

Evaluation of an open forecasting challenge to assess skill of West Nile virus neuroinvasive disease prediction

Karen M Holcomb (✉ kholcomb@cdc.gov)

National Oceanic and Atmospheric Administration

Sarabeth Mathis

Centers for Disease Control and Prevention

J Erin Staples

Centers for Disease Control and Prevention

Marc Fischer

Centers for Disease Control and Prevention

Christopher M Barker

University of California, Davis

Charles B Beard

Centers for Disease Control and Prevention

Randall J Nett

Centers for Disease Control and Prevention

Alexander C Keyel

Wadsworth Center

Matteo Marcantonio

Université Catholique de Louvain

Marissa L Childs

Stanford University

Morgan E Gorris

Los Alamos National Laboratory

Ilia Rochlin

Rutgers, The State University of New Jersey

Marco Hamins-Puértolas

University of California, San Francisco

Evan L Ray

Mount Holyoke College

Johnny A Uelmen

University of Illinois Urbana-Champaign

Nicholas DeFelice

Icahn School of Medicine at Mount Sinai

Andrew S Freedman

North Carolina State University

Brandon D Hollingsworth

Cornell University

Praachi Das

North Carolina State University

Dave Osthus

Los Alamos National Laboratory

John M Humphreys

Agricultural Research Service

Nicole Nova

Stanford University

Erin A Mordecai

Stanford University

Lee W Cohnstaedt

Agricultural Research Service

Devin Kirk

Stanford University

Laura D Kramer

Wadsworth Center

Mallory J Harris

Stanford University

Morgan P Kain

Stanford University

Emily MX Reed

Virginia Tech

Michael A Johansson

Centers for Disease Control and Prevention

Research Article

Keywords: : calibration, discriminatory power, forecasting, logarithmic score, multi-model assessment, West Nile virus, West Nile neuroinvasive disease, United States

Posted Date: August 26th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1992050/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: No competing interests reported.

Version of Record: A version of this preprint was published at Parasites & Vectors on January 12th, 2023.

See the published version at <https://doi.org/10.1186/s13071-022-05630-y>.

Abstract

Background: West Nile virus (WNV) is the leading cause of mosquito-borne illness in the continental United States. WNV occurrence has high spatiotemporal variation and current approaches for targeted control of the virus are limited, making forecasting a public health priority. However, little research has been done to compare strengths and weaknesses of WNV disease forecasting approaches on the national scale. We used forecasts submitted to the 2020 WNV Forecasting Challenge, an open challenge organized by the Centers for Disease Control and Prevention, to assess the status of WNV neuroinvasive disease (WNND) prediction and identify avenues for improvement.

Methods: We performed a multi-model comparative assessment of probabilistic forecasts submitted by 15 teams for annual WNND cases in US counties for 2020, and assessed forecast accuracy, calibration, and discriminatory power. In the evaluation, we included forecasts produced by comparison models of varying complexity as benchmarks of forecast performance. We also used regression analysis to identify modeling approaches and contextual factors that were associated with forecast skill.

Results: Simple models based on historical WNND cases generally scored better than more complex models and combined higher discriminatory power with better calibration of uncertainty. Forecast skill improved across updated forecast submissions submitted during the 2020 season. Among models using additional data, inclusion of climate or human demographic data was associated with higher skill, while inclusion of mosquito or land use data was associated with lower skill. We also identified population size, extreme minimum winter temperature, and interannual variation in WNND cases as county-level characteristics associated with variation in forecast skill.

Conclusions: Historical WNND cases were strong predictors of future cases with minimal increase in skill achieved by models that included other factors. Although opportunities might exist to specifically improve predictions for areas with large populations and low or high winter temperatures, areas with high case-count variability are intrinsically more difficult to predict. Also, the prediction of outbreaks, which are outliers relative to typical case numbers, remains difficult. Further improvements to prediction could be obtained with improved calibration of forecast uncertainty and access to real-time data streams (e.g., current weather and preliminary human cases).

Background

West Nile virus (WNV; *Flaviviridae*, *Flavivirus*) is the leading cause of mosquito-borne illness in the continental United States [1]. Symptomatic infections typically present as a febrile illness (approximately 20% of all infections). However, < 1% of all infections result in West Nile neuroinvasive disease (WNND) with manifestations including meningitis, encephalitis, or acute flaccid paralysis [2]. WNV was first detected in the United States in 1999 [3] and by 2005, had spread across the contiguous United States and up the Pacific coast [4]. From 1999–2020, the Centers for Disease Control and Prevention (CDC) reported a total of 26,683 non-neuroinvasive WNV disease cases and 25,849 WNND cases, resulting in

2,456 deaths [5]. During 2005–2020, a median of 409 (range 167–693) of 3,108 counties in the contiguous United States reported WNND cases each year. Even in counties that regularly report WNND cases, the number and location of WNND cases varies. Large spatial and temporal heterogeneity in annual WNND cases make accurate prediction of incidence both challenging and potentially valuable to guide prevention and control efforts.

The ecology of WNV is complex and spatially variable across the United States. The virus is maintained in an enzootic cycle between birds (predominantly passerines) and *Culex* mosquitoes [6–8], but can cause disease in horses and humans, which are dead-end hosts [9]. The vectors for WNV vary geographically [8]. In the east-central region (northeast, mid-Atlantic, and central United States), *Cx. pipiens* and *Cx. restuans* have been incriminated as the primary vectors with *Cx. salinarius* also playing an important role in maintenance and zoonotic transmission in coastal areas. In the southeast, *Cx. quinquefasciatus* has been implicated as the primary vector with *Cx. salinarius* and *Cx. nigripalpus* also capable of causing human disease. In western North America, *Cx. tarsalis* is largely responsible for zoonotic transmission, especially in more rural areas, while *Cx. pipiens* serves as the enzootic vector in urban areas in the more northern parts of the western United States (northern Great Plains, Rocky Mountains, and Pacific Northwest). In urban areas of the southwestern United States, *Cx. quinquefasciatus* can act as the dominant zoonotic vector. Other *Culex* mosquito species can have a secondary or localized importance in this region.

Meteorological factors like temperature and precipitation have a large impact on the transmission of WNV. Temperature influences mosquito survival and potential WNV transmission rates [10]. As temperatures warm, mosquito development and biting rates accelerate [10,11]. Additionally, with increasing temperature, the extrinsic incubation period for WNV decreases as viral replication rates increase [12–15]. Thus, with increasing temperature above the thermal minimum for mosquito survival and WNV replication [14,16], viral transmission and risk of zoonotic transmission increases. However, there is a thermal optimum (23.9–25.2°C [17]) above which transmission generally decreases due to negative impacts on mosquito survival and other traits. Variation in the interaction of climatic and landscape factors contributes to seasonal dynamics and spatial variation in the effect of temperature [8,18]. Increased precipitation generally increases the quantity of available larval habitat [19–21], but intense precipitation events can wash out immature mosquitoes from larval habitat such as catch basins [22]. The impact of precipitation varies broadly across the United States with a positive association between increased precipitation and above average-human cases in the western United States, but a negative association in the eastern United States, potentially due to difference in the mosquito species, their preferred egg-laying habitats, and other environmental factors present in each area [8,18,21]. Also, drought has been associated with WNV amplification and increased human cases, partially due to aggregation of hosts and vectors at dwindling water sources [23,24].

Statistical and mechanistic models have been developed to predict geographic or temporal dynamics of WNV transmission [25,26]. Models generally produce estimates on a single spatial and temporal scale aimed at guiding public health decisions or elucidating factors that enable increased transmission.

Models developed for prediction in one location often fail to perform well if applied to a different location due to variation in factors like ecology, primary mosquito species, and human behavior as well as availability of predictor data, like mosquito surveillance data [27]. Out-of-sample validation is often used to assess model performance, but no multi-model comparative assessment has been performed to assess the strengths and weaknesses of predictive WNV modeling at the local or national scale.

To systematically evaluate WNNND prediction across the continental United States, the CDC Epidemic Predictive Initiative and the Council for State and Territorial Epidemiologists launched an open West Nile virus Forecasting Challenge in 2020. The primary objective of the challenge was to predict the total number of WNNND cases for each county in the contiguous United States that would be reported to the national surveillance system for arboviral diseases, ArboNET, during the 2020 calendar year. In our evaluation of the Challenge, we 1) assessed whether some models had better predictive performance than others, 2) identified modeling approaches associated with better prediction, and 3) evaluated contextual factors of the counties (e.g., environmental, climatic, and historical WNV patterns) associated with variation in forecast skill.

Methods

Team participation

An announcement recruiting team participation in the 2020 WNV Forecasting Challenge was circulated widely by the CDC Epidemic Prediction Initiative through emails and postings on webpages starting in March 2020. Teams using any modeling approach were encouraged to participate.

Participating teams signed a data use agreement and were provided with annual WNNND case counts, by county for the contiguous United States and Washington DC during 2000–2018, from ArboNET, the national arboviral diseases surveillance system administered by the CDC. Provisional 2019 case data were provided to participants in early May 2020. Participants were allowed to use any other data source, like climate, weather, land use, mosquito surveillance, and human demographics, to develop their modeling approach.

Forecasting target

Teams predicted the total number of probable and confirmed WNNND cases that would be reported to ArboNET for all counties ($n = 3,108$) in the contiguous United States and Washington DC during 2020. WNNND cases were chosen as the outcome because the severe manifestations of the disease are more likely to be consistently recognized and reported compared with less severe, non-neuroinvasive WNV disease cases [28].

For each location, a forecast included both a point estimate and a binned probability distribution. The point estimate denoted the most likely number of cases. Fifteen bins were chosen to cover the range of cases from 0 to > 200 , with finer resolution for smaller numbers of expected cases (i.e., bins for 0, 1–5, 6–

10, ..., 46–50, 51–100, 101–150, 151–200, >200 cases). Teams assigned a probability between 0 and 1 to each bin, with a total probability equal to 1.0 across all bins per county.

Forecasts

The initial forecast due date was April 30, 2020, with submission to an online system (<https://predict.cdc.gov>). Additional, optional, updated submissions could be submitted by the following deadlines: May 31, June 30, and July 31, 2020. Further details are available through the online submission system.

Concurrently, we developed four additional models for comparison with the team forecasts: a naïve model, an always-absent model, a negative binomial model, and an ensemble model. The naïve model assigned equal probability to each of the bins (i.e., 1/15 probability). The always-absent model represented a universal expectation of zero cases by assigning a probability of 1.0 to the zero-case bin and zero probability to all other bins for each county. The negative binomial model was built to reflect a parsimonious probabilistic prediction relying exclusively on local historical data. For each county, we fitted a negative binomial distribution to historical WNND cases and extracted probabilities for each bin from the cumulative distribution function. The initial version of this forecast (April submission) used 2000–2018 case counts, while the May submission also incorporated the provisional 2019 data reported as of May 2020. Finally, we created a mean consensus ensemble using all team-submitted forecasts and the negative binomial forecast by averaging the probability bins for all forecasts at each location and submission deadline. For forecasts that were not updated, we used the last available forecast for each update of the ensemble.

We developed two additional models retrospectively as alternative baseline models: a first-order autoregressive model (i.e., AR(1)) and a first-order autoregressive model with a climate covariate (AR(1) Climate). For both models, we fitted log-transformed annual WNND case counts (2005–2019) using the *arima* function in the stats package in R (version 4.1.2; [29]). For the AR(1) Climate model, we considered seasonal aggregations of climate conditions (average temperature, mean minimum temperature, and total precipitation), using Parameter-elevation Regressions on Independent Slopes Model (PRISM) data [30]. We defined seasons as three-month periods for winter (Dec-Feb), spring (Mar-May), summer (Jun-Aug), and fall (Sep-Nov), using climate data from the previous winter to the concurrent year's spring to predict annual WNND case numbers (e.g., using seasonal climate data for Dec 2018-May 2020 to predict 2020 WNND cases). See Additional File 1: Text S1 for more details on the development of the autoregressive modeling framework.

Evaluation

As announced before the Challenge, we evaluated all forecasts using the logarithmic score, a proper scoring rule based on the probabilities assigned in each forecast in relation to the eventual observed case counts [31,32]. The score for each team was the average logarithm of the probability assigned to the observed outcome bin, the bin containing the reported number of WNND cases for 2020, per county. To avoid logarithmic scores of negative infinity for forecasts which assigned zero probability to the observed

outcome, we truncated binned predictions to have a minimum logarithmic score of -10. We compared mean logarithmic scores with ANOVA followed by Tukey post-hoc multiple comparisons to identify significant differences between forecast scores. We compared the forecasts for the final versions of team forecasts and comparison models, and between the initial and final versions of all forecasts.

We assessed probabilistic calibration by plotting forecasted probabilities versus observed frequencies for forecasts with each summarized in the following upper-bound inclusive probability bins: 0.0, 0.0-0.1, 0.1-0.2, ..., 0.9-1.0. We then calculated a metric of overall probabilistic calibration as the mean weighted squared difference of binned predicted probabilities versus the observed frequency of events;

$\frac{1}{N} \sum n_k (\bar{p}_k - \bar{o}_k)^2$, where N is the total number of a team's prediction, n_k is the number of predictions in bin k (e.g., between 0.2 and 0.3) with average probability \bar{p}_k , and \bar{o}_k is the frequency of those predictions being correct. Note that this considers calibration within the single forecast year, and provides no information on calibration of models across forecast years.

To assess discriminatory power, we used receiver-operator characteristic (ROC) curve analysis to assess the sensitivity and specificity of the probability of having at least one WNND case in each county. We then calculated the area under the curve (AUC) as the metric for discrimination.

Regression modeling

We used Bayesian regression modeling to identify high-level modeling approaches and contextual factors of counties associated with variation in skill. To assess the impact of modeling approach, we fitted generalized linear models to all team forecasts and the negative binomial comparison model (April and May versions) using the negative logarithmic score, or surprisal, as the outcome, assuming a Gamma distribution with the inverse link. We used the *stan_glm* function in the *rstanarm* package (version: 2.21.1, [33]) to fit the models. We assessed associations between surprisal and a suite of model-specific nominal covariates for a team's inclusion of data on climate, human demographics, land use, mosquito distributions/surveillance, and bird/equine infections, and if submissions were updated. To assess county-specific contextual factors, we fitted Bayesian generalized additive models (GAMs) to the ensemble forecasts using the *stan_gamm4* function in the *rstanarm* package (version: 2.21.1, [33]). We chose the ensemble forecast to capture the overall accuracy of all teams without the variation in performance between teams due to modeling approaches. Contextual factors investigated included environmental factors (e.g., land use, extreme minimum winter temperature, region), history of reported WNND cases (e.g., number of years and pattern of reported cases), and demographics (e.g., population size, population density, population > 65 years old). See Additional File 1: Text S1 for more details on methods, model selection, and a complete list of variables considered.

Results

Fifteen teams submitted binned probabilistic forecasts for the total number of WNND cases reported in each county using a variety of modeling approaches (see Additional File 1: Text S1 for team information

including model details and descriptions and Table S1 for model characteristics). Two teams (13%) included mechanistic model elements while the remainder used completely statistical approaches. Six teams (40%) used Bayesian frameworks for model fitting. We broadly categorized the modeling approaches teams used as machine learning (i.e., random forest, neural network), regression (i.e., maximum likelihood generalized linear models, generalized additive models), hurdle models (i.e., spatio-temporal hurdle models fit using integrated nested Laplace estimation), system of difference equations, or historical case distributions. Across the four submission timepoints, we received 30 unique forecast submissions (15 initial submissions, 5 teams that updated once, 2 that updated twice, and 2 that updated three times). Some teams used different data sources in different submissions. Across all submissions, 24 submissions (from 11 teams) used climatic data, 22 (from 11 teams) used human demographic data, 9 (from 5 teams) used land-use data, 12 (from 4 teams) used entomological data related to *Culex* mosquito species distributions or WNV infection prevalence in mosquitoes, 2 used data on avian WNV infections (1 team), and 2 used data on equine WNV infections (1 team).

The ensemble model assigned the highest probability to a non-zero bin for 115 counties, with the largest probabilities assigned to high numbers of WNND cases in highly urbanized counties: Los Angeles (CA, bin: 101–150 cases), Maricopa (AZ, bin: 51–100 cases), Cook (IL, bin: 51–100 cases), and Harris (TX, bin: 11–15 cases) (Fig. 1A); the other 111 counties assigned the highest probability to the 1–5 cases bin. The remaining 2,993 counties had the highest probability assigned to the zero-case bin. Uncertainty in ensemble predictions was greatest in more populous counties as well as in the southwest (CA, AZ, NV), in the Great Plains states, along the southern edges of the Great Lakes, and along the northeast coast (Fig. 1B).

Finalized case data for 2020 were released in November 2021 with 559 WNND cases reported in 181 counties. These counts were similar to totals reported annually during 2008–2011 and 2019 (Additional File 1: Table S2). The ratio of reported neuroinvasive to non-neuroinvasive cases was 3.25, the largest reported since 2001 (range for 2002–2019: 0.41–2.43).

Forecast skill, as measured by logarithmic score, generally increased across the submission timepoints with updated submissions (Fig. 2, Additional File 1: Table S3). Gains in skill for individual forecasting teams were typically abrupt and occurred at different times, presumably due to acquisition of new contextual data or revisions of modeling approaches. The ensemble forecast, which included all the most recent team forecasts and the negative binomial model at each time point, increased from a mean log score of -0.357 (April) to -0.253 (July), with the largest increase in skill occurring between the June and July submissions likely due to the dramatic improvement in the forecast by *UI*. Three teams (*MSSM*, *Stanford*, and *UNL*) and the negative binomial forecast consistently outscored the ensemble forecast with four teams (*MHC*, *NYSW*, *NYSW-CVD*, and *UCD*) outscoring the ensemble for at least one submission timepoint. The retrospectively implemented AR(1) and AR(1) Climate models (using mean winter temperature based on historical performance, Additional File 1: Fig S1) also consistently outperformed the ensemble.

Overall, models based only on historical distributions of cases had relatively high skill. The negative binomial comparison model, AR(1) comparison model, and an empirically weighted distribution (*MSSM*) were in the top five forecasts at each submission timepoint. Only the final forecast from *UCD* scored higher than the negative binomial model with a difference in mean logarithmic score of 0.007 ($P = 0.98$, Additional File 1: Fig S4).

Comparing high-level modeling approaches and controlling for submission date, we found variation in forecast skill was associated with the inclusion of some types of data (Additional File 1: Table S4). Skill was higher for teams that included climate (0.187, 95% CI: 0.174, 0.226) or demographic data (0.335, 95% CI: 0.326, 0.361). We found lower skill for forecasts that included land use (-0.100, 95% CI: -0.124, -0.031) or *Culex* mosquito geography (estimated ranges or WNV infection prevalence data, -0.114, 95% CI: -0.142, -0.048). We did not compare the association of skill with the inclusion of avian or equine WNV disease cases because only one team used each of these data types.

We next analyzed county-specific contextual factors that might be associated with varying forecast skill across modeling approaches by analyzing associations with ensemble forecast skill (Additional File 1: Fig S3). Average skill was highest in counties with mid-sized populations, low historical variation in annual WNND cases (permutation entropy), and relatively moderate winter minimum temperatures (-10° and 10°F, corresponding to the USDA Plant Hardiness Zones 6a to 7b). For extreme minimum winter temperatures, the ensemble had lower skill at extreme high and low values. For population size, the ensemble had lower skill at large sizes and a nonsignificant relationship at small sizes. Increased variation in interannual historic WNND cases (larger permutation entropy) was associated with decreased forecast skill with a plateau at permutation entropy above approximately 0.7.

Forecast calibration and the ability to predict whether WNND cases would occur (≥ 1 vs. 0 cases, i.e., discrimination) varied across teams (Fig. 3). Comparing forecasts to observations after binning by the forecasted probabilities (Additional File 1: Fig S5), we found that most forecasts were over-confident at lower probabilities and under-confident at higher probabilities. Expectations of the occurrence of cases, especially large numbers of cases, were commonly assigned low probabilities while the expectation of no reported cases was typically highly probable. The forecasts with the best calibration (i.e., reliable specification of probabilities) were those that did not assign any high probabilities (e.g., the naïve forecast), followed by the autoregressive (AR(1) and AR(1) Climate) and negative binomial models. We found that the discriminatory power of forecasts, assessed as the AUC comparing the probability of one or more cases in each county to whether at least one WNND case was reported, also varied widely across teams and comparison models (range of forecast AUC: 0.5-0.875, Additional File 1: Fig S6). The naïve and always-absent comparison models had the worst discriminatory performance, while the ensemble, the negative binomial, the AR(1), the AR(1) Climate forecasts, and several teams (*MHC*, *MSSM*, *NYSW*, *NYSW-CVD*, *Rutgers*, *Stanford*, and *UCD*) all had high discriminatory power. The forecasts with the highest overall skill combined good calibration and discrimination.

Discussion

Reliable early-warning of vector-borne disease outbreaks could offer new opportunities for effective prevention and control through targeting control to high-risk areas. We performed a multi-model evaluation of probabilistic forecasts for the total WNND cases reported by county in the contiguous United States and Washington DC in 2020. The comparison of forecast performance elucidated the current predictive capacity of WNND on this spatial and temporal scale, and avenues for improvement.

Although the COVID-19 pandemic caused dramatic changes in human behavior and challenges for health systems in 2020, it is not clear that the occurrence and reporting of WNND cases changed dramatically. The reported total number of WNND cases was similar to prior years with relatively low case numbers. The ratio of reported WNND to non-neuroinvasive cases for 2020 increased substantially, to the highest level since 2001, indicating likely under-detection and reporting of non-neuroinvasive cases. However, it remains unclear what impact COVID-19 may have had on human behavior and resulting exposure to WNV, treatment-seeking by infected individuals, or physicians' diagnosis and reporting of WNV disease.

Overall, simple models based on historical WNND cases (i.e., the negative binomial model) generally scored better than more complex models, combining discriminatory power and calibration of uncertainty. Only one team (*UCD*) had higher forecast skill than the negative binomial forecast model, and only by a small, nonsignificant margin. One explanation for the relatively strong performance of the negative binomial model is that the historical case distributions reflect the ecological differences across counties and therefore capture most of the inherent spatial variability in WNV transmission. Incorporating additional contextual factors explicitly might not necessarily improve prediction accuracy despite their importance. Also, matching case locations in space and time with available environmental data can introduce uncertainty in model predictions that consider environmental data on top of historical WNV data. For example, WNND data were available on the county-annual scale while environmental data were available at much finer spatial and temporal resolutions. Thus, decisions on aggregations or summaries of environmental data cannot fully capture the particular sequence of conditions precipitating zoonotic transmission.

Regressions to identify modeling approaches associated with variation in forecast skill confirmed an increase in score for later submissions after accounting for other differences. Changes in later forecast submissions were attributed largely to integration of updated data rather than changes in forecasting methods, so this score improvement highlights the value of including updated covariate data (e.g., reported updates included using recent weather data, newly released 2019 WNV data, and additional demographic data). Although we could not discern the relative contribution of each update on the change in score due to heterogeneity in the type of changes and number of submissions across teams, recent weather data appeared to have played some role in improving the predictive accuracy of forecasts. Improving access to real-time data streams could therefore improve predictive accuracy [26,34]. Moreover, these updates occurred before the majority of WNND cases were reported, indicating that although forecasts that provide early warning during the spring can allow for greater lead times for preventative actions, later updates that provide early detection of risk—even after some cases have begun to occur—could provide additional value [26]. From a practical standpoint, shifting forecast submission deadlines

by several days later could facilitate incorporating monthly aggregated data from the prior month when available.

The limited number of submissions prevented us from fully assessing the relative performance of different modeling approaches as models used different data inputs in addition to different methods. While the broad classifications we used provide some insight on general forecast skill, we could not assess the performance of specific model constructions because they varied in both methods and covariates included. It could be of interest to identify variation in predictive performance due to specific model constructions to guide the development and refinement of WNV prediction.

We found the inclusion of estimated mosquito distributions or mosquito surveillance data reduced forecast skill on average. This result seems counter-intuitive because the importance of key mosquito vectors and the relationship between entomological indicators of risk and WNV activity is clear [8,9,35–38]. One explanation is that mosquitoes are much more widespread than WNV cases, so it is difficult to discriminate counties with intense enzootic transmission without human involvement. An alternative explanation is that this finding might reflect model-specific limitations in how the data were incorporated or limited quality or availability of national datasets on mosquito distributions or entomological surveillance. Current distribution maps date back to the 1980s [39,40] with an update in 2021 using habitat suitability modeling [41]. Although the updated maps have increased spatial definition compared to earlier estimates, these distributions indicate relative habitat suitability rather than presence or absence. One publicly available surveillance database, ArboNET, maintains data on human disease and infections among presumptive viremic blood donors, veterinary disease cases, mosquitoes, dead birds, and sentinel animals for a variety of arboviruses. However, nonhuman arboviral surveillance is voluntary with large variation in spatial and temporal coverage between jurisdictions, and reported data are often incomplete [42] reducing the predictive utility of the database. The ensemble forecast had a higher forecasting skill (average logarithmic score) than most team forecasts, with better discriminatory power (ability to differentiate having at least one case) than any team forecast and better calibration (reliable uncertainty specification) than most. Previous forecasting efforts for influenza, dengue, and COVID-19 [43–46] demonstrated that ensemble approaches capitalize on the strengths of diverse models and balance uncertainty across modeling approaches to produce robust predictions. This general finding was replicated here. However, we also found a simple model based on historical data alone substantially outperformed both the ensemble and majority of team forecasts at every submission date for the 2020 Challenge. This indicates that even the strengths of a multi-modeling approach were not sufficient to improve prediction beyond historical trends for this year. However, we did not identify if forecasts performed better than others regionally. If we had weighted the ensemble based on regional performance, this might have improved the skill of the ensemble.

We found that heterogeneity in historic WNV cases had a significant impact on variation in forecast skill, and unsurprisingly, forecasts scored worse in locations of high historic heterogeneity. Improvement in forecast skill for these locations would likely be the most useful for vector control and public health officials, but the high variability also represents a significant challenge to forecasters.

Other intrinsic differences between counties associated with lower forecast skill could highlight areas that need improvement. By identifying local drivers in counties with relatively large populations and hotter or colder winters, forecast skill could be improved in these circumstances. For example, the ecological setting (i.e., *Culex* species present, composition of avian community, climate) would vary substantially between counties with “hot” or “cold” winter extremes and different drivers may need to be considered in each. Also, factors might interact together to impact zoonotic transmission, but due to the limited data and limited number of forecasts available for analysis, we were unable to investigate these.

Calibration across teams indicated other avenues for improving prediction. Overall, teams over-predicted the probability that cases would occur while correspondingly underestimating the probability that cases would not occur. Overestimating the probability of disease cases could lead to better preparedness but could also result in allocation of resources that are not ultimately needed. Moreover, repeated instances of non-events could lead public health officials or the public to doubt the accuracy of such forecasts. A forecast with demonstrated calibration is not immune to this type of perception but would be able to demonstrate over time or across locations that an 80% chance of an outbreak still results in no outbreak 20% of the time. Further work on refining calibration and identifying any relationship of modeling approach and calibration could improve the reliability and usability of forecasts.

The identification of climate factors predictive for WNV activity needs further refinement. Our analysis of modeling approaches indicated that teams that included climate data scored better than those that did not. However, the data source, climatic variables (e.g., minimum temperature, maximum temperature, total precipitation, variance in precipitation, Palmer Drought Severity score, dewpoint, soil moisture, anomalies in temperature or precipitation), and aggregation of the climate variable (e.g., number of days above or below a threshold; weekly average; average of 1–12 months; lagged values up to three year) varied widely among teams. Due to heterogeneity among teams and the limited number of total forecasts, we could not identify the most predictive subset of climatic factors nor the potential importance of variation in data quality among data sources. Similarly, the addition of any seasonal climatic variable in the autoregressive modeling framework we used to select the baseline climate model reduced the forecast skill relative to the AR(1) model. However, this model, which used a single climate variable nationally on a subjectively prescribed three-month season, could not capture spatial variation in climatic zones. Previous studies have also demonstrated challenges in identifying a single environmental driver for predicting WNV activity [47–51]. The essential role of climate in WNV transmission likely varies substantially across different ecological areas, with geographic heterogeneity in which combination of environmental factors, avian populations (composition and seropositivity), and mosquito species drive local transmission.

The forecasts generated here provide some important insight on the challenges with current capabilities and opportunities for improvement, but also on potential uses. As in other forecasting efforts, an ensemble was more accurate than most individual component forecasts. However, in this case, a model based on historical data had more forecast skill and could be considered as a benchmark for a national-scale early warning system even though the current best indicator of high risk is a past history of larger

outbreaks. The use of heuristic principles, like historic outbreaks, can be useful, but sometimes leads to severe and systematic errors [52]. Early indications of high risk can support preparedness across scales, such as resource planning and allocation at the state or local scale. Forecasts at finer spatio-temporal resolution (e.g., two-week forecast on the neighborhood scale) could be even more useful to directly guide effective vector control within counties within seasons [26]. Additional targets like onset or peak week of transmission could also guide vector control activities. There might also be opportunities to frame and communicate forecasts more effectively. Here, we have focused on binned probabilities of different levels of incidence. However, forecasts could also be framed as the probability of above average incidence or predicted range of case numbers (e.g., a 90% prediction interval) that might be actionable in different ways.

Conclusions

The 2020 WNV Forecasting Challenge highlighted the current state of large-scale, early-warning prediction capacity for WNND cases in the United States. Simple models based on previous WNND cases generally performed better than more complex forecasts. The forecasts evaluated therefore indicate that historical incidence provides a relatively reliable indicator of future risk, but substantial uncertainty remains and future models can build upon findings here to improve forecasting as well as providing insight on the probability that the next season will be different from previous seasons. Among models using additional data, inclusion of climate or human demographic data was associated with higher skill, while inclusion of mosquito or land use data was associated with lower skill. These differences indicate that WNV forecasts can benefit by considering location-specific historical data and incorporating additional covariates with caution. Forecast skill was also associated with intrinsic differences among counties, with lower skill in counties with relatively large populations, cold or hot winters, and high variability in yearly case counts. High case count variability likely indicates counties that are intrinsically more difficult to predict, but there may be opportunities to specifically improve predictions for areas with large populations and low or high winter temperatures. Most forecasts, including the highest skill forecasts, also showed patterns of calibration that could potentially be improved. In addition to improved forecast models, increased data collection, data sharing, and real-time data access (e.g., meteorological observations, avian immunity to WNV, mosquito surveillance (abundance and infection rates), mosquito control activities) may support improved predictions. These findings lay the foundation for improving future WNV forecasts.

Abbreviations

CDC

Centers for Disease Control and Prevention

WNND

West Nile virus neuroinvasive disease

WNV

Declarations

Acknowledgements

We thank all those who were involved with data collection, reporting, and data cleaning of WNND cases in ArboNET. We also thank everyone who participated in developing forecasts, including Oliver Elison Timm, Ania Kawiecki, Pascale Stiles, and Sarah Abusaa. Thank you also to Maria Diuk-Wasser and Maria del Pilar Fernandez for their contribution of TickApp data for the NYSW-CVD model. We also thank Stanley Benjamin (National Oceanic and Atmospheric Administration, NOAA), Evan Kalina (NOAA), Georg Grell (NOAA), and Hunter Jones (NOAA) for their hearty discussion around WNV prediction. We also thank Sarah Abusaa (University of California, Davis) for her insights and discussions around forecasting Challenges.

Funding

KMH was a NOAA-CDC climate and health postdoc supported by the NOAA - Climate Adaptation and Mitigation Program and administered by UCAR's Cooperative Programs for the Advancement of Earth System Science (CPAESS) under awards #NA16OAR4310253, #NA18OAR4310253B, and #NA20OAR4310253C. CMB acknowledges funding support from the Pacific Southwest Center of Excellence in Vector-Borne Diseases funded by the U.S. Centers for Disease Control and Prevention (Cooperative Agreement 1U01CK000516). EAM, DK, MPK, and NN were supported by the National Institutes of Health (R35GM133439). EAM was supported by the National Science Foundation (DEB-2011147 with Fogarty International Center), the Stanford Woods Institute for the Environment, King Center on Global Development, and Center for Innovation in Global Health. MLC was supported by the Illich-Sadowsky Fellowship through the Stanford Interdisciplinary Graduate Fellowship. NN was supported by the Stanford Data Science Scholars Program and the Center for Computational, Evolutionary and Human Genomics Predoctoral Fellowship. MJH was supported by the Knight-Hennessy Scholars Program. ACK was supported by cooperative agreement 1U01CK000509-01, funded by the Centers for Disease Control and Prevention. MEG gratefully acknowledges support from a Los Alamos National Laboratory, Laboratory Directed Research and Development, Director's Postdoc Fellowship.

None of the funding bodies had a role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript. The findings and conclusions in this report are those of the author(s) and do not necessarily represent the views of the Centers for Disease Control and Prevention or the Department of Health and Human Services.

Availability of data and materials

The datasets used and/or analyzed during the current study are available in the WNV-forecast-project-2020 repository, <https://github.com/cdcepi/WNV-forecast-project-2020>.

Authors' contributions

KMH: Formal analysis; Writing – original draft; Writing – review & editing; Visualization

MAJ: Conceptualization; Writing – original draft; Writing – review & editing

CMB: Writing – original draft; Writing – review & editing

JES: Data Curation; Writing – review & editing

RJN: Writing – review & editing

CBB: Conceptualization; Writing – review & editing

MLC: Methodology; Software; Validation; Formal Analysis; Investigation

DK: Methodology; Software; Validation; Formal Analysis; Investigation

ELR: Methodology; Software; Writing – review & editing

MJH: Methodology; Software; Validation; Formal Analysis; Investigation; Writing – review & editing

NN: Methodology; Software; Validation; Formal Analysis; Investigation; Writing – review & editing

MPK: Methodology; Software; Validation; Formal Analysis; Investigation; Writing – review & editing

EAM: Writing – review & editing; Supervision; Funding Acquisition

ACK: Methodology; Software; Formal Analysis; Investigation; Writing – review & editing

JMH: Methodology; Writing – review & editing

LWC: Methodology; Writing – review & editing

BDH: Methodology; Software; Validation; Formal Analysis; Investigation; Writing – review & editing

MHP: Methodology; Software; Validation; Formal Analysis; Investigation; Writing – review & editing

MEG: Methodology; Writing – review & editing

MM: Methodology; Writing – review & editing

SM: Project administration; Writing – review & editing

MF: Conceptualization; Data curation; Writing – review & editing

JAU: Methodology; Software; Validation; Formal Analysis; Investigation; Writing – review & editing

ND: Methodology; Software; Validation; Formal Analysis; Investigation; Writing – review & editing

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

References

1. Rosenberg R, Lindsey NP, Fischer M, Gregory CJ, Hinckley AF, Mead PS, et al. Vital signs: Trends in reported vectorborne disease cases – United States and territories, 2004–2016. *Morb Mortal Wkly Rep.* 2018;67(17):496–501.
2. Mostashari F, Bunning ML, Kitsutani PT, Singer DA, Nash D, Cooper MJ, et al. Epidemic West Nile encephalitis, New York, 1999: Results of a household-based seroepidemiological survey. *Lancet.* 2001;358(9278):261–4.
3. Nash D, Mostashari F, Fine A, Miller J, O’Leary D, Murray K, et al. The outbreak of West Nile virus infection in the New York City area in 1999. *N Engl J Med.* 2001;344(24):1807–14.
4. Kramer LD, Ciota AT, Kilpatrick AM. Introduction, spread, and establishment of West Nile virus in the Americas. *J Med Entomol.* 2019;1–8.
5. Centers of Disease Control and Prevention. West Nile virus disease cases and deaths reported to CDC by year and clinical presentation, 1999–2020. Final cumulative maps & data for 1999–2020. 2021. Available from: <https://www.cdc.gov/westnile/statsmaps/cumMapsData.html#three>
6. McLean RG, Ubico SR, Docherty DE, Hansen WR, Sileo L, McNamara TS. West Nile virus transmission and ecology in birds. *Ann N Y Acad Sci.* 2001;951:54–7.
7. Kilpatrick AM, LaDeau SL, Marra PP. Ecology of West Nile virus transmission and its impact on birds in the western hemisphere. *Auk.* 2007;124(4):1121–36.
8. Rochlin I, Faraji A, Healy K, Andreadis TG. West Nile virus mosquito vectors in North America. *J Med Entomol.* 2019;1–16.
9. Kramer LD, Styer LM, Ebel GD. A global perspective on the epidemiology of West Nile virus. *Annu Rev Entomol.* 2008;53:61–81.
10. Ciota AT, Matakchiero AC, Kilpatrick AM, Kramer LD. The effect of temperature on life history traits of *Culex* mosquitoes. *J Med Ent.* 2014;51(1):55–62.

11. Reisen WK. Effect of temperature on *Culex tarsalis* (Diptera: Culicidae) from the Coachella and San Joaquin valleys of California. *J Med Entomol.* 1995;32(5):636–45.
12. Dohm DJ, O’guinn ML, Turell MJ. Effect of environmental temperature on the ability of *Culex pipiens* (Diptera: Culicidae) to transmit West Nile virus. Vol. 39, *J. Med. Entomol.* 2002.
13. Kilpatrick AM, Meola MA, Moudy RM, Kramer LD. Temperature, viral genetics, and the transmission of West Nile virus by *Culex pipiens* mosquitoes. *PLoS Pathog.* 2008;4(6):e1000092.
14. Reisen WK, Fang Y, Martinez VM. Effects of temperature on the transmission of West Nile virus by *Culex tarsalis* (Diptera: Culicidae). *J Med Entomol.* 2006;43(2):309–17.
15. Cornel AJ, Jupp PG, Blackburn NK. Environmental temperature on the vector competence of *Culex univittatus* (Diptera: Culicidae) for West Nile Virus. *J Med Entomol.* 1993;30(2):449–56.
16. Goddard LB, Roth AE, Reisen WK, Scott TW. Extrinsic incubation period of West Nile virus in four California *Culex* (Diptera: Culicidae) species. *Proc Pap Mosq Control Assoc Calif.* 2003;71:70–5.
17. Shocket MS, Verwillow AB, Numazu MG, Slamani H, Cohen JM, El Moustaid F, et al. Transmission of West Nile and five other temperate mosquito-borne viruses peaks at temperatures between 23°C and 26°C. *eLife.* 2020;9:1–67.
18. Hahn MB, Monaghan AJ, Hayden MH, Eisen RJ, Delorey MJ, Lindsey NP, et al. Meteorological conditions associated with increased incidence of West Nile virus disease in the United States, 2004–2012. *Am J Trop Med Hyg.* 2015;92(5):1013–22.
19. Shaman J, Harding K, Campbell SR. Meteorological and hydrological influences on the spatial and temporal prevalence of West Nile Virus in *Culex* mosquitoes, Suffolk County, New York. *J Med Entomol.* 2011 Jul;48(4):867–75.
20. Shaman J, Day JF, Komar N. Hydrologic conditions describe West Nile Virus risk in Colorado. *Int J Env Res Public Heal.* 2010;7:494–508.
21. Landesman WJ, Allan BF, Langerhans RB, Knight TM, Chase JM. Inter-annual associations between precipitation and human incidence of West Nile virus in the United States. *Vector-Borne Zoonotic Dis.* 2007;7(3):337–43.
22. Gardner AM, Hamer GL, Hines AM, Newman CM, Walker ED, Ruiz MO. Weather variability affects abundance of larval *Culex* (Diptera: Culicidae) in storm water catch basins in suburban Chicago. *J Med Entomol.* 2012 Mar;49(2):270–6.
23. Johnson BJ, Sukhdeo MVK. Drought-induced amplification of local and regional West Nile virus infection rates in New Jersey. *J Med Entomol.* 2013 Jan;50(1):195–204.
24. Paull SH, Horton DE, Ashfaq M, Rastogi D, Kramer LD, Diffenbaugh NS, et al. Drought and immunity determine the intensity of West Nile virus epidemics and climate change impacts. *Proc R Soc B.* 2017;284(20162078):1–10.
25. Reiner RC, Perkins TA, Barker CM, Niu T, Chaves LF, Ellis AM, et al. A systematic review of mathematical models of mosquito-borne pathogen transmission: 1970–2010. *J R Soc Interface.* 2013;10(81).

26. Barker CM. Models and surveillance systems to detect and predict West Nile virus outbreaks. *J Med Entomol.* 2019;56:1508–15.
27. Keyel AC, Gorris ME, Rochlin I, Uelmen JA, Chaves LF, Hamer GL, et al. A proposed framework for the development and qualitative evaluation of West Nile virus models and their application to local public health decision-making. Viennet E, editor. *PLoS Negl Trop Dis.* 2021 Sep 9;15(9):e0009653.
28. McDonald E, Mathis S, Martin SW, Staples JE, Fischer M, Lindsey NP. Surveillance for West Nile virus disease - United States, 2009–2018. *MMWR Surveill Summ.* 2021;70(No. SS-1):1–15.
29. R Core Team. R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing; 2021. Available from: <https://www.r-project.org/>
30. PRISM Climate Group, Oregon State University. Monthly mean temperature, minimum temperature, and total precipitation datasets. 2021. Available from: <https://prism.oregonstate.edu>
31. Gneiting T, Raftery AE. Strictly proper scoring rules, prediction, and estimation. *J Am Stat Assoc.* 2007;102(477):359–78.
32. Rosenfeld R, Grefenstette J, Burke D. A proposal for standardized evaluation of epidemiological models. 2012; Available from: https://delphi.midas.cs.cmu.edu/files/StandardizedEvaluation_Revised_12-11-09.pdf
33. Goodrich B, Gabry J, Ali I, Brilleman S. rstanarm: Bayesian applied regression modeling via Stan. 2020. R package version 2.21.1. Available from: <https://mc-stan.org/rstanarm>
34. DeFelice NB, Birger R, DeFelice N, Gagner A, Campbell SR, Romano C, et al. Modeling and surveillance of reporting delays of mosquitoes and humans infected with West Nile virus and associations with accuracy of West Nile virus forecasts. *JAMA Netw open.* 2019;2(4):e193175.
35. Danforth ME, Snyder RE, Lonstrup ETN, Barker CM, Kramer VL. Evaluation of the effectiveness of the California mosquito-borne virus surveillance and response plan, 2009–2018. Rasgon JL, editor. *PLoS Negl Trop Dis.* 2022 May 9;16(5):e0010375.
36. Winters AM, Bolling BG, Beaty BJ, Blair CD, Eisen RJ, Meyer AM, et al. Combining mosquito vector and human disease data for improved assessment of spatial West Nile virus disease risk. *Am J Trop Med Hyg.* 2008;78(4):654–65.
37. Bolling BG, Barker CM, Moore CG, Pape WJ, Eisen L. Seasonal patterns for entomological measures of risk for exposure to *Culex* vectors and West Nile virus in relation to human disease cases in Northeastern Colorado. *J Med Entomol.* 2009;46(6):1519–31.
38. Kilpatrick AM, Pape WJ. Predicting human West Nile virus infections with mosquito surveillance data. *Am J Epidemiol.* 2013;178(5):829–35.
39. Darsie RF, Ward RA. Review of new Nearctic mosquito distributional records north of Mexico, with notes on additions and taxonomic changes of the fauna, 1982-89. *J Am Mosq Control Assoc.* 1989 Dec;5(4):552–7.
40. Darsie RFJ, Ward RA. Identification and geographic distribution of mosquitoes of North America, north of Mexico. *Supplements to mosquito systematics.* Fresno: American Mosquito Control Association; 1981. 1–313 p.

41. Gorris ME, Bartlow AW, Temple SD, Romero-Alvarez D, Shutt DP, Fair JM, et al. Updated distribution maps of predominant *Culex* mosquitoes across the Americas. *Parasites and Vectors*. 2021 Dec 1;14(1).
42. Lindsey NP, Brown JA, Kightlinger L, Rosenberg L, Fischer M. State Health Department Perceived Utility of and Satisfaction with ArboNET, the U.S. National Arboviral Surveillance System. *Public Health Rep*. 2012;127:383–90.
43. Reich NG, Brooks LC, Fox SJ, Kandula S, McGowan CJ, Moore E, et al. A collaborative multiyear, multimodel assessment of seasonal influenza forecasting in the United States. *Proc Natl Acad Sci*. 2019 Feb 19;116(8):3146–54.
44. Reich NG, McGowan CJ, Yamana TK, Tushar A, Ray EL, Osthus D, et al. Accuracy of real-time multimodel ensemble forecasts for seasonal influenza in the U.S. *PLoS Comput Biol*. 2019;15(11).
45. Johansson MA, Apfeldorf KM, Dobson S, Devita J, Buczak AL, Baugher B, et al. Correction for Johansson et al., An open challenge to advance probabilistic forecasting for dengue epidemics. *Proc Natl Acad Sci*. 2019 Dec 17;116(51):26087–8.
46. Cramer EY, Ray EL, Lopez VK, Bracher J, Brennen A, Castro Rivadeneira AJ, et al. Evaluation of individual and ensemble probabilistic forecasts of COVID-19 mortality in the United States. *Proc Natl Acad Sci*. 2022 Apr 12;119(15).
47. Lockaby G, Noori N, Morse W, Zipperer W, Kalin L, Governo R, et al. Climatic, ecological, and socioeconomic factors associated with West Nile virus incidence in Atlanta, Georgia, U.S.A. *J Vector Ecol*. 2016;41(2):232–43.
48. Wimberly MC, Lamsal A, Giacomo P, Chuang TW. Regional variation of climatic influences on West Nile virus outbreaks in the United States. *Am J Trop Med Hyg*. 2014;91(4):677–84.
49. Degroote JP, Sugumaran R, Ecker M. Landscape, demographic and climatic associations with human West Nile virus occurrence regionally in 2012 in the United States of America. *Fac Publ*. 2014;3.
50. Poh KC, Chaves LF, Reyna-nava M, Roberts CM, Fredregill C, Bueno R, et al. The influence of weather and weather variability on mosquito abundance and infection with West Nile virus in Harris County, Texas, USA. *Sci Total Environ*. 2019;675:260–72.
51. Yoo EH, Chen D, Diao C. The effects of weather and environmental factors on West Nile virus mosquito abundance in Greater Toronto area. *Earth Interact*. 2016;20(3):1–22.
52. Tversky A, Kahneman D. Judgment under uncertainty: Heuristics and biases. *Science* (80). 1974;185(4157):1124–31.

Figures

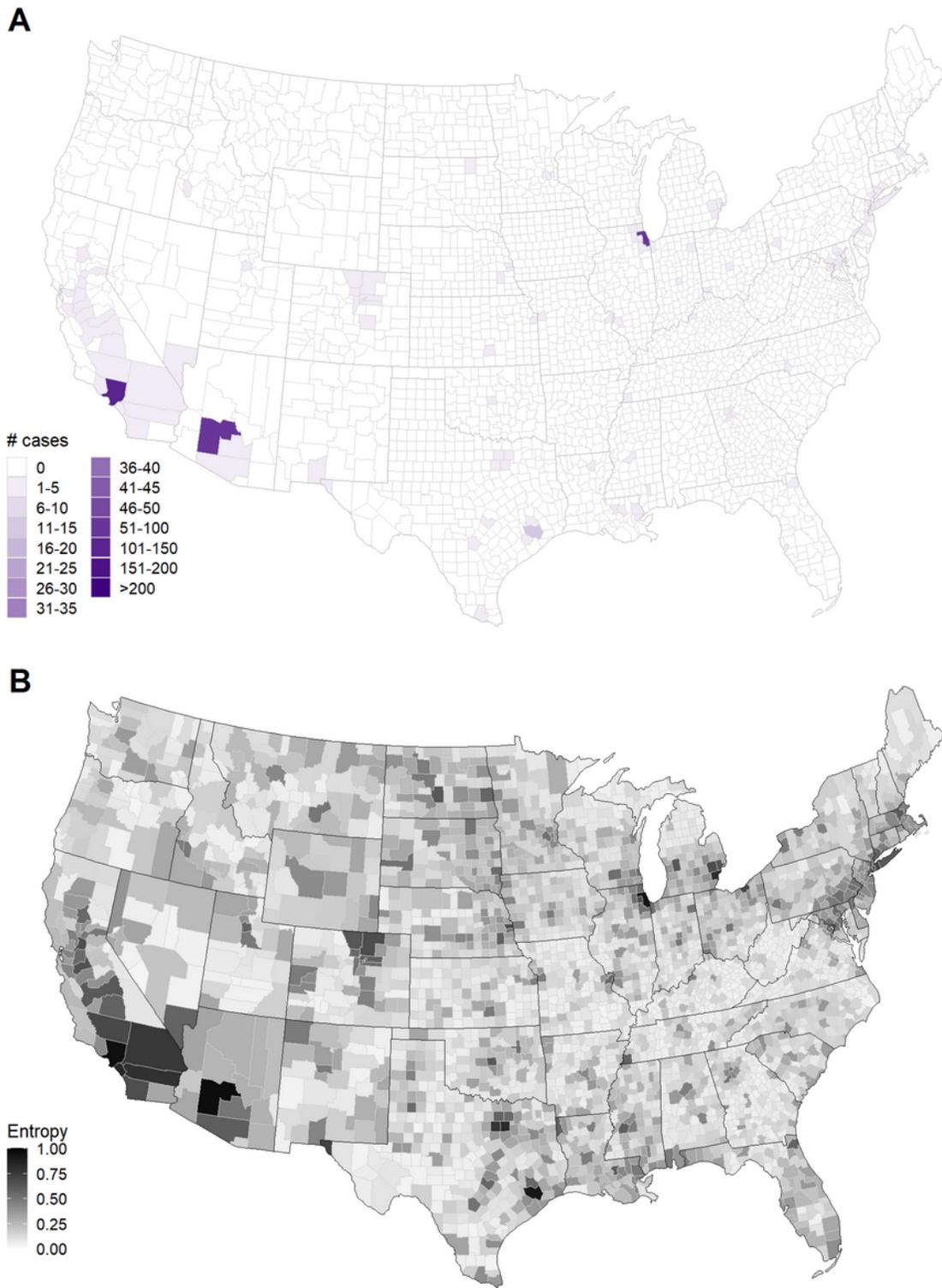


Figure 1

Ensemble forecast for July submissions. A) Most likely number of WNND cases from and B) uncertainty (Shannon entropy) of ensemble model forecast. Mean ensemble model built using the July of last submitted forecasts of all teams and negative binomial model (2000-2019 data). Shannon entropy measures the spread of probability across the binned case counts with a value of zero indicating high

certainty in prediction (all probability in a single bin) and a value of one indicating high uncertainty in prediction (probability equally spread across all bins).

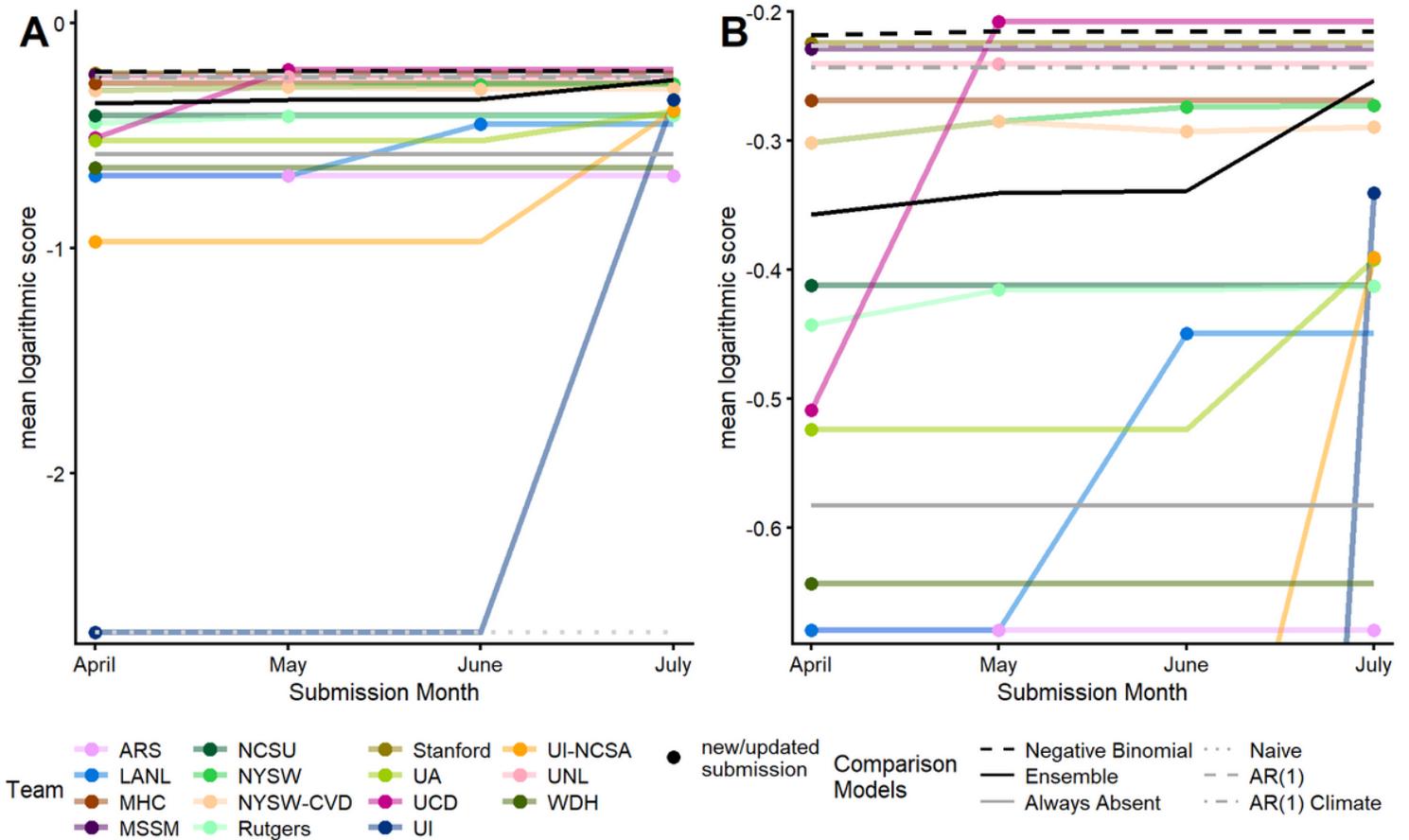


Figure 2

Mean logarithmic score of submissions from teams and comparison models. A) Full range of mean scores and B) vertically truncated range to visualize differences in score among top models for each submission timepoint. If a team did not submit a new forecast at a submission timepoint, we used the previously submitted forecast to calculate the score. See Additional File 1: Table S3 for individual forecast mean logarithmic scores.

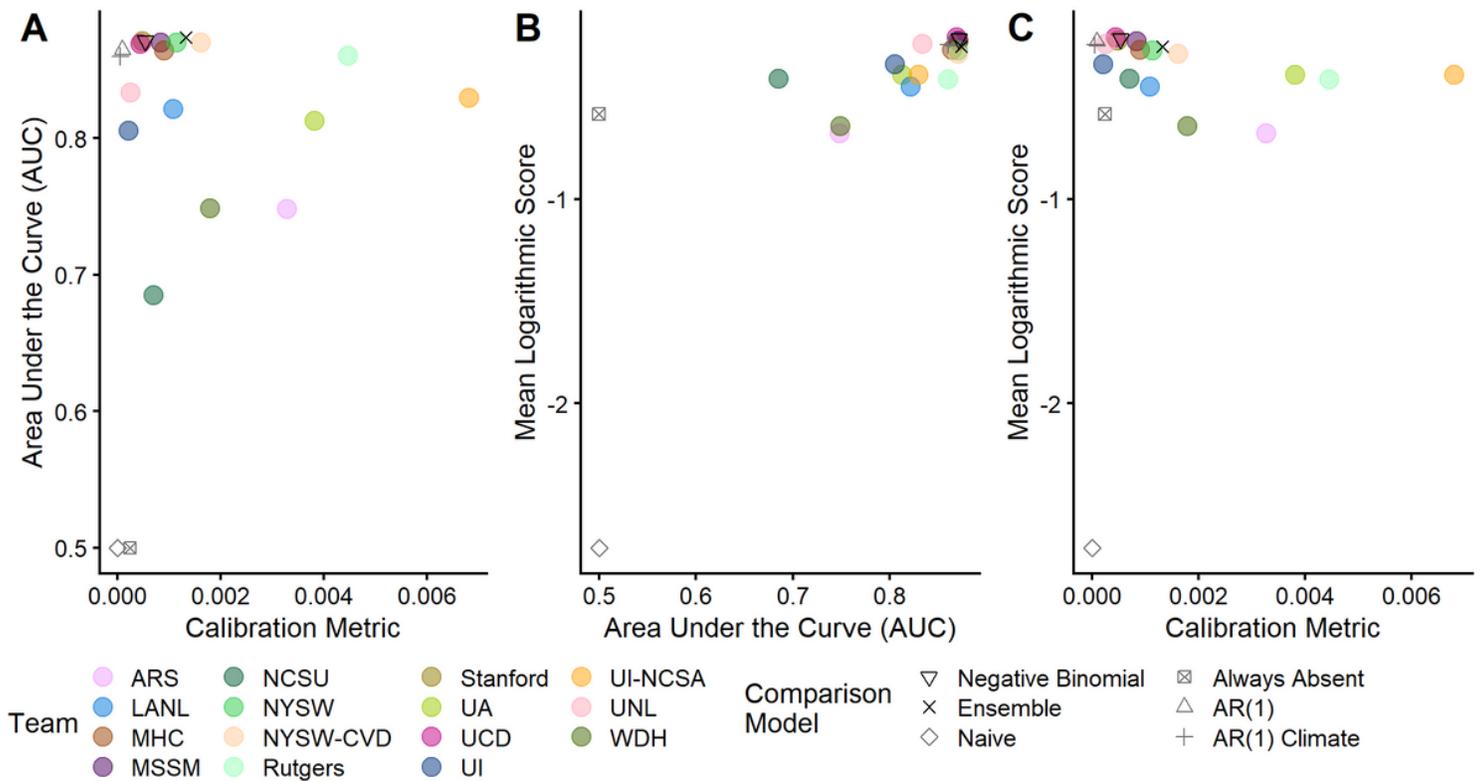


Figure 3

Discrimination, calibration, and mean logarithmic score of final forecasts by teams and comparison models. Area under the curve (AUC) was used to measure a forecast’s ability to discriminate situations with reported WNV cases vs. no cases (AUC of 1.0 would indicate perfect discrimination). Calibration was calculated as the mean weighted squared difference of binned predicted probabilities vs. observed frequency of events (metric of zero is perfectly calibrated). Mean logarithmic score of zero indicates perfect prediction accuracy. Top-performing models are in the top left (A), top right (B), and top left (C). See Additional File 1: Table S3 and Fig S5-S6 for individual forecast score, calibration, and discrimination.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [WNVChallengeAppendix20220720edits.docx](#)