

# Water Resources Research®



## RESEARCH ARTICLE

10.1029/2023WR035901

### Key Points:

- Aggregated lake temperature forecast skill was higher for multi-model ensemble (MME) forecasts than individual model forecasts
- Including baseline empirical models (day-of-year, persistence) with process models improved MME forecast performance
- MME forecasts improved forecast skill by “hedging,” as no individual model performed best at all horizons or depths

### Supporting Information:

Supporting Information may be found in the online version of this article.

### Correspondence to:

F. Olsson,  
[fre Yao@vt.edu](mailto:fre Yao@vt.edu)

### Citation:

Olsson, F., Moore, T. N., Carey, C. C., Breef-Pilz, A., & Thomas, R. Q. (2024). A multi-model ensemble of baseline and process-based models improves the predictive skill of near-term lake forecasts. *Water Resources Research*, 60, e2023WR035901. <https://doi.org/10.1029/2023WR035901>

Received 21 JULY 2023

Accepted 2 FEB 2024

### Author Contributions:

**Conceptualization:** Freya Olsson, Cayelan C. Carey, R. Quinn Thomas  
**Data curation:** Adrienne Breef-Pilz  
**Formal analysis:** Freya Olsson  
**Funding acquisition:** Cayelan C. Carey, R. Quinn Thomas  
**Investigation:** Freya Olsson, R. Quinn Thomas  
**Methodology:** Freya Olsson, Tadhg N. Moore, R. Quinn Thomas  
**Resources:** Adrienne Breef-Pilz  
**Software:** Tadhg N. Moore, Cayelan C. Carey, R. Quinn Thomas  
**Supervision:** Cayelan C. Carey, R. Quinn Thomas  
**Visualization:** Freya Olsson  
**Writing – original draft:** Freya Olsson

© 2024. The Authors.

This is an open access article under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

## A Multi-Model Ensemble of Baseline and Process-Based Models Improves the Predictive Skill of Near-Term Lake Forecasts

Freya Olsson<sup>1</sup> , Tadhg N. Moore<sup>1</sup> , Cayelan C. Carey<sup>1</sup> , Adrienne Breef-Pilz<sup>1</sup> , and R. Quinn Thomas<sup>1,2</sup> 

<sup>1</sup>Department of Biological Sciences, Virginia Tech, Blacksburg, VA, USA, <sup>2</sup>Department of Forest Resources and Environmental Conservation, Virginia Tech, Blacksburg, VA, USA

**Abstract** Water temperature forecasting in lakes and reservoirs is a valuable tool to manage crucial freshwater resources in a changing and more variable climate, but previous efforts have yet to identify an optimal modeling approach. Here, we demonstrate the first multi-model ensemble (MME) reservoir water temperature forecast, a forecasting method that combines individual model strengths in a single forecasting framework. We developed two MMEs: a three-model process-based MME and a five-model MME that includes process-based and empirical models to forecast water temperature profiles at a temperate drinking water reservoir. We found that the five-model MME improved forecast performance by 8%–30% relative to individual models and the process-based MME, as quantified using an aggregated probabilistic skill score. This increase in performance was due to large improvements in forecast bias in the five-model MME, despite increases in forecast uncertainty. High correlation among the process-based models resulted in little improvement in forecast performance in the process-based MME relative to the individual process-based models. The utility of MMEs is highlighted by two results: (a) no individual model performed best at every depth and horizon (days in the future), and (b) MMEs avoided poor performances by rarely producing the worst forecast for any single forecasted period (<6% of the worst ranked forecasts over time). This work presents an example of how existing models can be combined to improve water temperature forecasting in lakes and reservoirs and discusses the value of utilizing MMEs, rather than individual models, in operational forecasts.

## 1. Introduction

In the face of increased ecosystem variability, researchers are developing new methods for forecasting freshwater quality and quantity at near-term (sub-daily to decadal) horizons (Lofton et al., 2023). Here, we define a forecast as a prediction of a future state of a variable with quantified uncertainty (Lewis et al., 2022). Forecasts of freshwater variables have considerable potential for improving management and guiding ecosystem service provision as environmental conditions increasingly exceed the historical envelope due to climate and land use change (Bradford et al., 2020; Dietze et al., 2018; IPCC, 2023). Despite the urgent need for freshwater forecasts, however, the optimal modeling approach for developing forecasts remains unresolved across different spatial and temporal scales. One promising forecasting approach that has emerged from other disciplines is multi-model ensembles (MMEs), in which more than one model is used to simultaneously forecast the same variable into the future (Chandler, 2013; Clark et al., 2022; Humphries et al., 2018; Kirtman et al., 2014; Long et al., 2021; Velázquez et al., 2011). To date, MMEs have not been applied to near-term freshwater forecasting (reviewed by Lofton et al., 2023), motivating the need to understand how an MME forecast performs relative to individual models, as well as how the structure of the different models in the MME influences forecast performance.

Water temperature forecasting in lakes and reservoirs is an ideal application for testing the performance of MMEs. First, water temperature forecasts can be useful for the management of inland waters (Lofton et al., 2023). For example, water temperature forecasts are used to optimise downstream water release from reservoirs (Huang et al., 2011; Jackson-Blake et al., 2022; Weber et al., 2017; Zwart et al., 2023), guide water quality management related to lake mixing events (Carey et al., 2022b; Thomas et al., 2020), as well as underpin the development of other water quality and ecological forecasts (Huang et al., 2011; Page et al., 2018; Weber et al., 2017), given the importance of water temperature for determining metabolism, water chemistry, and biological growth (Carey et al., 2012; Kraemer et al., 2017; Yvon-Durocher et al., 2015). Second, a wide range of models have been

**Writing – review & editing:**

Freya Olsson, Tadhg N. Moore, Cayelan C. Carey, Adrienne Breef-Pilz, R. Quinn Thomas

developed to predict lake and reservoir water temperatures, thereby providing an excellent opportunity for examining the sensitivity of an MME's performance to the identity and structure of multiple component models.

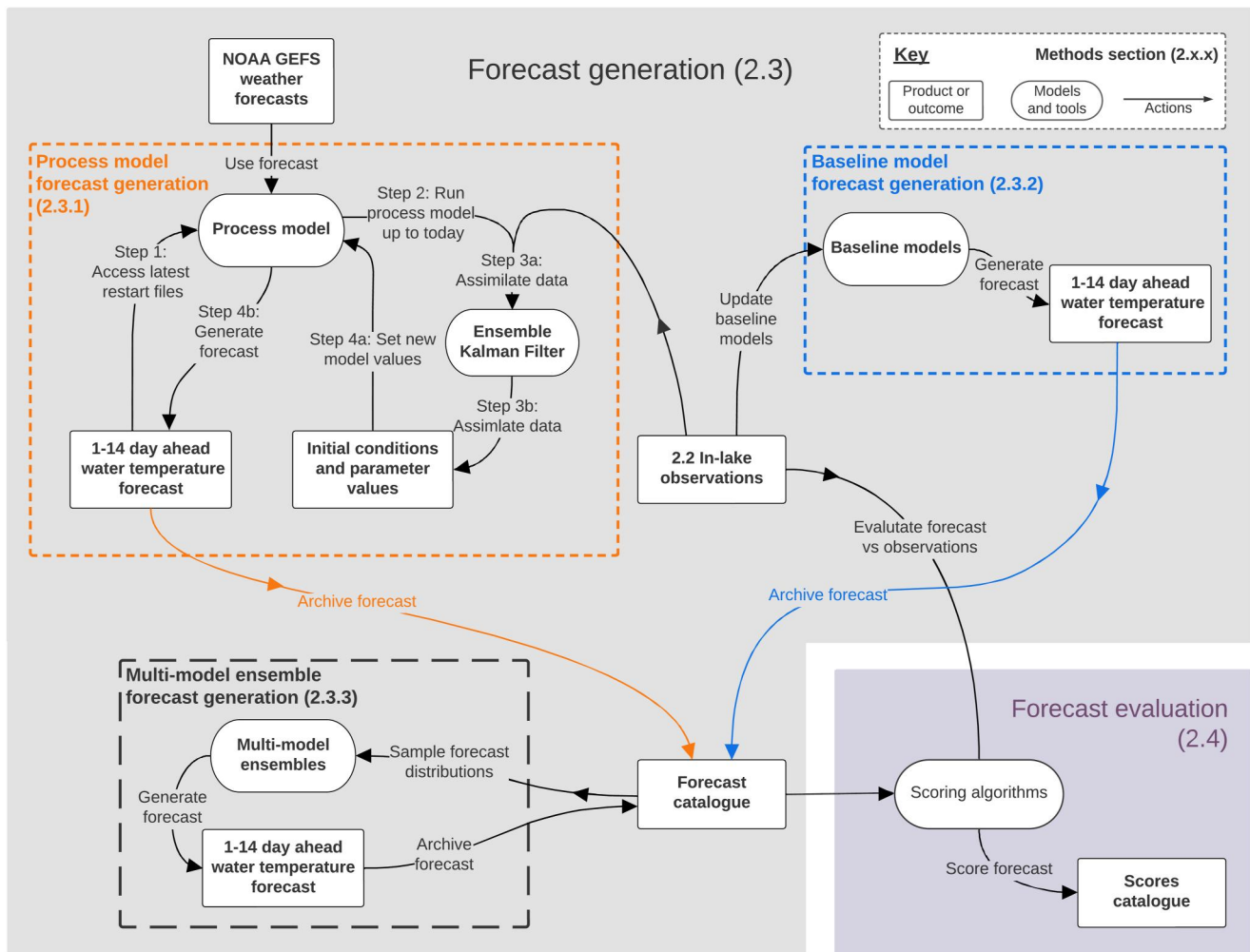
To date, process-based models (Baracchini et al., 2020; Clayer et al., 2023; Mercado-Bettín et al., 2021; Thomas et al., 2020), machine learning and data-driven models (Di Nunno et al., 2023; Read et al., 2019; Zwart et al., 2023), as well as “hybrid” approaches (e.g., Saber et al., 2020) have been used to forecast near-term dynamics (days to seasons ahead) in lake and reservoir water temperatures, with varying levels of performance (reviewed by Lofton et al., 2023). Of these modeling approaches, process-based models (hereafter, PMs) have shown substantial promise, especially in near-term forecast horizons (Baracchini et al., 2020; Carey et al., 2022b; Mercado-Bettín et al., 2021; Thomas et al., 2020), with a performance of 0.4–1.4°C RMSE (daily root mean square error) for reservoir water temperature forecasted 1–16 days-ahead (Thomas et al., 2020). However, the skill of these models is often limited by the skill of other forecasts (e.g., weather and inflow discharge) needed as model driver data (Mercado-Bettín et al., 2021; Thomas et al., 2020). Moreover, PMs also often demonstrate substantial differences in skill among forecasted sites (Thomas et al., 2023) and depths (Thomas et al., 2020), as well as at different times of year (e.g., in thermally stratified vs. mixed conditions; Thomas et al., 2020; Wander et al., 2024).

Despite their simplicity, simple empirical models such as persistence and climatology (historical day-of-year mean and variance) models can also provide useful forecasts (Ward et al., 2014). Often used as null models to test the skill of emerging forecasting approaches (Lofton et al., 2023; Pappenberger et al., 2015), these simple baseline models include information on current conditions and seasonal trends that influence lake temperature dynamics. For example, a persistence model can be useful for forecasting dynamics in systems with high inertia that exhibit small changes across the forecast horizon (i.e., time into the future; Ward et al., 2014), which is common in lakes and reservoirs that exhibit seasonal thermal stratification. Additionally, climatology forecasts exhibit high performance at longer horizons (e.g., months to years), for which repeatable seasonal cycles dominate the dynamics (Pappenberger et al., 2015).

Multi-model ensembles (MMEs) that integrate both PMs and simple baseline models may be particularly effective for forecasting lake and reservoir water temperatures. This type of MME may be able to overcome the limitations of individual process and baseline models that are unable to consistently forecast all environmental conditions with high accuracy across space (i.e., multiple depths in a lake), time (i.e., different seasons within a year), and forecast horizons. Implementation in other disciplines has overwhelmingly found that MMEs often produce more skillful forecasts, on average, than individual model forecasts (Atiya, 2020; Clark et al., 2022; Humphries et al., 2018; Velázquez et al., 2011). Using MMEs also leads to greater diversity in forecast predictions, potentially increasing decision-making success (Boettiger, 2022). Although predictions from individual models can outperform the aggregated prediction from the MME locally, at a specific depth, time, or horizon (Abrahart & See, 2002; Atiya, 2020), it is often not known a priori which forecast model will be best at any given future timestep, especially for forecasts of sites with substantial spatial and temporal heterogeneity. MMEs are ideally suited for these situations, because they integrate information from different model structures into a single forecast, enabling the forecaster to “hedge” (i.e., minimise risk of incorrect forecasts by assigning non-zero probability to a wide range of possible outcomes) and provide a more comprehensive and accurate representation of the potential forecasted outcomes than individual models (Abrahart & See, 2002; Atiya, 2020). MMEs have been successfully applied to a diverse range of ecological and environmental forecasting applications, including ticks (Clark et al., 2022), sea level (Long et al., 2021), penguins (Humphries et al., 2018), and river flow (Abrahart & See, 2002; Velázquez et al., 2011), suggesting that their application for forecasting freshwater ecosystems has promise.

To the best of our knowledge, no one has applied an MME approach to forecasting lake and reservoir temperatures with specified uncertainty. While MMEs for water temperatures have been applied to long-term projections (Almeida et al., 2022; Feldbauer et al., 2022; La Fuente et al., 2022; Wynne et al., 2023), or as model inter-comparisons (Golub et al., 2022), the utility of MMEs for real-time water temperature forecasting remains unknown. This gap may exist because ensemble near-term forecasts have, to date, focused on using ensembles of multiple driver data sets (e.g., weather forecasts; Mercado-Bettín et al., 2021) and parameter sets (e.g., Thomas et al., 2020) to partition and quantify uncertainty (Clayer et al., 2023; Thomas et al., 2020), rather than using multiple models to generate more skillful operational forecasts.

Here, we developed a near-term forecasting system that integrates an MME of lake PMs, baseline empirical models, and data assimilation algorithms in an automated forecasting approach. We used this MME to produce



**Figure 1.** Multi-model ensemble (MME) and individual forecast generation (gray shading) and forecast evaluation workflow (purple shading), with each corresponding text section number in parentheses (e.g., 2.3). Boxes represent tools, objects, and/or products and lines represent actions. The parallel workflows of individual model forecast generation are shown in the orange (process models (PMs)) and blue (baseline models) boxes. Within the PM forecast workflow, the steps correspond to the text in Section 2.3.1. Each individual model forecast is archived into the Forecast Catalog from which the distributions are sampled and combined in the MME (black dashed box). The MME forecasts are also archived in the Forecast Catalog. From this catalog, forecasts are evaluated against in-lake observations via several scoring algorithms to generate a Scores Catalog, which is subsequently analysed (Section 2.5).

weekly, 1–14 days-ahead forecasts of water temperature profiles for two years in a small, temperate, drinking water reservoir. We aimed to understand how MME approaches may improve near-term forecast performance and how forecast performance varies over different spatial scales and forecast horizons. We used the MME forecasts to answer the research questions: (a) How does the forecast performance of the PM MME compare to the individual PMs? (b) How does the addition of the baseline models into the MME affect forecast performance? and (c) How does the forecast performance of the individual models and MMEs vary across horizons and depths? Our goal was to determine if MMEs can improve freshwater water quality forecasting to guide the development of operational forecasting workflows.

## 2. Methods

### 2.1. Overview of Forecasting System

Here, we summarise the automated MME forecasting framework (Figure 1) that leverages the state-of-the-art Forecasting Lake And Reservoir Ecosystems (FLARE) water forecasting system (Thomas et al., 2020). FLARE uses in situ water temperature sensor data, which are wirelessly transmitted directly from the waterbody

to the cloud, in a data assimilation algorithm to update model initial conditions and to calibrate model parameters (Figure 1; Daneshmand et al., 2021). FLARE's ensemble-based forecasting algorithm generates forecasts using process hydrodynamic models that quantify the uncertainty from driver data (weather forecasts from National Oceanic and Atmospheric Administration's (NOAA's) Global Ensemble Forecasting System (GEFS; Hamill et al., 2022), initial conditions, model process, and model parameters and then samples from these sources of uncertainty to generate probability distributions for water temperature at multiple lake or reservoir depths (see Thomas et al., 2020).

Instead of using a single PM, as has been done in previous implementations of FLARE (Carey et al., 2022b; Thomas et al., 2020, 2023), we used three different PMs, implemented via integration with *LakeEnsemblR* R software (LER; Moore et al., 2021), to answer question 1. These PMs were run inside the FLARE framework to generate a multi-model ensemble (MME) from the output (Figure 1), hereafter referred to as the PM MME forecast. To answer questions 2 and 3, two baseline models were also included to produce the full MME forecast (full MME, hereafter), which consisted of five individual models ( $n = 3$  PMs and  $n = 2$  baseline models). Finally, these forecasts are evaluated using the in-situ water temperature observations (Figure 1) via a suite of metrics, described below.

## 2.2. Site Description and Data Collection

We generated water temperature forecasts for Falling Creek Reservoir (FCR), a eutrophic reservoir located in Vinton, Virginia, USA (37.30°N, 79.84°W). FCR is managed by the Western Virginia Water Authority as a drinking water source. The reservoir has a mean depth of 4 m and a maximum depth of 9.3 m, with a surface area of 0.12 km<sup>2</sup> (Carey et al., 2022a). A dimictic system, FCR generally stratifies from May to October and has intermittent ice-cover from December to March (Carey & Breef-Pilz, 2023). The reservoir has one primary inflow and water level is maintained to be generally constant over time.

FCR is monitored by a series of high-frequency sensors deployed at fixed depths in the water column at its deepest site near the dam. Water temperature data were collected using T-Node FR thermistors (NexSens, Fairborn, OH, USA) from March 2019 to March 2023 (Carey et al., 2023; Olsson et al., 2023a), with minor data gaps due to sensor maintenance (see metadata in Carey et al., 2023), across 10 depths in the water column (0.1, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0, 9.0 m). Additional temperature data were collected at 1.6 m with a YSI EXO2 sonde (Xylem Inc., Yellow Springs, OH, USA), and at 5.0 and 9.0 m using RDO PRO-X Dissolved Oxygen Sensors (In-Situ Inc., Fort Collins, CO, USA). For sensor accuracy information, see Text S1 in Supporting Information S1 and Carey et al. (2023). All measurements were collected at a 10-min frequency and averaged to an hourly timestep. Observations were binned into 0.25 m depth intervals to match the vertical resolution of the FLARE forecast output. Measurements within the same bin were averaged together. These data were used in FLARE data assimilation and PM parameter tuning, as well as inputs to the two baseline models (see Section 2.3), and in forecast evaluation (Figure 1, Section 2.4).

Dates of the mixed and stratified periods were calculated based on the density difference between 1 m from surface (1 m) and ~1 m above the bottom (8 m) of the lake, with a minimum density difference of 0.1 kg m<sup>-2</sup> indicating that the lake was stratified (Wilson et al., 2020). The stratified period was defined as the summer period when continuous stratified conditions occurred and the mixed period as any time outside of this. In addition, the observed thermocline depth was estimated using the *LakeAnalyzer* R package (Read et al., 2011).

## 2.3. Forecast Generation

### 2.3.1. Process Model Forecasts

In this application, FLARE generated forecasts for each of the three PMs every seven days for 1 to 14 days-ahead horizons over 2 years (March 2021–March 2023), resulting in a total of  $n = 104$  forecasts per model. For each PM, the weekly forecasts were generated using the following steps (Figure 1): Step 1) access the individual FLARE forecasts (Figure 1, step 1) for 1-week ago in a prior FLARE run (or, in the case of the first forecast, following a spin-up described below); Step 2) use this prediction to initialise each PM FLARE run that starts 1-week ago and runs to current day (Figure 1, step 2); Step 3) use a data assimilation algorithm (the ensemble Kalman filter; Evensen, 2003) to assimilate new observations collected over the past week (Figure 1, step 3) to update that

model's states and parameters (Figure 1, step 3); and Step 4) use the updated states and parameters as initial conditions for a 1–14-day ahead forecast that starts today (Figure 1, step 4).

Each forecast included 256 simulations (ensemble members) that quantified the uncertainty from driver data (weather forecasts), initial conditions, model process, and model parameters. First, for driver uncertainty, each of the 256 FLARE ensemble members was assigned one of the 31 NOAA GEFS ensemble members to drive the PM. Second, initial conditions uncertainty was based on the spread of model states on Day 0 of the forecast that was set by spread in the 256 ensemble members following data assimilation on Day 0. Third, model process uncertainty was generated by adding random noise to the FLARE predictions, drawing from a normal distribution with a standard deviation (SD) of 0.75°C (after Thomas et al., 2020). Finally, parameter uncertainty was generated using a unique parameter value assigned to each of the 256 ensemble members that was determined through data assimilation. Process model parameters that were not included in the data assimilation algorithm had fixed values that represent typical values used in that model. Overall, the 256-member FLARE ensemble for each PM included each of these sources of uncertainty. Additional information about FLARE configuration can be found in Thomas et al. (2020) and Thomas et al. (2023).

Within the FLARE framework, three PMs were implemented using the *LakeEnsemblR* R package (LER; Moore et al., 2021), and underwent data assimilation within FLARE as described above. This R package facilitates the running of up to five one-dimensional hydrodynamic lake models simultaneously using the same driving data and configuration files (see Moore et al., 2021). The three PMs we included in the PM MME were the General Lake Model (GLM; Hipsey et al., 2019), General Ocean Turbulence Model (GOTM; Umlauf et al., 2005), and Simstrat (Goudsmit et al., 2002), hereafter referred to as PM1, PM2, and PM3, respectively. The other two PMs implemented in LER (FLake, MyLake) were not included because our aim was to apply data assimilation and iteratively forecast full water column temperature profiles. Specifically, FLake simulates lake systems using a two-layer representation (Mironov, 2021) that does not simulate a full profile, and MyLake is not able to “restart” daily (Saloranta & Andersen, 2007), as needed for iterative forecasting with data assimilation.

All three PMs require forecasted meteorological driving data to produce water temperature forecasts. To make near-term predictions of water temperature, we used weather forecasts for FCR from the NOAA's GEFS (Hamill et al., 2022). The NOAA GEFS weather forecast consists of a set of 31 simulations, at a spatial resolution of 0.5°, and a forecast horizon of 1–16 days-ahead, which we used to produce 1–14 days-ahead water temperature forecasts from the midnight UTC data product.

We followed the standardised FLARE configuration for forecasting (Thomas et al., 2023). All PMs were run at an hourly time step with the midnight output as the daily forecast. A spin-up of all models was run from 1 October 2020 to 1 March 2021, the date of the first forecast. During this spin-up, each model's parameters were individually tuned by the ensemble Kalman filter within FLARE (see Table S1 and Figure S1 in Supporting Information S1). Each model used default parameters to initialise the forecast run and two sensitive parameters were tuned in the data assimilation process of FLARE (See Table S1 and Figure S1 in Supporting Information S1). The sensitive parameters selected, based on initial investigation and configuration in other lakes, were the sediment temperature and incoming shortwave radiation scaling factor for GLM, and the wind scaling and incoming shortwave radiation factors for Simstrat and GOTM (see Text S2 in Supporting Information S1).

### 2.3.2. Baseline Models

Two simple, empirical baseline models were also used to generate forecasts (Figure 1). The persistence model uses the last observation for each specific depth as a prediction of future conditions and a day-of-year model uses a long-term day-of-year mean as the daily forecast (Hyndman & Athanasopoulos, 2021; Jolliffe & Stephenson, 2012); both are described in detail below.

#### 2.3.2.1. Persistence Model

A persistence model assumes that, on average, the forecasted state (in this case, water temperature) on average does not change over the forecast horizon, with uncertainty driven by a random walk process (Hyndman & Athanasopoulos, 2021):

$$y_{T+1} = y_T + e_{T+1} \quad (1)$$



where  $y_T$  is today's observation or forecast,  $e_{T+1}$  is random noise, and  $y_{T+1}$  is the next day's forecast. The uncertainty ( $e_{T+1}$ ) in the persistence model forecasts were generated using a bootstrapping method, as a normal distribution could not be assumed. The method assumes future uncertainty will be drawn from the same distribution of the residual error in the fit to historical data. The persistence model was iteratively fitted to all observations from the beginning of the data collection up to the day that the forecast was initiated, using the RW (random walk) function in the *fable* R package (version 0.3.2; O'Hara-Wild et al., 2022). From this fitted model the residual, or model error, was calculated between the model and the observation. At each forecast timestep a value of  $e_{T+1}$  was sampled from this distribution of past residuals for each ensemble member. Forecasts were generated using the *generate* function from *fable* with a bootstrap value of  $n = 256$  ensemble members to match the number of simulations as the PMs.

### 2.3.2.2. Day-Of-Year Model

A day-of-year (DOY) model, based on historic observations, was used to generate a second baseline forecast. We assume that the DOY model had a normal distribution with the mean equal to the mean for each day-of-year from the past 2 years. We analysed observations between March 2019–March 2021 from FCR to calculate the day-of-year mean water temperatures at each depth. We chose this period because the thermistor sensors were deployed in summer 2018 and we wanted to ensure that each day-of-year mean water temperature was derived from the same number of historical observations. The SD of the DOY forecasts was equal to the residuals in a linear model between the 2 years of observations (i.e., the observations for each DOY in year 1 were used to predict the observations on the corresponding DOY in year 2). Separate linear models were fit for each depth. We used the residuals from the linear model to estimate the SD, rather than estimating the SD from the observations for each DOY, because we only had 2 years of observations. Lakes with more years of observations can use the latter approach for generating the DOY model. We generated the probabilistic DOY forecasts by sampling from a normal distribution with the obtained mean and SD, generating  $n = 256$  ensemble members.

### 2.3.3. Multi-Model Ensembles (MMEs)

As described above, we generated two MMEs: the PM MME (containing PM1, PM2, PM3;  $n = 3$  models total) and the full MME that also included the two baseline models (persistence, DOY;  $n = 5$  models total). To create the full MME forecasts, the  $n = 256$  ensemble members from each of the three individual PMs and two baseline models were combined into a new MME (Figure 1). We sampled from the pool of individual model simulations to generate MMEs with  $n = 256$  ensemble members, with each model equally represented. The number of simulations from each individual model in the MME forecasts is given as  $256/n$ , where  $n$  is the number of models in the MME. For example, in the full MME there were 5 models represented in the forecast, giving 51 simulations ( $256/5$ ) from each individual model. This sub-sampling of ensemble members to generate 256 simulations for each forecast allowed for comparison among the individual model and MME forecasts, because the number of ensemble members can affect forecast skill; a doubling of ensemble size can result in a non-trivial improvement in forecast skill, especially at longer forecast horizons (Machete & Smith, 2016).

## 2.4. Forecast Evaluation

Forecasts from both the individual models and MMEs were evaluated using four evaluation metrics calculated on each forecast-observation pair. We used multiple evaluation metrics because each metric provides complementary information about the performance of the forecast. First, we calculated the mean absolute bias (absolute difference in mean forecasted water temperature and observed water temperature). Forecasts with lower absolute bias indicate increased forecast accuracy. Second, we calculated the SD of the ensemble using the predictions from each ensemble member to understand uncertainty in the forecasts. We expect uncertainty to increase across the forecast horizon as confidence in future conditions decreases. We also expect to see larger SD in the MME forecasts than individual model forecasts as they reflect a greater diversity of predictions. Both metrics are useful for determining how the forecast accuracy (bias) and precision (using SD as a metric of uncertainty) vary independently and are commonly calculated metrics for forecast performance (Jolliffe & Stephenson, 2012).

Third, we evaluated the models using the ignorance score (IGN), which uses the full ensemble distribution to assess both the accuracy and the precision of the forecasts in its evaluation, and describes the probability placed by the forecast on the observed outcome (Smith et al., 2015). IGN was calculated using the *scoringRules* R package

(Jordan et al., 2019), in which larger values represent a lower probability placed on the observed outcome and lower forecast performance. IGN, originally proposed by Good (1952), is defined as:

$$\text{IGN}(p(x), X) = -\log(p(X)) \quad (2)$$

where  $p(X)$  is the probability density assigned by the ensemble forecast to the observed water temperature ( $X$ ) at that forecasted timestep. The forecast was transformed to a probability density function ( $p$ ) from the ensemble forecast using kernel density estimation, with a Gaussian kernel that was calculated using the default method in the *scoringRules* R package (Jordan et al., 2019). In kernel density estimation we used the default bandwidth produced by the *bw.nrd* function in the *stats* R package (R Core Team, 2021).

IGN penalises forecasts that place very low probabilities on the observed outcome and gives an infinitely large score if a forecast places zero probability on an outcome that is ultimately observed (Smith et al., 2015). We selected the IGN score as a focal evaluation metric because differences in scores between models represent the additional probability placed on the observed outcome in the more skillful forecast (Smith et al., 2015). The exponent of the difference in IGN scores between two models can be used to calculate the probability difference between the models (Smith et al., 2015). For example, an IGN score difference of 0.5 units between two models corresponds to the better model placing  $e^{0.5}$ , or 1.41 times more probability, on the more skillful forecast. Thus, in this example, there is a confidence gain of 41% in the better model compared to the other model (Smith et al., 2015).

Finally, we calculated shadowing time, which quantifies the time that the forecast is able to “shadow” the observations, given an estimate of observational uncertainty (Gilmour & Smith, 1997; Smith et al., 2010). The shadowing time is the maximum number of consecutive days, starting from forecast initiation, that at least one simulation (ensemble member) tracks the mean observation, within a specified observation uncertainty. Here, we define a simulation as shadowing when it falls within the 95% confidence interval of each observation (assuming a normal distribution centered on the observation). Observational uncertainty (SD) was estimated at 0.2°C, based on an analysis of the variation in observations within each day and depth (see Text S3 and Figure S2 in Supporting Information S1). Shadowing time is a useful metric to determine how well the forecast models can replicate the dynamics of a system, rather than the statistics of the forecast (Gilmour & Smith, 1997; Smith et al., 2010).

## 2.5. Analyses

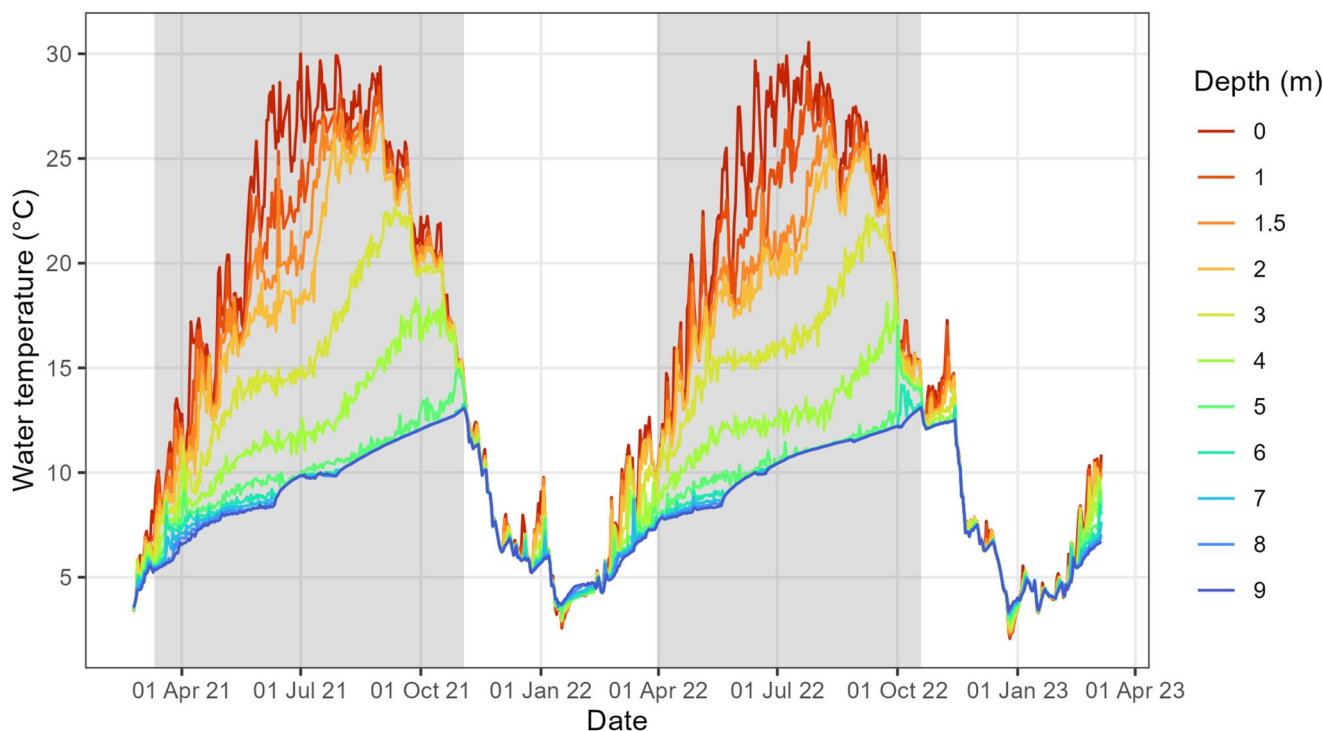
First, to address question 1, we compared the evaluation metrics among the individual PM forecasts and the PM MME forecast. Second, to address question 2, we compared the full MME with the PM MME, and the performance of the five individual PMs and baseline models. Aggregated performance metrics were assessed by calculating the mean of the metric for all forecast dates and depths at each forecast horizon. To understand how and why the MME forecasts might be able to outperform individual models, we also calculated the Pearson correlation coefficient ( $r$ ) on daily forecast bias for each forecast-observation pair. Third, to address question 3, we compared the forecast metrics at different depths and forecast horizons. We also determined each model's rank (out of the 7 forecasts from the five individual and 2 MMEs) for each individual forecast-observation pair across depth and horizon using the IGN score. All analyses were conducted using R statistical software (v.4.2.1; R Core Team, 2021).

All data and code are archived and available in the Zenodo repository (Olsson et al., 2023a, 2023b) or the Environmental Data Initiative repository (Carey & Breef-Pilz, 2023; Carey et al., 2023). Instructions on reproducing the individual model forecasts as well as the MME are available in Olsson et al. (2023b). In addition, the forecasts and scores can be accessed here reproduce the manuscript figures (Olsson et al., 2023a).

## 3. Results

### 3.1. Observed and Forecasted Temperature Dynamics at FCR

FCR exhibited typical seasonal dynamics during the 2-year forecasting period. Continuous summer thermal stratification lasted from 11 March–3 November 2021 and 31 March–19 October 2022. Outside of these periods, there were transient periods of mixing and stratification during spring and autumn as well as brief periods with inverse thermal profiles in winter (Figure 2). We hereafter refer to the period outside of the main summer stratified period as “mixed”. Mean thermocline depth during the summer stratified period was 2.7 m in 2022 and 3.1 m in 2023.



**Figure 2.** Observed high-frequency water temperatures across 11 depths at Falling Creek Reservoir from March 2021 to March 2023. The gray shaded areas show the periods of continuous summer stratification and white shaded areas show the mixed periods.

Our workflow (Figure 1) was able to successfully produce weekly 1–14 days-ahead forecasts for the 2-year forecasting period for all five individual models and the two MMEs (Figure S3 in Supporting Information S1). In general, mean forecast performance was highest at the beginning of the forecast horizon and decreased further into the 14-day-horizon (Figure S3 in Supporting Information S1). Across all depths and horizons, the IGN score of the individual models (other than DOY) increased by 80%–170% from 1 to 14 days-ahead, representing lower performance. Forecast uncertainty, shown by the 95% confidence intervals (Figure S3 in Supporting Information S1), also increased across the 14-day horizon for all models except for DOY (by >100%).

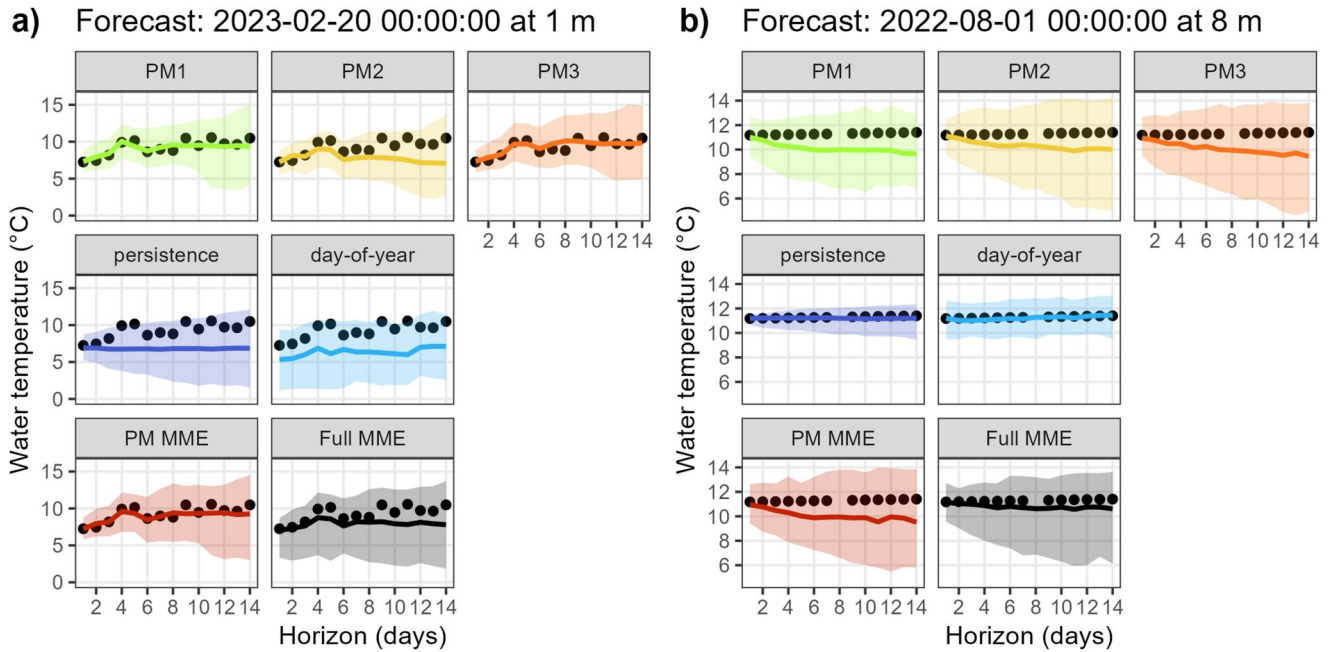
Two examples highlight how the forecasts generated by the individual models exhibited differences in how well they reproduced observations across depths and times (Figure 3). First, forecasts generated during the mixed period at 1 m depth (20 February 2023) show that PM1 and PM3 forecasts closely followed observations throughout the 1 to 14-day ahead horizon, with PM2 diverging from observations after the eighth day of the forecast horizon. In contrast, the DOY and persistence baseline models consistently underestimated water temperature. In a second period with stratified water temperatures (1 August 2022), forecasts generated at 8 m depth show that PM3, and to a lesser extent PM1 and PM2, underestimated the water temperature. The two baseline models were most skillful for this particular forecast.

The variable performances of the individual models are reflected in the performance of the two MME forecasts (Figure 3). For example, in the first example forecasts for 20 February 2023 at 1 m depth, the PM MME performed better than the full MME because of the superior performance of the PMs than the baseline models. Likewise, in the second example forecasts for 1 August 2022 at 8 m, the full MME performed better than the PM MME because of the strength of the baseline models.

### 3.2. Question 1: How Does the Performance of the Process Model MME Compare to the Individual Process Models?

Overall, the PM MME exhibited a higher performance, as determined by the lowest mean absolute bias and mean ignorance score, than the individual PMs (Table 1; Figure 4), when aggregated across all forecast dates, horizons, and depths. The mean absolute bias of the PM MME was similar to PM1, highlighting how the addition of the





**Figure 3.** Two example water temperature forecasts from the five individual models and the two multi-model ensembles (MMEs): one generated on 20 February 2023 at 1 m depth (mixed period; left panels) and one generated on 1 August 2022 at 8 m depth (stratified period; right panels). The top row shows the individual forecasts from the three process models (PMs), the middle row shows the individual forecasts from the two baseline models, and the bottom row shows the MMEs (PM and full MMEs). Shaded areas show the 95% confidence interval around the median forecast (line) and the filled points are the observed water temperatures. The colors for the different forecasts are consistent throughout.

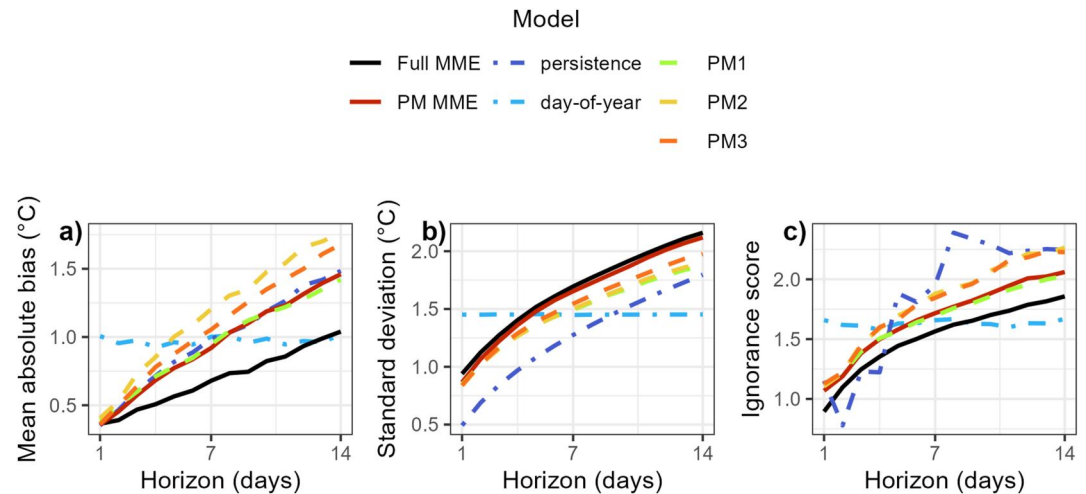
other two PMs with slightly higher absolute bias did not increase bias in the MME (Table 1). The bias increased over the 1–14 days forecast horizon for all PMs and the PM MME, with bias increasing less for the PM MME and PM1 (Figure 4a). In contrast to bias, the PM MME had a larger mean forecast uncertainty (SD) than any of the individual models, especially at longer horizons (Figure 4b). SD increased over the forecast horizon at a faster rate for the PM MME than any individual PM (Figure 4b). At 1 day-ahead, the SD was similar for all PM forecasts (1.1°C), but by 14 days-ahead, the PM MME had 0.2°C higher SD than the best individual PM forecast (Figure 4b).

When using the IGN metric to evaluate performance, which combines accuracy and precision, the PM MME performance was similar but slightly lower than the performance of PM1 (Table 1). This result highlights the penalty given by the IGN score to the higher SD in the PM MME. The best performing PM only placed 14% more probability on the observed outcome than the worst PM forecast on average, and 1% more probability than the PM

**Table 1**  
*The Mean Ignorance Score (IGN), Absolute Bias, Standard Deviation (SD), and Shadowing Time Aggregated for All Forecasts Across All Depths, Times, and Horizons for Each Forecast Model, Individual and Multi-Model Ensemble (MME) Across the Two-Year Forecasting Period*

| Forecast model    | IGN         | Absolute bias (°C) | SD (°C)     | Shadowing time (days) |
|-------------------|-------------|--------------------|-------------|-----------------------|
| Full MME          | <b>1.52</b> | <b>0.69</b>        | 1.66        | 7.3                   |
| Day-of-year       | 1.63        | 0.98               | 1.45        | 3.2                   |
| PM1               | 1.67        | 0.95               | 1.47        | 4.5                   |
| Process model MME | 1.68        | 0.94               | 1.62        | 4.2                   |
| PM2               | 1.80        | 1.16               | 1.49        | 3.7                   |
| PM3               | 1.81        | 1.09               | 1.53        | 4.0                   |
| Persistence       | 1.89        | 0.98               | <b>1.26</b> | <b>7.9</b>            |

*Note.* Models are sorted by most to least skillful, based on IGN, with the “best” forecast based on each metric in bold.



**Figure 4.** (a) Mean absolute bias, (b) standard deviation, and (c) ignorance score across the 14-day forecast horizon for the three individual process models, two baseline models (day-of-year (DOY) and persistence), and the PM multi-model ensemble and full multi model ensemble.

MME (Equation 2). IGN increased over the forecast horizon at a similar rate for both the PM MME and most skillful individual forecast (PM1, Figure 4). At 14 days-ahead, the PM MME placed 13% more probability in the observed outcome and PM1 16% more probability than the two other individual forecast models. This change in probability demonstrates that the MME is not penalised strongly for including the “worse” models overall (Figure 4).

Using the shadowing time metric, the PM MME did not show increased ability to replicate observed water temperature dynamics relative to the individual models. The mean shadowing time for the PM MME (4.2 days) was slightly shorter than the best PM (4.5 days; Table 1). Shadowing time for the other PMs (PM2 and PM3) were shorter than the PM MME but all were between 3.7 and 4.5 days, less than half of the total forecast horizon.

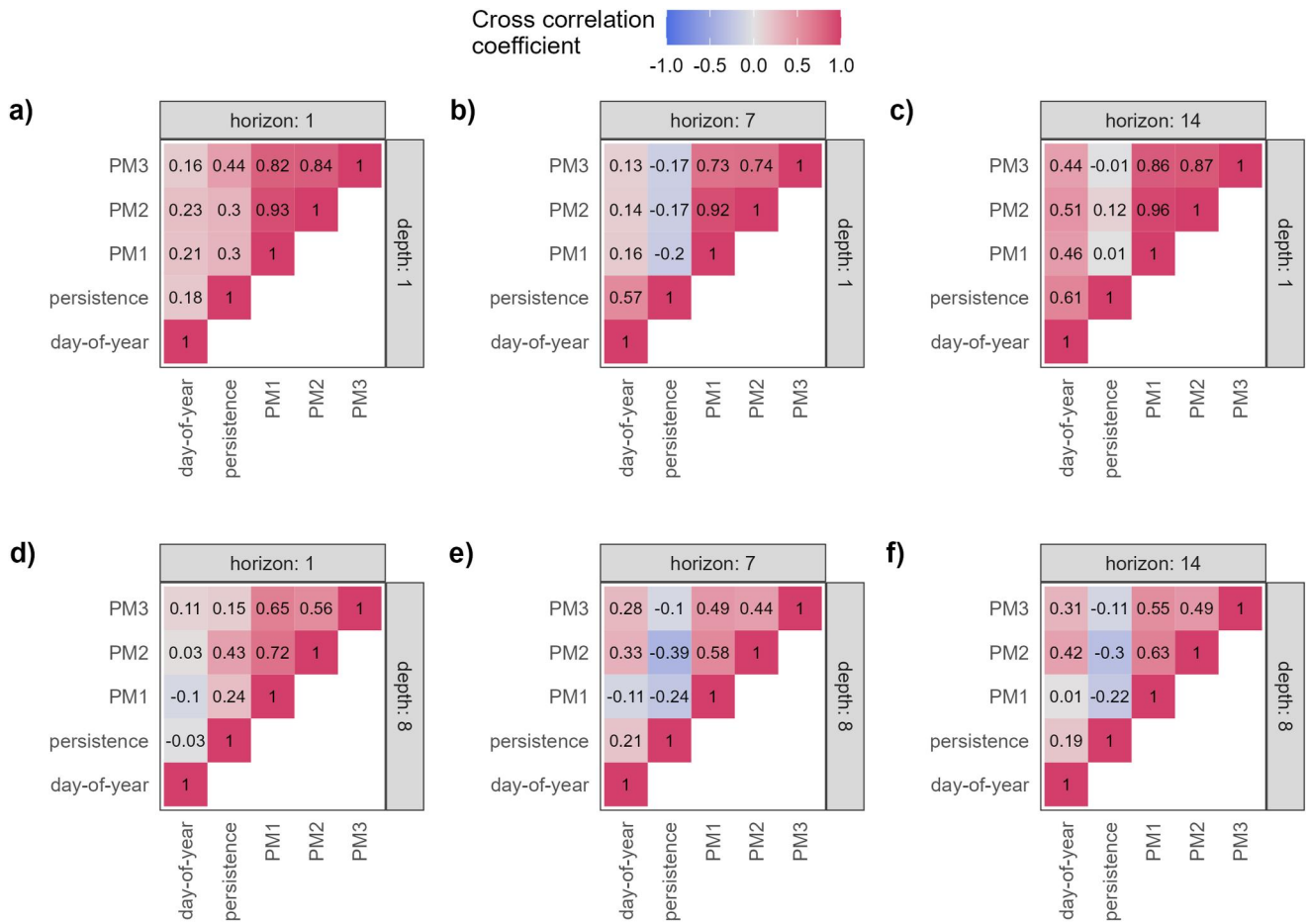
### 3.3. Question 2: How Does the Addition of the Baseline Models Into the Full MME Affect Forecast Performance?

Altogether, the full MME had the lowest IGN, lowest bias, and highest SD of any individual model or MME (Table 1). Aggregated across all depths, times of year, and horizons, the inclusion of the two baseline models into the full MME decreased bias by 26% but only increased the SD by 2% (Table 1), relative to the PM MME. This large reduction in bias led to a lower IGN for the full MME (vs. the PM MME) despite the slight increase in uncertainty. Using the difference in the IGN metric, the full MME placed 16% more probability on the observed outcome than the PM MME. Overall, the improvement in performance of the full MME relative to the PM MME increased throughout the forecast horizon (Figure 4a), a 6% improvement at 2 days-ahead compared to 15% at 14 days-ahead.

The shadowing time of the full MME (7.3 days) was longer than the PM MME (4.2 days; Table 1). This improvement in shadowing time was due to the inclusion of the persistence model in the full MME. The persistence model had the longest shadowing time of any individual model or MME (7.9 days).

The individual PM forecasts exhibited high covariance with other PM forecasts and low covariance with the baseline model forecasts (Figure 5), based on the correlation (Pearson  $r$ ) among individual model forecasts' bias. At 1 m, the PM models exhibited strong positive correlations at 1, 7, and 14 days-ahead horizons ( $r = 0.73$  to  $0.96$ ), with PM1 and PM2 being most correlated at these three horizons. At 8 m, the covariance among the PM models was lower but still showed moderate positive correlation at 1, 7, and 14 days-ahead horizons ( $0.44$ – $0.72$ ) and was lowest at 7 days-ahead horizons ( $r = 0.44$  to  $0.58$ ).

In contrast to the individual PM models, the individual baseline models generally showed low covariance between each other and with the PM models (Figure 5). A few exceptions to this pattern were at the 1 day-ahead horizon, when the persistence model showed a moderate correlation with PM3 ( $r = 0.45$ ) at 1 m and with PM2 at 8 m



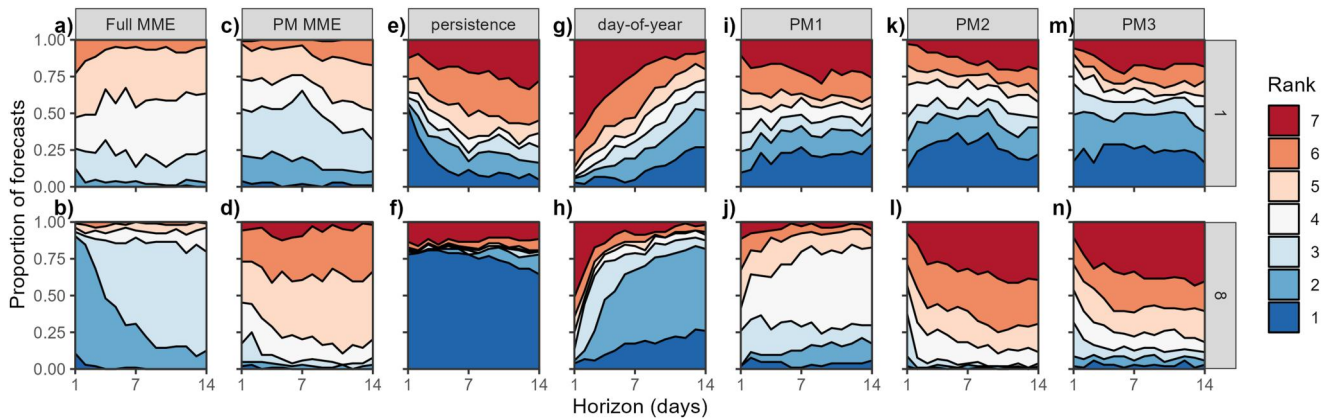
**Figure 5.** Correlation (Pearson  $r$ ) of bias among individual model forecasts. The correlation coefficient between individual models was calculated for the mean forecast bias (mean - observations) at 1, 7, and 14 days-ahead horizons for 1 and 8 m. Red indicates a strong positive correlation and blue indicates a strong negative correlation.

( $r = 0.43$ ). Similarly, at 7 and 14 days-ahead horizons, the persistence and DOY models showed a positive correlation at 1 m ( $r = 0.58$  and  $r = 0.60$ , respectively), although little correlation occurred between the two baseline models at 8 m. Additionally, at the 14 days-ahead horizons, the DOY model was positively correlated with all other models for the 1 m forecasts and showed some correlation with PM2. The correlations among the PMs were always higher than any correlation involving a baseline model (Figure 5).

### 3.4. Question 3: How Does the Forecast Performance of the Individual Models and MMEs Vary Across Horizons and Depths?

The ranking of models demonstrates the hedging that occurs when using MMEs to forecast at different depths and horizons. The individual model forecasts were more likely to be ranked the “worst” of the seven forecasts (Figure 6, Figure S4 in Supporting Information S1) than the two MME forecasts. Out of all  $n = 104$  forecasts generated, the full MME had  $<1\%$  ( $n = 1$  forecast) of rank 7 (worst) forecasts across both 1 and 8 m and 10% of rank 1 (best) forecasts ( $n = 11$ ). At 1 m, the full MME was most often ranked in the middle (65%–95% of forecasts ranked 3–5, respectively; Figure 6a). At 8 m, the full MME was more often ranked as the second-best forecast, especially at shorter horizons (Figure 6b), with more than 50% of forecasts at ranks 1 or 2 up to 4 days-ahead. Despite the decrease in high-ranking forecasts (ranks 1–2) at longer horizons, there was no appreciable increase in the proportion of worst-ranking forecasts (ranks 6–7), remaining between 2% and 6% of forecasts at most horizons.

The individual PM forecasts were dominated by rankings of either the best (1) or worst (7) performance, whereas the PM MME had fewer of these extreme ranks. At 1 m, the individual PM models had almost equal proportions of rank 1 and rank 7 forecasts across the full horizon (Figures 6i, 6k, and 6m), with over 40% of forecasts ranked at



**Figure 6.** Proportion of total forecasts ( $n = 104$ ) with each rank, from 1 (best) to 7 (worst), out of the five individual models and two multi-model ensembles (MMEs) (process model (PM) and full MME). Ranks were calculated for each individual forecast ( $n = 104$ ) and each horizon (1–14 days-ahead) based on the ignorance score at 1 m (top row) and 8 m (bottom row) depths.

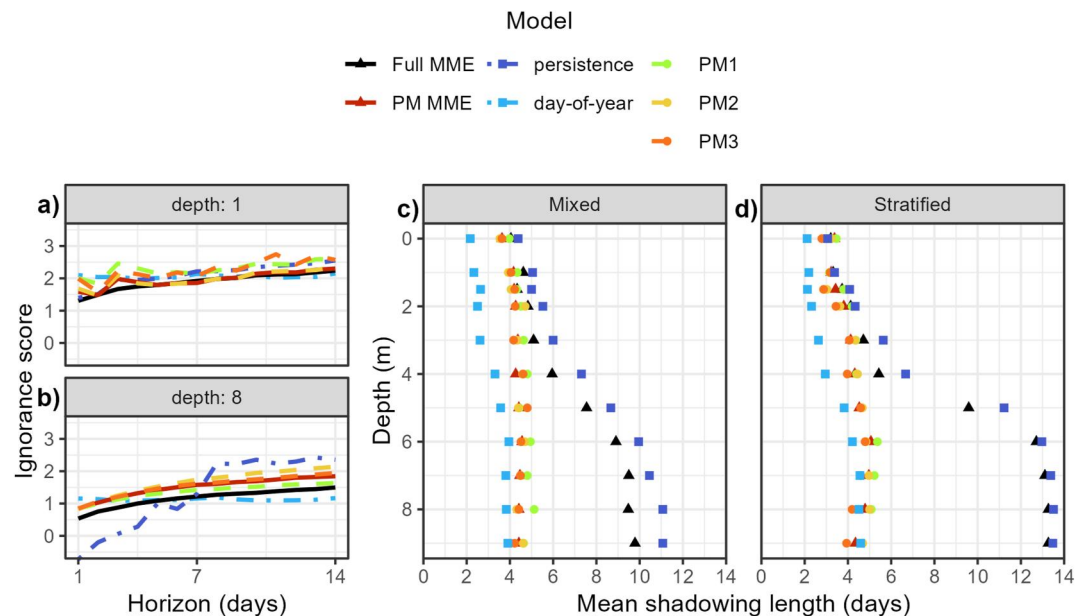
one of these extremes, compared to only 2% of the PM MME forecasts exhibiting one of these extreme ranks. At 8 m, the individual PM models were more often at an intermediate rank than at 1 m (Figures 6j, 6k, 6l, and 6n), although PM2 and PM3 had more than 40% of the worst forecast, whereas PM1 had up to 56% of forecasts with an intermediate rank and fewer very poor forecasts (rank = 7).

The ranks of the baseline models varied substantially at different depths and horizons. At 1 m, the persistence model had more than 50% of forecasts in rank 1 for 1 day-ahead forecasts, which declined steeply to only 10% at horizons >5 days-ahead (Figure 6e). Concurrently, the proportion of forecasts for which persistence was the worst forecast also increased across the forecast horizon, with more than 50% of the forecasts having persistence at ranks 6 or 7 forecast at 13–14 days ahead (Figure 6e). At 8 m, the persistence forecasts dominated the best performing rank across the whole horizon (Figure 6f), only decreasing marginally from around 80%–65% of total forecasts by 14 days-ahead (Figure S4 in Supporting Information S1). The DOY model demonstrated strengths at longer horizons at both 1 and 8 m. The proportion of DOY forecasts at 1 m with a rank 1 increased across the forecast horizon, from <5% at 1 day-ahead to 26% of forecasts at 14 days-ahead (Figure 6g). However, DOY was frequently the least skillful forecast at 1 m, especially at 1 day-ahead (Figure S4 in Supporting Information S1; 65% of forecasts). At horizons between 3 and 10 days-ahead, 40% of the DOY forecasts were either the first or second ranked forecast, which increased to 80% at horizons >10 days-ahead.

Inspection of the disaggregated forecast scores further demonstrates that there was no one consistently best-performing model or MME at all horizons and all depths, as determined by IGN scores and shadowing time (Figures 7a and 7b). At 1 m, the two MME forecasts had the highest skill across the total horizon (Table S2 in Supporting Information S1), although they were outperformed at certain horizons by the individual PM2 model and, beyond 10 days-ahead, the DOY model (Figure 7a). Conversely, at 8 m, the persistence model had the best performance for 2 days-ahead, the full MME exhibited the best performance 1 and 3–5 days-ahead, and then the DOY model had the highest skill up to 14 days-ahead (Figure 7b).

As with the aggregated shadowing time, including the baseline models in the full MME extended the shadowing time compared to the PM MME at almost all depths during both the stratified and mixed periods (Figures 7c and 7d). The persistence model had the longest shadowing time across all forecasts (mean = 7.9 days, Table 1), which was consistent across depths, except for forecasts at the surface (0 m) during the stratified period (Figure 7d). The persistence model showed significantly better shadowing ability than the other individual model forecasts, especially at depths deeper than 4 m, which corresponded to depths below the thermocline, calculated at a depth of 2.7–3.1 m during the forecast period. For example, at 5 m, the shadowing time of the persistence forecast during the stratified period was 2.5 times longer than the next best individual model (PM1). The shadowing time of the PM MME did not improve on the best individual model (PM1), although all PM showed low shadowing ability (<6 days at all depths) relative to the persistence and full MME. At 8 m, both the persistence model and the full MME were able to almost shadow the full horizon (Figure 7d; mean shadowing time = 13.5 and 13.2 days, respectively).





**Figure 7.** Disaggregated forecast performance (ignorance score) at 1 (a) and 8 m (b) for each and mean shadowing time at each observed depth in the water column in the mixed (c) and stratified periods (d) for the three individual process models (PMs), two baseline models (day-of-year and persistence), and the PM multi-model ensemble (MME) and full MME.

#### 4. Discussion

Reservoir water temperature forecasts generated using a MME consisting of process and baseline models performed better overall than using individual models or a PM MME. Our results support previous research that shows that MME methods often outperform individual models (Atiya, 2020; Johansson et al., 2019; Viboud et al., 2018). For example, in a large diverse forecasting competition of multiple finance and demography variables, 70% of the most accurate forecasts were MMEs (Atiya, 2020). Our results showed that no individual model performed best at all depths and horizons, as the best models at 1 m (the individual PMs) were the worst performers at 8 m. In contrast to this finding, the full MME was rarely the worst-performing forecast, highlighting the hedging ability of MMEs to prevent very poor forecast performance (Atiya, 2020). MMEs incorporate the strengths of multiple models given that all models are likely imperfect representations of reality (Atiya, 2020) as well as acknowledging the between-model uncertainty (Humphries et al., 2018). Below, we examine some of the implications for using MME forecasts and highlight ways to further improve MME forecasts for other applications.

##### 4.1. No One Individual Model Is Optimal for All Forecast Horizons or Depths

For individual 1–14 days-ahead forecasts at specific horizons and depths, individual models outperformed the MMEs (Figure 7), accounting for >96% of the best forecasts at 1 m and >91% at 8 m (Figure S4 in Supporting Information S1). Each model captures slightly different dynamics of the mechanistic processes controlling reservoir water temperature and therefore performed optimally under different conditions (Lapeyrolerie & Boettiger, 2023). This was also observed in a multi-model river forecasting study in which individual models alternately performed best in predicting different stages, phases, or mechanisms of rainfall-runoff (Abrahart & See, 2002) and a penguin population forecasting study in which a range of models differentially captured inter-annual and inter-species variability (Humphries et al., 2018). Altogether, our study contributes to the evidence that combining forecasts from different models provides a more comprehensive and accurate representation of the forecasted system than one model alone.

In our analysis, the optimal model varied by depth and horizon, demonstrating the individual strengths of each model. The persistence forecast was significantly better across all horizons at 8 m than other models (ranked best in 66%–82% of all forecasts, Figure 6c), but generally performed poorly at 1 m at horizons beyond 1 to 2



days-ahead (Figure 7). This finding is in agreement with a previous water temperature forecast study at the same reservoir, which found high forecast skill from a persistence model deeper in the lake and higher skill from a PM at the surface (Thomas et al., 2020). Individual PMs have been shown to be successful at forecasting water temperature dynamics at the lake surface at short horizons (Thomas et al., 2020; Wander et al., 2024). As weather forecast skill degrades further into the future, there is a subsequent reduction in water temperature forecasting skill at these shallower depths at longer horizons (Carey et al., 2022b; Thomas et al., 2020). This pattern is likely because meteorological driver data uncertainty has been shown to be the primary source of uncertainty in surface water temperature forecasts, due to the sensitivity of surface water temperatures to atmospheric forcing (Thomas et al., 2020).

One promising approach for better utilizing the strengths of the individual models is to weight the individual models within the MME based on their historical forecast performance. Weighting the individual models may further increase MME skill (reviewed by Wang et al., 2022), as these methods seek to exploit the inherent benefits of each individual model represented in the MME (Abrahart & See, 2002). MME blending methods that weight accurate models more highly and adjust weights dynamically may leverage the strengths of the models whilst minimizing their weaknesses (Chandler, 2013; Spence et al., 2018). For example, Abrahart and See (2002) used a fuzzy logic approach to use the previous forecast performance to weight the models used in the next forecast MME when forecasting river flow. However, selecting the optimal MME blending method was dependent on the dynamics of the flow conditions (Abrahart & See, 2002). Wang et al. (2022) note in their review that simple combination methods, such as the linear pooling with equal weights (as done here) or simple averaging, are some of the most robust approaches for model blending and that improvements from optimised weights can be outweighed by the error added by estimating these parameter values (Dormann et al., 2018). In short, estimating the weighting parameters adds another source of uncertainty to the forecasts whereas simple averaging is robust and easier to implement (Barrow & Kourentzes, 2016). A potential alternative to applying a weighting method would be to identify a suitable pool of models to use in the MME and omit the worst performing ones, thus diminishing the worst predictions within each individual model forecast (Abrahart & See, 2002; Dormann et al., 2018), unless it is very diverse from the other models (Atiya, 2020). However, in real-time forecasting applications it is not possible to know which models will perform best or worst under future conditions. Therefore, model performance can only be assessed based on *past* forecasts, which may or may not continue to perform similarly under future conditions, limiting the potential of such weighting methods.

#### 4.2. The PM MME Did Not Significantly Improve on the Best Individual Process Model

When aggregating the ignorance score across all forecasts, the PM MME performed slightly worse than the best individual PM model. However, the PM MME had many fewer individual forecasts when it was ranked as the least skillful model (Figure 6). This result demonstrates the value of hedging through MMEs. Even when the aggregate forecast skill of the PM MME is not significantly improved compared to its individual models, the process MME still provides value by preventing the generation of poorly performing forecasts that can occur from individual models (Doblas-Reyes et al., 2005; Hagedorn et al., 2005).

Overall, the performance of the individual PMs was highly positively correlated (Figure 5), limiting the amount of unique information provided by individual models to the MME. Others have found that MME forecasts were most skillful when the covariance among models was low (Dormann et al., 2018; Renwick et al., 2018), as well as when models exhibit diverging bias in their mean predictions (Dormann et al., 2018; Petropoulos et al., 2022). This finding supports the need for more diverse model structures to fully optimise the MME forecasts. In this study, high covariance among PMs was likely caused by three key drivers. First, the three PMs were all 1-D hydrodynamic models. Examining whether adding more complex PMs (e.g., 3-D models) or simpler PMs (e.g., Hanson et al., 2023) could help reduce inter-model covariance is another opportunity for further research. Second, the three PMs all used the same forecasted weather from the NOAA Global Ensemble Forecasting System as driver data. Future work could include models that use alternative weather drivers, such as different weather forecast products (e.g., Buizza & Richardson, 2017) or historical weather climatology. Third, all three PMs applied the same data assimilation algorithm (an ensemble Kalman filter). Future work could explore the influence of the diversity of data assimilation algorithms on MME forecasts by including alternative data assimilation approaches, such as a particle filter (Fearnhead & Künsch, 2018).

### 4.3. Including Baseline Models in the MME Improved Forecast Skill

Results from the full MME demonstrate that more model diversity within an MME increases forecast skill (Figure 4c; Table 1). The most model diversity was added to the MME by including the two baseline models that represent end members of empirical models. Specifically, the persistence model represents the most recent data and or DOY represents the long-term historical average for the forecasted system. By including these baseline empirical models, water temperature forecast performance was substantially increased compared to the PM MME (Table 1). The improvement in forecast performance was particularly evident in the full MME's ability to shadow observations, caused by the inclusion of the persistence model. The persistence model had longer shadowing time than the other individual models, especially in the more stable depths below the thermocline, where turbulent mixing is less frequent (Cannon et al., 2021), preventing rapid changes in temperature. The uncertainty in the persistence model can capture the slower changes in temperatures at 8 m, thus shadowing observations. Including ensemble members from the persistence model in the MME therefore increased the model's ability to also capture this dynamic.

Including baseline models in an MME presents a relatively easy approach, with low computational costs, to improve forecast performance if data are readily available for constructing the baseline models. While many forecasting studies use baseline models as null models to evaluate forecasts, here we show their value as a component of the forecast themselves. These baseline models, despite their simplicity, provide additional forecast information that the complex PMs do not, and highlight that model complexity does not necessarily translate to forecast skill (Viboud et al., 2018; Ward et al., 2014). Even simple models, lacking any domain expertise, can provide useful information to an MME (Wang et al., 2022). For example, forecasting of penguin populations showed that simpler domain-agnostic time series models produced better forecasts than complex domain-specific population models (Humphries et al., 2018).

### 4.4. Recommendations and Next Steps

Identifying a set of models with low covariance is likely to increase aggregated forecast skill from an MME relative to its individual models. In advance of producing an MME forecast, a model selection process would help ensure that the MME will improve skill relative to individual models, based on among-model covariance and individual model variance and bias (Dormann et al., 2018; Hagedorn et al., 2005). It is likely that the optimal set of models to include in the MME will be specific to individual sites, given how individual models perform differently among lakes (e.g., Bruce et al., 2018). For example, the same forecast model performed better at some lakes than others in a multi-site comparison (Thomas et al., 2023), with similar differences in model performance found among sites when forecasting phytoplankton (Page et al., 2018; Rouso et al., 2020).

Further ways to improve forecast skill should also focus on constraining uncertainty. The full MME had the highest variance of any of the forecast models, which undermines some of the improvement in bias from the model averaging and leaves a forecast that is likely underconfident (Wang et al., 2022). Methods such as boosting, dimensionality reduction, and trimming can optimise bias-variance trade-offs (Wang et al., 2022). For example, trimming the tails (exterior) of the individual forecast distributions has been shown to increase confidence in the MME by reducing the variance of the individual model forecasts before being combined into an MME (Howerton et al., 2023; Zhao et al., 2022). Previous results showed that MMEs were more successful when their component model forecasts were overconfident (low variance) (Hagedorn et al., 2005; Wang et al., 2022; Weigel et al., 2008).

Finally, our results demonstrate the value of calculating multiple evaluation metrics when assessing the skill of forecasting methods, as each metric highlights potential areas to improve overall skill. For example, the forecast SD evaluation showed that uncertainty was much larger for the MMEs than any individual model (Figures 7c and 7d). Simultaneously, the MMEs had the lowest bias (Figures 7a and 7b). The IGN score was able to combine these two evaluation components into a single metric of statistical performance, highlighting that improvements in overall performance would likely come from reducing forecast uncertainty. Although shadowing time is a metric infrequently used in freshwater forecast evaluation (Lofton et al., 2023), it is potentially valuable, given its focus on the model's ability to replicate actual dynamics, rather than just the statistics of the forecast (Gilmour & Smith, 1997; Petropoulos et al., 2022) or the shape of the distribution (Smith et al., 2015), providing information on likely lead times at which a forecast will have utility (Smith et al., 2010). Improving the capacity of the PMs to have longer shadowing time may help improve their overall representation of lake and reservoir dynamics.

## 5. Conclusions

This work has demonstrated the usefulness of MMEs in improving water temperature forecasts. A five-model MME had the highest forecast skill among all of the forecasts generated by individual models or a three-model MME, which is likely due to hedging: the five-model MME was able to avoid generating very bad forecasts despite being unable to provide the most skillful forecast at many individual horizons or depths. The addition of two baseline models, which had low covariance with the PM models, into the MME provided useful shadowing ability and complementary forecast information. Our results present an example of how existing models can be combined to improve water temperature forecasting in lakes and reservoirs. Future work could focus on including additional forecasting model structures to further increase the diversity of predictions included in the MME and investigate optimal methods to blend predictions and constrain model variance. Altogether, we highlight the value of including simple baseline models (which may in some cases be already calculated as null models for forecast evaluation) into MMEs for forecasting to improve forecasting skill effectively and efficiently with little additional effort.

## Data Availability Statement

All data and code to produce the forecasts, scores, and figures presented in this manuscript are available in the Zenodo repositories (Olsson et al., 2023a, 2023b) or the Environmental Data Initiative repositories (Carey & Breef-Pilz, 2023; Carey et al., 2023).

## Acknowledgments

We thank the Reservoir Group for their help with data collection and helpful feedback, and Lenny Smith for many constructive discussions of forecast evaluation that enhanced this work substantially. This work was financially supported by US National Science Foundation grants DEB-1926388, DBI-1933016, DBI-1933102, and DEB-2327030.

## References

- Abrahart, R. J., & See, L. (2002). Multi-model data fusion for river flow forecasting: An evaluation of six alternative methods based on two contrasting catchments. *Hydrology and Earth System Sciences*, 6(4), 655–670. <https://doi.org/10.5194/hess-6-655-2002>
- Almeida, M. C., Shevchuk, Y., Kirillin, G., Soares, P. M. M., Cardoso, R. M., Matos, J. P., et al. (2022). Modeling reservoir surface temperatures for regional and global climate models: A multi-model study on the inflow and level variation effects. *Geoscientific Model Development*, 15(1), 137–197. <https://doi.org/10.5194/gmd-2021-64>
- Atiya, A. F. (2020). Why does forecast combination work so well? *International Journal of Forecasting*, 36(1), 197–200. <https://doi.org/10.1016/j.ijforecast.2019.03.010>
- Baracchini, T., Wüest, A., & Bouffard, D. (2020). Meteolakes: An operational online three-dimensional forecasting platform for lake hydrodynamics. *Water Research*, 172, 115529. <https://doi.org/10.1016/j.watres.2020.115529>
- Barrow, D. K., & Kourentzes, N. (2016). Distributions of forecasting errors of forecast combinations: Implications for inventory management. *International Journal of Production Economics*, 177, 24–33. <https://doi.org/10.1016/j.ijpe.2016.03.017>
- Boettiger, C. (2022). The forecast trap. *Ecology Letters*, 25(7), 1655–1664. <https://doi.org/10.1111/ele.14024>
- Bradford, J. B., Weltzin, J. F., McCormick, M., Baron, J., Bowen, Z., Bristol, S., et al. (2020). *Ecological forecasting—21st century science for 21st century management*. U.S. Geological Survey Open-File Report 2020-1073. <https://doi.org/10.3133/ofr20201073>
- Bruce, L. C., Frassl, M. A., Arhonditsis, G. B., Gal, G., Hamilton, D. P., Hanson, P. C., et al. (2018). A multi-lake comparative analysis of the General Lake Model (GLM): Stress-testing across a global observatory network. *Environmental Modelling & Software*, 102, 274–291. <https://doi.org/10.1016/j.envsoft.2017.11.016>
- Buizza, R., & Richardson, D. (2017). 25 Years of ensemble prediction. *ECMWF Newsletter*, 153, 20–31. <https://doi.org/10.21957/bv4180>
- Cannon, D. J., Troy, C., Bootsma, H., Liao, Q., & MacLellan-Hurd, R. (2021). Characterizing the seasonal variability of hypolimnetic mixing in a large, deep Lake. *Journal of Geophysical Research: Oceans*, 126(11), 1–21. <https://doi.org/10.1029/2021JC017533>
- Carey, C. C., & Breef-Pilz, A. (2023). *Ice cover data for Falling Creek Reservoir and Beaverdam Reservoir, Vinton, Virginia, USA for 2013–2023 ver. 1*. Environmental Data Initiative. Retrieved from <https://portal-s.edirepository.org/nis/mapbrowse?scope=edi&identifier=1076&revision=1>
- Carey, C. C., Breef-Pilz, A., & Woelmer, W. M. (2023). *Time series of high-frequency sensor data measuring water temperature, dissolved oxygen, pressure, conductivity, specific conductance, total dissolved solids, chlorophyll a, phycocyanin, fluorescent dissolved organic matter, and turbidity at discrete dept*. Environmental Data Initiative. <https://doi.org/10.6073/pasta/f6bb4f5f602060dec6652ff8eb555082>
- Carey, C. C., Ibelings, B. W., Hoffmann, E. P., Hamilton, D. P., & Brookes, J. D. (2012). Eco-physiological adaptations that favour freshwater cyanobacteria in a changing climate. *Water Research*, 46(5), 1394–1407. <https://doi.org/10.1016/j.watres.2011.12.016>
- Carey, C. C., Lewis, A. S. L., Howard, D. W., Woelmer, W. M., Gantzer, P. A., Bierlein, K. A., et al. (2022a). *Bathymetry and watershed area for Falling Creek Reservoir, Beaverdam Reservoir, and Carvins Cove Reservoir*. Environmental Data Initiative. <https://doi.org/10.6073/pasta/352735344150f7e77d2bc18b69a22412>
- Carey, C. C., Woelmer, W. M., Lofton, M. E., Figueiredo, R. J., Bookout, B. J., Corrigan, R. S., et al. (2022b). Advancing lake and reservoir water quality management with near-term, iterative ecological forecasting. *Inland Waters*, 12(1), 107–120. <https://doi.org/10.1080/20442041.2020.1816421>
- Chandler, R. E. (2013). Exploiting strength, discounting weakness: Combining information from multiple climate simulators. *Philosophical Transactions of the Royal Society A: Mathematical, Physical & Engineering Sciences*, 371(1991), 20120388. <https://doi.org/10.1098/rsta.2012.0388>
- Clark, N. J., Proboste, T., Weerasinghe, G., & Magalhães, R. J. S. (2022). Near-term forecasting of companion animal tick paralysis incidence: An iterative ensemble model. *PLoS Computational Biology*, 18(2), e1009874. <https://doi.org/10.1371/journal.pcbi.1009874>
- Clayer, F., Jackson-Blake, L., Mercado-Bettín, D., Shikhani, M., French, A., Moore, T., et al. (2023). Sources of skill in lake temperature, discharge and ice-off seasonal forecasting tools. *Hydrology and Earth System Sciences*, 27(6), 1361–1381. <https://doi.org/10.5194/hess-27-1361-2023>

- Daneshmand, V., Breef-Pilz, A., Carey, C. C., Jin, Y., Ku, Y.-J., Subratie, K. C., et al. (2021). Edge-to-cloud virtualized cyberinfrastructure for near real-time water quality forecasting in lakes and reservoirs. In *2021 IEEE 17th International Conference on eScience (eScience)* (pp. 138–148). IEEE. <https://doi.org/10.1109/eScience51609.2021.00024>
- Dietze, M. C., Fox, A., Beck-Johnson, L. M., Betancourt, J. L., Hooten, M. B., Jarnevich, C. S., et al. (2018). Iterative near-term ecological forecasting: Needs, opportunities, and challenges. *Proceedings of the National Academy of Sciences*, *115*(7), 1424–1432. <https://doi.org/10.1073/pnas.1710231115>
- Di Nunno, F., Zhu, S., Ptak, M., Sojka, M., & Granata, F. (2023). A stacked machine learning model for multi-step ahead prediction of lake surface water temperature. *Science of the Total Environment*, *890*, 164323. <https://doi.org/10.1016/j.scitotenv.2023.164323>
- Doblas-Reyes, F. J., Hagedorn, R., & Palmer, T. N. (2005). The rationale behind the success of multi-model ensembles in seasonal forecasting – II. Calibration and combination. *Tellus A: Dynamic Meteorology and Oceanography*, *57*(3), 234–252. <https://doi.org/10.3402/tellusa.v57i3.14658>
- Dormann, C. F., Calabrese, J. M., Guillera-Aroita, G., Matechou, E., Bahn, V., Bartoń, K., et al. (2018). Model averaging in ecology: A review of Bayesian, information-theoretic, and tactical approaches for predictive inference. *Ecological Monographs*, *88*(4), 485–504. <https://doi.org/10.1002/ecm.1309>
- Evensen, G. (2003). The Ensemble Kalman Filter: Theoretical formulation and practical implementation. *Ocean Dynamics*, *53*(4), 343–367. <https://doi.org/10.1007/s10236-003-0036-9>
- Fearnhead, P., & Künsch, H. R. (2018). Particle filters and data assimilation. *Annual Review of Statistics and Its Application*, *5*(1), 421–449. <https://doi.org/10.1146/annurev-statistics-031017-100232>
- Feldbauer, J., Ladwig, R., Mesman, J. P., Moore, T. N., Zündorf, H., Berendonk, T. U., & Petzoldt, T. (2022). Ensemble of models shows coherent response of a reservoir's stratification and ice cover to climate warming. *Aquatic Sciences*, *84*(4), 50. <https://doi.org/10.1007/s00027-022-00883-2>
- Gilmour, I., & Smith, L. A. (1997). Enlightenment in shadows. In *Nonlinear dynamics and stochastic systems near the millenium* (pp. 335–340). AIP.
- Golub, M., Thiery, W., Marcé, R., Pierson, D., Vanderkelen, I., Mercado-Bettin, D., et al. (2022). A framework for ensemble modelling of climate change impacts on lakes worldwide: The ISIMIP Lake Sector. *Geoscientific Model Development*, *15*(11), 4597–4623. <https://doi.org/10.5194/gmd-15-4597-2022>
- Good, I. J. (1952). Rational decisions. *Journal of the Royal Statistical Society: Series B*, *14*(1), 107–114. <https://doi.org/10.1111/j.2517-6161.1952.tb00104.x>
- Goudsmit, G.-H., Burchard, H., Peeters, F., & Wüest, A. (2002). Application of k-ε turbulence models to enclosed basins: The role of internal seiches. *Journal of Geophysical Research*, *107*(C12), 3230. <https://doi.org/10.1029/2001JC000954>
- Hagedorn, R., Doblas-Reyes, F. J., & Palmer, T. N. (2005). The rationale behind the success of multi-model ensembles in seasonal forecasting – I. Basic concept. *Tellus A: Dynamic Meteorology and Oceanography*, *57*(3), 219. <https://doi.org/10.3402/tellusa.v57i3.14657>
- Hamill, T. M., Whitaker, J. S., Shlyaeva, A., Bates, G., Fredrick, S., Pegion, P., et al. (2022). The reanalysis for the global ensemble forecast system, version 12. *Monthly Weather Review*, *150*(1), 59–79. <https://doi.org/10.1175/MWR-D-21-0023.1>
- Hanson, P. C., Ladwig, R., Buelo, C., Albright, E. A., Delany, A. D., Carey, C., et al. (2023). Legacy phosphorus and ecosystem memory control future water quality in a eutrophic lake. Preprint. <https://doi.org/10.22541/essoar.168677211.11983579/v1>
- Hipsey, M. R., Bruce, L. C., Boon, C., Busch, B., Carey, C. C., Hamilton, D. P., et al. (2019). A General Lake Model (GLM 3.0) for linking with high-frequency sensor data from the Global Lake Ecological Observatory Network (GLEON). *Geoscientific Model Development*, *12*(1), 473–523. <https://doi.org/10.5194/gmd-12-473-2019>
- Howerton, E., Runge, M. C., Bogich, T. L., Borchering, R. K., Inamine, H., Lessler, J., et al. (2023). Context-dependent representation of within- and between-model uncertainty: Aggregating probabilistic predictions in infectious disease epidemiology. *Journal of the Royal Society Interface*, *20*(198), 20220659. <https://doi.org/10.1098/rsif.2022.0659>
- Huang, B., Langpap, C., & Adams, R. M. (2011). Using instream water temperature forecasts for fisheries management: An application in the Pacific northwest. *Journal of the American Water Resources Association*, *47*(4), 861–876. <https://doi.org/10.1111/j.1752-1688.2011.00562.x>
- Humphries, G. R. W., Che-Castaldo, C., Bull, P. J., Lipstein, G., Ravia, A., Carrión, B., et al. (2018). Predicting the future is hard and other lessons from a population time series data science competition. *Ecological Informatics*, *48*, 1–11. <https://doi.org/10.1016/j.ecoinf.2018.07.004>
- Hyndman, R. J., & Athanasopoulos, G. (2021). Forecasting: Principles and practice. Retrieved from [OTexts.com/fpp3](https://otexts.com/fpp3)
- IPCC (Intergovernmental Panel on Climate Change). (2023). *Technical summary. Climate change 2021 – The physical science basis*. Cambridge University Press. <https://doi.org/10.1017/9781009157896.002>
- Jackson-Blake, L. A., Clayer, F., Haande, S., Sample, J. E., & Moe, S. J. (2022). Seasonal forecasting of lake water quality and algal bloom risk using a continuous Gaussian Bayesian network. *Hydrology and Earth System Sciences*, *26*(12), 3103–3124. <https://doi.org/10.5194/hess-26-3103-2022>
- Johansson, M. A., Apfeldorf, K. M., Dobson, S., Devita, J., Buczak, A. L., Baugher, B., et al. (2019). An open challenge to advance probabilistic forecasting for dengue epidemics. *Proceedings of the National Academy of Sciences of the United States of America*, *116*(48), 24268–24274. <https://doi.org/10.1073/pnas.1909865116>
- Jolliffe, I. T., & Stephenson, D. B. (2012). In I. T. Jolliffe & D. B. Stephenson (Eds.), *Forecast verification: A practitioner's guide in atmospheric science* (2nd ed.). Wiley Blackwell.
- Jordan, A., Krüger, F., & Lerch, S. (2019). Evaluating probabilistic forecasts with scoring rules. *Journal of Statistical Software*, *90*(12), 1–37. <https://doi.org/10.18637/jss.v090.i12>
- Kirtman, B. P., Min, D., Infanti, J. M., Kinter, J. L., Paolino, D. A., Zhang, Q., et al. (2014). The North American multimodel ensemble: Phase-1 seasonal-to-interannual prediction; Phase-2 toward developing intraseasonal prediction. *Bulletin of the American Meteorological Society*, *95*(4), 585–601. <https://doi.org/10.1175/BAMS-D-12-00050.1>
- Kraemer, B. M., Chandra, S., Dell, A. I., Dix, M., Kuusisto, E., Livingstone, D. M., et al. (2017). Global patterns in lake ecosystem responses to warming based on the temperature dependence of metabolism. *Global Change Biology*, *23*(5), 1881–1890. <https://doi.org/10.1111/gcb.13459>
- La Fuente, S., Jennings, E., Gal, G., Kirillin, G., Shatwell, T., Ladwig, R., et al. (2022). Multi-model projections of future evaporation in a subtropical lake. *Journal of Hydrology*, *615*, 128729. <https://doi.org/10.1016/j.jhydrol.2022.128729>
- Lapeyrolerie, M., & Boettiger, C. (2023). Limits to ecological forecasting: Estimating uncertainty for critical transitions with deep learning. *Methods in Ecology and Evolution*, *14*(3), 785–798. <https://doi.org/10.1111/2041-210X.14013>
- Lewis, A. S. L., Woelmer, W. M., Wander, H. L., Howard, D. W., Smith, J. W., McClure, R. P., et al. (2022). Increased adoption of best practices in ecological forecasting enables comparisons of forecastability. *Ecological Applications*, *32*(2), e02500. <https://doi.org/10.1002/eap.2500>
- Lofton, M. E., Howard, D. W., Thomas, R. Q., & Carey, C. C. (2023). Progress and opportunities in advancing near-term forecasting of freshwater quality. *Global Change Biology*, *29*(7), 1691–1714. <https://doi.org/10.1111/gcb.16590>



- Long, X., Widlansky, M. J., Spillman, C. M., Kumar, A., Balmaseda, M., Thompson, P. R., et al. (2021). Seasonal forecasting skill of sea-level anomalies in a multi-model prediction framework. *Journal of Geophysical Research: Oceans*, 126(6), e2020JC017060. <https://doi.org/10.1029/2020JC017060>
- Machete, R. L., & Smith, L. A. (2016). Demonstrating the value of larger ensembles in forecasting physical systems. *Tellus A: Dynamic Meteorology and Oceanography*, 68(1), 28393. <https://doi.org/10.3402/tellusa.v68.28393>
- Mercado-Bettín, D., Clay, F., Shikhani, M., Moore, T. N., Frías, M. D., Jackson-Blake, L., et al. (2021). Forecasting water temperature in lakes and reservoirs using seasonal climate prediction. *Water Research*, 201, 117286. <https://doi.org/10.1016/j.watres.2021.117286>
- Mironov, D. V. (2021). *Parameterization of lakes in numerical weather prediction description of a lake model*. COSMO Technical Report. Offenbach am Main. [https://doi.org/10.1007/978-3-030-58292-0\\_60177](https://doi.org/10.1007/978-3-030-58292-0_60177)
- Moore, T. N., Mesman, J. P., Ladwig, R., Feldbauer, J., Olsson, F., Pilla, R. M., et al. (2021). LakeEnsembleR: An R package that facilitates ensemble modelling of lakes. *Environmental Modelling & Software*, 143, 105101. <https://doi.org/10.1016/j.envsoft.2021.105101>
- O'Hara-Wild, M., Hyndman, R., & Wang, E. (2022). fable: Forecasting models for tidy time series. R package version 0.3.2. Retrieved from <https://cran.r-project.org/package=fable>
- Olsson, F., Moore, T. N., Carey, C. C., Breef-Pilz, A., & Thomas, R. Q. (2023a). A multi-model ensemble of baseline and process-based models improves the predictive skill of near-term lake forecasts: Data, forecasts, and scores [Dataset]. Zenodo. <https://doi.org/10.5281/zenodo.8136960>
- Olsson, F., Moore, T. N., Carey, C. C., Breef-Pilz, A., & Thomas, R. Q. (2023b). OlssonF/FCRE-forecast-code: A multi-model ensemble of baseline and process-based models improves the predictive skill of near-term lake forecasts: code (v1.1.0). Zenodo. <https://doi.org/10.5281/zenodo.10679591>
- Zwart, J. A., Oliver, S. K., Watkins, W. D., Sadler, J. M., Appling, A. P., Corson-Dosch, H. R., et al. (2023). Near-term forecasts of stream temperature using deep learning and data assimilation in support of management decisions. *Journal of the American Water Resources Association*, 59(2), 317–337. <https://doi.org/10.1111/1752-1688.13093>
- Page, T., Smith, P. J., Beven, K. J., Jones, I. D., Elliott, J. A., Maberly, S. C., et al. (2018). Adaptive forecasting of phytoplankton communities. *Water Research*, 134, 74–85. <https://doi.org/10.1016/j.watres.2018.01.046>
- Pappenberger, F., Ramos, M. H., Cloke, H. L., Wetterhall, F., Alfieri, L., Bogner, K., et al. (2015). How do I know if my forecasts are better? Using benchmarks in hydrological ensemble prediction. *Journal of Hydrology*, 522, 697–713. <https://doi.org/10.1016/j.jhydrol.2015.01.024>
- Petropoulos, F., Apiletti, D., Assimakopoulos, V., Babai, M. Z., Barrow, D. K., Ben Taieb, S., et al. (2022). Forecasting: Theory and practice. *International Journal of Forecasting*, 38(3), 705–871. <https://doi.org/10.1016/j.ijforecast.2021.11.001>
- R Core Team. (2021). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Retrieved from <https://www.r-project.org>
- Read, J. S., Hamilton, D. P., Jones, I. D., Muraoka, K., Winslow, L. A., Kroiss, R., et al. (2011). Derivation of lake mixing and stratification indices from high-resolution lake buoy data. *Environmental Modelling & Software*, 26(11), 1325–1336. <https://doi.org/10.1016/j.envsoft.2011.05.006>
- Read, J. S., Jia, X., Willard, J., Appling, A. P., Zwart, J. A., Oliver, S. K., et al. (2019). Process-Guided Deep Learning predictions of lake water temperature. *Water Resources Research*, 55(11), 9173–9190. <https://doi.org/10.1029/2019WR024922>
- Renwick, K. M., Curtis, C., Kleinhesselink, A. R., Schlaepfer, D., Bradley, B. A., Aldridge, C. L., et al. (2018). Multi-model comparison highlights consistency in predicted effect of warming on a semi-arid shrub. *Global Change Biology*, 24(1), 424–438. <https://doi.org/10.1111/gcb.13900>
- Rouso, B. Z., Bertone, E., Stewart, R., & Hamilton, D. P. (2020). A systematic literature review of forecasting and predictive models for cyanobacteria blooms in freshwater lakes. *Water Research*, 182, 115959. <https://doi.org/10.1016/j.watres.2020.115959>
- Saber, A., James, D. E., & Hayes, D. F. (2020). Long-term forecast of water temperature and dissolved oxygen profiles in deep lakes using artificial neural networks conjugated with wavelet transform. *Limnology & Oceanography*, 65(6), 1297–1317. <https://doi.org/10.1002/lno.11390>
- Saloranta, T. M., & Andersen, T. (2007). MyLake—A multi-year lake simulation model code suitable for uncertainty and sensitivity analysis simulations. *Ecological Modelling*, 207(1), 45–60. <https://doi.org/10.1016/j.ecolmodel.2007.03.018>
- Smith, L. A., Cuéllar, M. C., Du, H., & Judd, K. (2010). Exploiting dynamical coherence: A geometric approach to parameter estimation in nonlinear models. *Physics Letters, Section A: General, Atomic and Solid State Physics*, 374(26), 2618–2623. <https://doi.org/10.1016/j.physleta.2010.04.032>
- Smith, L. A., Suckling, E. B., Thompson, E. L., Maynard, T., & Du, H. (2015). Towards improving the framework for probabilistic forecast evaluation. *Climatic Change*, 132(1), 31–45. <https://doi.org/10.1007/s10584-015-1430-2>
- Spence, M. A., Blanchard, J. L., Rossberg, A. G., Heath, M. R., Heymans, J. J., Mackinson, S., et al. (2018). A general framework for combining ecosystem models. *Fish and Fisheries*, 19(6), 1031–1042. <https://doi.org/10.1111/faf.12310>
- Thomas, R. Q., Figueiredo, R. J., Daneshmand, V., Bookout, B. J., Puckett, L. K., & Carey, C. C. (2020). A near-term iterative forecasting system successfully predicts reservoir hydrodynamics and partitions uncertainty in real time. *Water Resources Research*, 56(11), e2019WR026138. <https://doi.org/10.1029/2019WR026138>
- Thomas, R. Q., McClure, R. P., Moore, T. N., Woelmer, W. M., Boettiger, C., Figueiredo, R. J., et al. (2023). Near-term forecasts of NEON lakes reveal gradients of environmental predictability across the US. *Frontiers in Ecology and the Environment*, 21(5), 220–226. <https://doi.org/10.1002/fee.2623>
- Umlauf, L., Burchard, H., & Bolding, K. (2005). GOTM – Sourcecode and test case documentation. Retrieved from <http://www.gotm.net/pages/documentation/manual/stable/pdf/a4.pdf>
- Velázquez, J. A., Anctil, F., Ramos, M. H., & Perrin, C. (2011). Can a multi-model approach improve hydrological ensemble forecasting? A study on 29 French catchments using 16 hydrological model structures. *Advances in Geosciences*, 29, 33–42. <https://doi.org/10.5194/adgeo-29-33-2011>
- Viboud, C., Sun, K., Gaffey, R., Ajelli, M., Fumanelli, L., Merler, S., et al. (2018). The RAPIDD ebola forecasting challenge: Synthesis and lessons learnt. *Epidemics*, 22, 13–21. <https://doi.org/10.1016/j.epidem.2017.08.002>
- Wander, H. L., Thomas, R. Q., Moore, T. N., Lofton, M. E., Breef-Pilz, A., & Carey, C. C. (2024). Data assimilation experiments inform monitoring needs for near-term ecological forecasts in a eutrophic reservoir. *Ecosphere*, 15(2), e4752. <https://doi.org/10.1002/ecs2.4752>
- Wang, X., Hyndman, R. J., Li, F., & Kang, Y. (2022). Forecast combinations: An over 50-year review. *International Journal of Forecasting*, 39(4), 1–56. <https://doi.org/10.1016/j.ijforecast.2022.11.005>
- Ward, E. J., Holmes, E. E., Thorson, J. T., & Collen, B. (2014). Complexity is costly: A meta-analysis of parametric and non-parametric methods for short-term population forecasting. *Oikos*, 123(6), 652–661. <https://doi.org/10.1111/j.1600-0706.2014.00916.x>



- Weber, M., Rinke, K., Hipsey, M. R., & Boehrer, B. (2017). Optimizing withdrawal from drinking water reservoirs to reduce downstream temperature pollution and reservoir hypoxia. *Journal of Environmental Management*, *197*, 96–105. <https://doi.org/10.1016/j.jenvman.2017.03.020>
- Weigel, A. P., Liniger, M. A., & Appenzeller, C. (2008). Can multi-model combination really enhance the prediction skill of probabilistic ensemble forecasts? *Quarterly Journal of the Royal Meteorological Society*, *134*(134), 241–260. <https://doi.org/10.1002/qj>
- Wilson, H. L., Ayala, A. I., Jones, I. D., Rolston, A., Pierson, D., de Eyto, E., et al. (2020). Variability in epilimnion depth estimations in lakes. *Hydrology and Earth System Sciences*, *24*(11), 5559–5577. <https://doi.org/10.5194/hess-24-5559-2020>
- Wynne, J. H., Woelmer, W., Moore, T. N., Thomas, R. Q., Weathers, K. C., & Carey, C. C. (2023). Uncertainty in projections of future lake thermal dynamics is differentially driven by lake and global climate models. *PeerJ*, *11*, e15445. <https://doi.org/10.7717/peerj.15445>
- Yvon-Durocher, G., Allen, A. P., Cellamare, M., Dossena, M., Gaston, K. J., Leitao, M., et al. (2015). Five years of experimental warming increases the biodiversity and productivity of phytoplankton. *PLoS Biology*, *13*(12), e1002324. <https://doi.org/10.1371/journal.pbio.1002324>
- Zhao, F., Zhan, X., Xu, H., Zhu, G., Zou, W., Zhu, M., et al. (2022). New insights into eutrophication management: Importance of temperature and water residence time. *Journal of Environmental Sciences*, *111*, 229–239. <https://doi.org/10.1016/j.jes.2021.02.033>