  $f_m(z|\mathbf{x})$  is the predicted density function from model  $m$  for some target  $Z$ , conditional on some data  $\mathbf{x}$  and  $z^*$  is the observed value of the target  $Z$ .

While a true log score only evaluates the probability assigned to the exact observed value  $z^*$ , the CDC uses a modified log score that classifies additional values as “accurate”. For predictions of season onset and peak week, probabilities assigned to the week before and after the observed week are included as correct, so the modified log score becomes  $\log \int_{z^*-1}^{z^*+1} f_m(z|\mathbf{x}) dz$ . For season peak percentage and the short-term forecasts, probabilities assigned to wILI values within 0.5 units of the observed values are included as correct, so the modified log score becomes  $\log \int_{z^*-0.5}^{z^*+0.5} f_m(z|\mathbf{x}) dz$ . We refer to these modified log scores as simply log scores hereafter.

Individual log scores can be averaged across different combinations of forecast regions, target, weeks, or seasons. Formally, each model  $m$  has a large number of region-, target-, season-, and week-specific log scores, and we represent a specific scalar log score as  $\log f_{m,r,t,s,w}(z^*|\mathbf{x})$ . These individual log scores can be averaged across combinations of regions, targets, seasons, and weeks to compare model performance.

As other forecasting efforts have used mean square error (MSE) or root mean square error (RMSE) as an evaluation method, we additionally evaluated the prospective forecasts received during the 2017-2018 season using RMSE. The submitted point forecast was used to score each component, and a point forecast was generated for each FSNetwork model by taking the median of the predicted distribution. For each model  $m$ , we calculated  $RMSE_{m,r,t,w}$  for each region  $r$ , target  $t$ , and week  $w$  as  $RMSE_{m,r,t,w} = \sqrt{(\hat{y}_{m,r,t,w} - y_{r,t,w})^2}$ , where  $\hat{y}_{m,r,t,w}$  is the point prediction of model  $m$  for observed value  $y_{r,t,w}$ .

### 4.4 Ensemble components

To provide training data for the ensemble, four teams submitted between 1 and 9 models each, for a total of 21 ensemble components. \*(NOTE: Could refer to comparison manuscript here if it's out)\* Teams submitted out-of-sample forecasts for the 2010/2011 through 2016/2017 influenza seasons. Teams constructed their forecasts in a prospective fashion, using only data that were available at the time of the forecast. For some data sources (i.e. wILI prior to the 2014/2015 influenza season), data as they were published at the time were not available. In such cases, teams were still allowed to use those data sources while making efforts to only use data available at the time forecasts would have been made.

For each influenza season, teams submitted weekly forecasts from epidemic week 40 (EW40) of the first year through EW20 of the following year, using standard CDC definitions for epidemic week (citation). If a season contained EW53, forecasts were submitted for that week as well. In total, teams submitted 233 individual forecast files representing forecasts across the seven influenza seasons. Once submitted, the forecast files were not updated except in four instances where explicit programming bugs had resulted in numerical issues in the forecast. Teams were explicitly discouraged from re-tuning or adjusting their models for different prior seasons to avoid issues with over-fitting.

Teams utilized a variety of methods and modeling approaches in the construction of their submissions. Seven of the models used a compartmental structure in their models (i.e. Susceptible-Infectious-Recovered) to model the

disease transmission process in some way, while other models used more statistical approaches to directly model the observed wILI curve. Six of the models explicitly incorporate additional data sources beyond previous wILI data, including weather data and Google search data.

## 4.5 Distinction between standalone models and ensemble components

It is important to distinguish ensemble components from standalone forecasting models. Standalone models are optimized to be as accurate as possible on their own by, among other things, using proper smoothing. Ensemble components might be designed to be accurate on their own, or else they may be included merely to complement weak spots in other components, i.e. to reduce the ensemble's variance. Because we had sufficient cross-validation data to estimate ensemble weights for several dozen components, some groups contributed non-smoothed "complementing" components for that purpose. Such components may perform poorly on their own, yet their contribution to overall ensemble accuracy may still be significant.

## 4.6 Ensemble Construction

All ensemble models can be represented with the same notation. Let  $f_c(y_{t,r,w})$  represent the predictive density of ensemble component  $c$  for the value of the target  $Y_{t,r,w}$ , where  $t$  indexes the particular target,  $r$  indexes the region, and  $w$  indexes the week. We combine these components together into an ensemble model  $f(y_t)$  as follows:

$$f(y_{t,r,w}) = \sum_{c=1}^C \pi_{c,t,r} f_c(y_{t,r,w}) \quad (1)$$

where  $\pi_{c,t,r}$  is the weight assigned to component  $c$  for predictions of target  $t$  in region  $r$ . We require  $\sum_{c=1}^C \pi_{c,t,r} = 1$  and thereby ensure that  $f(y_{t,r,w})$  remains a valid probability distribution.

A total of five weighted ensemble models were considered, with varying complexity and number of estimated weights (Table 1).

- Equal Weight (FSNetwork-EW):  $\pi_{c,t,r} = 1/C$ . This model consisted of assigning all components the same weight regardless of performance and is equivalent to the arithmetic mean of the components.
- Constant Weight model (FSNetwork-CW):  $\pi_{c,t,r}$  varies across components but has the same value for all targets and regions, for a total of 21 weights.
- Target Type Weight model (FSNetwork-TTW):  $\pi_{c,t,r}$  is estimated separately for short-term targets and seasonal targets with no variation across regions, resulting in a total of 42 weights.
- Target Weight model (FSNetwork-TW):  $\pi_{c,t,r}$  is estimated separately for each of the seven targets for each component with no variation across regions, resulting in 147 weights
- Target-Region Weight model (FSNetwork-TRW): The most complex model considered, which involved estimating  $\pi_{c,t,r}$  separately for each component-target-region combination, resulting in 1617 unique weights.

[EM algorithm/weighting/density stacking]

Component weights were trained using a leave-one-season out cross-validation approach on component forecasts from the 2010/2011 through 2016/2017 seasons. Given the limited number of seasons, we used data from all other seasons as training data to estimate weights for a given test season, even if the training season occurred chronologically after the test season of interest.

Based on the results of the cross-validation study, we selected one ensemble model as the official FluSight Network entry to the CDC's 2017/2018 influenza forecasting challenge. Component weights for that model were estimated using all seven seasons of training data. Over the course of the 2017/2018 influenza season, participating teams submitted weekly forecasts from each component, which were combined using the estimated weights into the FluSight Network model and submitted to the CDC. The component weights for the submitted model were unchanged throughout the course of the season.

It should be noted that ensemble weights are not a measure of ensemble components' standalone accuracy nor do they measure the overall contribution of a particular model to the ensemble accuracy. For example, consider a setting where a duplicate of a identical (or highly similar) ensemble component with weight  $\pi^*$  is added to a given ensemble. The accuracy of the original ensemble can be maintained in a number of ways, including (a) assigning each copy a weight of  $\pi^*/2$ , or (b) assigning the first copy a weight of  $\pi^*$  and the second copy a weight of 0. In both of these weightings, at least one high accuracy ensemble component would be assigned significantly lower weight based on the presence of another identical or similar component. Additionally, components can be assigned small weights but have a large impact on ensemble accuracy compared to an ensemble excluding it.

This effort shows that collaborative efforts between research teams to develop ensemble forecasting approaches bring measurable improvements in accuracy and reductions in variability. We therefore are moving substantially closer to forecasts that can and should be used to inform routine, ongoing public health surveillance of infectious diseases. With the promise of new, real-time data sources and continued methodological innovation for both component models and ensemble approaches, there is good reason to believe that infectious disease forecasting will continue to mature and improve in upcoming years. As modeling efforts become more commonplace in the support of public health decision-making worldwide, it will be critical to develop infrastructure so that multiple models can more easily be brought online to create ensemble forecasts as well as to develop our understanding of how best to communicate the forecasts and their uncertainty to decision-makers and the general public. Efforts such as this, that emphasize real-time testing and evaluation of forecasting models and facilitate the close collaboration between public health officials and modeling researchers, are critical to improving our understanding of how best to use forecasts to improve public health response to seasonal and emerging epidemic threats.

## References

- [1] Estimated influenza illnesses, medical visits, hospitalizations, and deaths averted by vaccination in the united states. <https://www.cdc.gov/flu/about/disease/2016-17.htm>, May 2018.
- [2] Overview of influenza surveillance in the united states. <https://www.cdc.gov/flu/weekly/overview.htm>, 2017.
- [3] Jeffrey Shaman, Alicia Karspeck, Wan Yang, James Tamerius, and Marc Lipsitch. Real-time influenza forecasts during the 2012–2013 season. *Nature Communications*, 4, 2013.



- [4] Wan Yang, Alicia Karspeck, and Jeffrey Shaman. Comparison of Filtering Methods for the Modeling and Retrospective Forecasting of Influenza Epidemics. *PLoS Computational Biology*, 10(4):e1003583, apr 2014.
- [5] Shihao Yang, Mauricio Santillana, and S C Kou. Accurate estimation of influenza epidemics using Google search data via ARGO. *Proceedings of the National Academy of Sciences of the United States of America*, 112(47):14473–8, nov 2015.
- [6] Jean-Paul Chretien, Dylan George, Jeffrey Shaman, Rohit A. Chitale, and F. Ellis McKenzie. Influenza forecasting in human populations: A scoping review. *PLOS ONE*, 9(4):1–8, 04 2014.
- [7] Sasikiran Kandula, Daniel Hsu, and Jeffrey Shaman. Subregional Nowcasts of Seasonal Influenza Using Search Trends. *Journal of medical Internet research*, 19(11):e370, nov 2017.
- [8] Dave Osthus, James Gattiker, Reid Priedhorsky, and Sara Y. Del Valle. Dynamic Bayesian Influenza Forecasting in the United States with Hierarchical Discrepancy. *arXiv*, aug 2017.
- [9] Logan C. Brooks, David C. Farrow, Sangwon Hyun, Ryan J. Tibshirani, and Roni Rosenfeld. Nonmechanistic forecasts of seasonal influenza with iterative one-week-ahead distributions. *PLOS Computational Biology*, 14(6):e1006134, jun 2018.
- [10] Sen Pei, Sasikiran Kandula, Wan Yang, and Jeffrey Shaman. Forecasting the spatial transmission of influenza in the United States. *Proceedings of the National Academy of Sciences of the United States of America*, 115(11):2752–2757, mar 2018.
- [11] Teresa K. Yamana, Sasikiran Kandula, and Jeffrey Shaman. Individual versus superensemble forecasts of seasonal influenza outbreaks in the United States. *PLOS Computational Biology*, 13(11):e1005801, nov 2017.
- [12] Evan L. Ray and Nicholas G. Reich. Prediction of infectious disease epidemics via weighted density ensembles. *PLOS Computational Biology*, 14(2):e1005910, feb 2018.
- [13] Morgan E Smith, Brajendra K Singh, Michael A Irvine, Wilma A Stolk, Swaminathan Subramanian, T Déirdre Hollingsworth, and Edwin Michael. Predicting lymphatic filariasis transmission and elimination dynamics using a multi-model ensemble framework. *Epidemics*, 18:16–28, 2017.
- [14] Cécile Viboud, Kaiyuan Sun, Robert Gaffey, Marco Ajelli, Laura Fumanelli, Stefano Merler, Qian Zhang, Gerardo Chowell, Lone Simonsen, and Alessandro Vespignani. The RAPIDD ebola forecasting challenge: Synthesis and lessons learnt. *Epidemics*, aug 2017.
- [15] DARPA. CHIKV Challenge Announces Winners, Progress toward Forecasting the Spread of Infectious Diseases. <https://www.darpa.mil/news-events/2015-05-27>, 2015.
- [16] NOAA and CDC. Dengue Forecasting.
- [17] Matthew Biggerstaff, Krista Kniss, Daniel B Jernigan, Lynnette Brammer, Joseph Bresee, Shikha Garg, Erin Burns, and Carrie Reed. Systematic assessment of multiple routine and near-real time indicators to classify the severity of influenza seasons and pandemics in the united states, 2003–04 through 2015–2016. *American Journal of Epidemiology*, 187:1040–1050, 2018.
- [18] Sen Pei and Jeffrey Shaman. Counteracting structural errors in ensemble forecast of influenza outbreaks. *Nature Communications*, 8(1):925, dec 2017.

- [19] Logan C Brooks, David C Farrow, Sangwon Hyun, Ryan J Tibshirani, and Roni Rosenfeld. epiforecast: Tools for forecasting semi-regular seasonal epidemic curves and similar time series. <https://github.com/cmu-delphi/epiforecast-R>, 2015.
- [20] Logan C. Brooks, David C. Farrow, Sangwon Hyun, Ryan J. Tibshirani, and Roni Rosenfeld. Flexible Modeling of Epidemics with an Empirical Bayes Framework. *PLOS Computational Biology*, 11(8):e1004382, aug 2015.
- [21] Evan L. Ray, Krzysztof Sakrejda, Stephen A. Lauer, Michael A. Johansson, and Nicholas G. Reich. Infectious disease prediction with kernel conditional density estimation. *Statistics in Medicine*, sep 2017.
- [22] Nicholas G Reich, Logan Brooks, Spencer Fox, Sasikiran Kandula, Craig McGowan, Evan Moore, Dave Osthus, Evan L Ray, Abhinav Tushar, Teresa Yamana, Matthew Biggerstaff, Michael A Johansson, Roni Rosenfeld, and Jeffrey Shaman. Forecasting seasonal influenza in the u.s.: A collaborative multi-year, multi-model assessment of forecast performance. *bioRxiv*, 2018.
- [23] David C Farrow, Logan C Brooks, Sangwon Hyun, Ryan J Tibshirani, Donald S Burke, and Roni Rosenfeld. A human judgment approach to epidemiological forecasting. *PLoS computational biology*, 13(3):e1005248, 2017.