

Journal of Hydrology

Enhancing the Accuracy and Generalizability of Reference Evapotranspiration Forecasting in California Using Deep Global Learning --Manuscript Draft--

Manuscript Number:	HYDROL62303
Article Type:	Research paper
Keywords:	Time series forecasting; Global learning; Deep learning; Reference evapotranspiration; Water management
Corresponding Author:	Arman Ahmadi University of California Berkeley Berkeley, CA UNITED STATES
First Author:	Arman Ahmadi
Order of Authors:	Arman Ahmadi Andre Daccache Minxue He Peyman Namadi Alireza Ghaderi Bafti Prabhjot Sandhu Zhaojun Bai Richard Snyder Tariq Kadir
Abstract:	<p>Reference evapotranspiration (ETO) indicates the atmospheric water demand and is decisive in regional to global water cycles. Like other hydrometeorological time series at monthly scales, monthly ETO time series are primarily driven by seasonality. A reliable forecast of these time series is crucial for sustainable water resources planning and management. Although the current and previous research on hydrometeorological time series forecasting focuses on local learning (i.e., training a forecasting model on a single time series and using the trained model for future time steps), our study points to the superiority of the global learning scheme. In global learning, the forecasting model is trained over a pool of multiple time series and tested on new instances. To quantify how deep learning (DL) models can benefit from global learning in hydrometeorological forecasting, our study uses monthly ETO time series from 55 standardized weather stations in the Central Valley of California. This study 1) compares the performance of statistical and deep forecasting models in local learning, 2) quantifies the performance improvement of DL models in global versus local learning, and 3) automatically optimizes hyperparameters of the best-performing DL models to achieve state-of-the-art forecasting accuracy. Our findings reveal that while statistical models such as Holt-Winters outperform DL models in local learning, global learning can unleash the true potential of high-capacity DL models such as N-BEATS and N-HiTS. This approach results in RMSE values below 10 mm/month for one-year-ahead forecasts on unseen stations. In addition to superior accuracies, global learning enhances the generalizability of DL models, making them applicable to ungauged locations and recently established weather stations. Our findings also point to the benefits of automatic hyperparameter optimization in deep global forecasting.</p>
Suggested Reviewers:	<p>Mojtaba Sadegh, Ph.D. Associate Professor, Boise State University mojtabasadegh@boisestate.edu Expert in hydrology and machine learning</p> <p>Amir AghaKouchak, Ph.D. Professor, University of California Irvine amir.a@uci.edu Expert in hydrometeorology</p> <p>John Abatzoglou, Ph.D. Professor, University of California Merced jabatzoglou@ucmerced.edu</p>

	<p>Mohammad Reza Alizadeh, Ph.D. Postdoctoral Associate, Massachusetts Institute of Technology alizade@mit.edu Expert in hydrology and machine learning</p>
	<p>Yun Yang, Ph.D. Assistant Professor, Cornell University yy2356@cornell.edu Expert in hydrology and water resources</p>

1

2

Enhancing the Accuracy and Generalizability of Reference
Evapotranspiration Forecasting in California Using Deep Global Learning

3

4

Arman Ahmadi^{1,*}, Andre Daccache², Minxue He³, Peyman Namadi³, Alireza Ghaderi Bafti⁴,
Prabhjot Sandhu³, Zhaojun Bai⁵, Richard L. Snyder⁶, Tariq Kadir³

5

6

7

¹ Department of Environmental Science, Policy and Management, University of California, Berkeley, Berkeley, CA, USA.

8

² Department of Biological and Agricultural Engineering, University of California, Davis, Davis, CA, USA.

9

³ California Department of Water Resources, Sacramento, CA, USA.

10

⁴ Department of Civil Engineering, Faculty of Engineering, Ferdowsi University of Mashhad, Mashhad, Iran.

11

⁵ Department of Computer Science, University of California, Davis, Davis, CA, USA.

12

⁶ Department of Land, Air and Water Resources, University of California, Davis, CA 95616, USA

13

* Corresponding Author

14

Abstract

15

16

Reference evapotranspiration (ET_O) indicates the atmospheric water demand and is decisive in regional to global water cycles. Like other hydrometeorological time series at monthly scales,

monthly ET_0 time series are primarily driven by seasonality. A reliable forecast of these time series is crucial for sustainable water resources planning and management. Although the current and previous research on hydrometeorological time series forecasting focuses on local learning (i.e., training a forecasting model on a single time series and using the trained model for future time steps), our study points to the superiority of the global learning scheme. In global learning, the forecasting model is trained over a pool of multiple time series and tested on new instances. To quantify how deep learning (DL) models can benefit from global learning in hydrometeorological forecasting, our study uses monthly ET_0 time series from 55 standardized weather stations in the Central Valley of California. This study 1) compares the performance of statistical and deep forecasting models in local learning, 2) quantifies the performance improvement of DL models in global versus local learning, and 3) automatically optimizes hyperparameters of the best-performing DL models to achieve state-of-the-art forecasting accuracy. Our findings reveal that while statistical models such as Holt-Winters outperform DL models in local learning, global learning can unleash the true potential of high-capacity DL models such as N-BEATS and N-HiTS. This approach results in RMSE values below 10 mm/month for one-year-ahead forecasts on unseen stations. In addition to superior accuracies, global learning enhances the generalizability of DL models, making them applicable to ungauged locations and recently established weather stations. Our findings also point to the benefits of automatic hyperparameter optimization in deep global forecasting.

Keywords:

Time series forecasting, Global learning, Deep learning, Reference evapotranspiration, Water management

1. Introduction

Hydrometeorological time series forecasting, particularly at the monthly scale, is crucial in water resources planning and management (Milly et al., 2008; Li et al., 2017). Accurate forecasts of variables such as precipitation, streamflow, and evapotranspiration are fundamental for effective decision-making in various sectors, including agriculture, hydropower generation, reservoir operations, and environmental conservation (Gleick, 2003; Chen et al., 2024; Quinn et al., 2024). Monthly time series forecasts are critical as they align with many water management activities' operational and planning timescales, balancing shorter-term weather predictions and longer-term climate projections (Mehdizadeh et al., 2019; Le et al., 2024).

Reference evapotranspiration (ET_0) is the amount of water that vaporizes from a virtual 0.12 m tall, vegetated surface using empirical equations to estimate canopy and aerodynamic resistances in a modified Penman-Monteith equation (Allen et al., 2005, 2006). In reality, ET_0 is approximately equal to the evapotranspiration of a broad expanse of a healthy, well-watered, cool-season grass pasture close to 0.12 m height (Allen et al., 1998). Assuming the weather data are measured over a large expanse of well-watered pasture maintained at approximately 0.12 m height, ET_0 is a meteorological variable driven by air temperature, humidity, solar radiation, and wind speed (Allen et al., 1998, 2005, 2006; Ahmadi et al., 2022). The standardized ET_0 is used as a direct measure of atmospheric evaporative power or water demand, making it a significant component in regional to global water cycles. Like other hydrometeorological time series at

monthly time scales, monthly ET_O time series are primarily seasonality-driven, with a trivial trend component. This study focuses on one year-ahead monthly ET_O forecasting as an instance of hydrometeorological time series forecasting useful for water resources planning and management.

The Central Valley of California is chosen as the case study for this research. The Central Valley is one of the world's most productive agricultural regions, providing over a third of the vegetables and two-thirds of the fruits and nuts consumed in the United States (Hanak et al., 2011). Due to its extensive agricultural activities and limited in-season precipitation, the region is primarily dependent on irrigation for crop production. Optimal irrigation planning is imperative for sustainable water resources and food production in the Central Valley, as excessive irrigation depletes valuable ground and surface water resources and exacerbates soil salinity issues, further threatening agricultural sustainability. Consequently, accurate forecasts of ET_O are crucial for efficient water management and irrigation practices in this region.

A considerable amount of research exists on using data-driven approaches to estimate ET_O (Zhang et al., 2020; Chen et al., 2020; Dong et al., 2022; Ahmadi et al., 2024). Numerous studies have focused on forecasting ET_O with data-driven methodologies. For instance, Gocić et al. (2015) explored soft computing methods for monthly ET_O forecasting, while Karbasi et al. (2022) investigated a hybrid deep learning model for weekly ET_O forecasting. Ferreira and da Cunha (2020) utilized deep learning (DL) models for multi-step ahead daily ET_O forecasts, and Chia et al. (2022) employed deep neural networks for long-term forecasting of monthly mean ET_O . In a recent study, Ahmadi et al. (2023) analyzed the accuracy, complexity, and data efficiency of several statistical, machine learning (ML), and DL methods for monthly ET_O

forecasting. The current study is a follow-up of this research, focusing on global learning and cutting-edge DL models for 12-month ahead ET_0 forecasting.

Contrary to local learning, which trains a forecasting model on a single time series, global learning schemes use a pool of multiple time series from different sources as the training set (Bandara et al., 2020; Smyl, 2020; Salinas et al., 2020). Global learning is specifically useful for hydrometeorological time series forecasting when several stations in a climatically homogeneous region measure the variable of interest. For example, this study uses the ET_0 time series from 47 weather stations in the Central Valley as its training set and tests the performance of the globally learned DL models on another eight unseen stations that constitute the test set. By leveraging data from where it is available, global learning enables reliable forecasting in poorly gauged regions and recently established stations.

This research aims to quantitatively analyze the effects of global learning on the performance of deep forecasting models in monthly ET_0 forecasting. The research objectives are: 1) to assess the forecasting accuracy of statistical and advanced DL models in local learning scheme, 2) to compare the DL models' performance in local versus global learning, and 3) to achieve state-of-the-art deep forecasting accuracy through automatic hyperparameter optimization. In addition to quantifying performance improvements, our study discusses how global learning can enhance the generalizability of deep forecasting models. We hypothesize that globally learned DL models outperform both locally learned DL models and the best-performing statistical models while automatically tuned deep forecasting models achieve the highest accuracy. We also hypothesize that globally learned models are more generalizable to unseen instances.

2. Study area and dataset

The Central Valley of California is the case study of this research (figure 1). The Central Valley, relying on extensive irrigation, is one of the most agriculturally productive regions worldwide (Hanak, 2011). This study uses monthly ET_O time series at 55 locations in the Central Valley from the *California Irrigation Management Information System (CIMIS)* program. CIMIS comprises over 145 automated and standardized weather stations to aid irrigators and water managers in planning and decision-making. CIMIS employs the Penman-Monteith equation and a modified version of Penman's equation to calculate ET_O . Hourly weather data is used to determine hourly ET_O , which is then summed over 24 hours (midnight to midnight local time) to estimate daily ET_O . The monthly ET_O values reported by the CIMIS portal are aggregates of these daily values in metric units (mm). Further details about CIMIS data and the Penman-Monteith equation can be found in Ahmadi et al. (2022).

According to the CIMIS website (accessed July 2023), 42 inactive and 38 active stations are in the Central Valley (<https://cimis.water.ca.gov/>). After reviewing the maintenance and data quality reports on the CIMIS website, we eliminated eight stations: three stations had poor data quality, two stations had bare reference surfaces, two stations had alfalfa reference surfaces, and one was a non- ET_O site. Of the remaining 72 stations, 17 were excluded due to insufficient time series data. Ultimately, 55 stations with more than six consecutive years of data remained. Monthly data from January 1986 to June 2023 is downloaded for all stations. Data from eight active stations (i.e., stations with data available until June 2023) are used as the test set, and the remaining 47 stations constitute the training set (figure 1). The zoning classification of figure 1 defines California's homogeneous ET_O zones. More information about these zones is found in Ahmadi et al. (2022).

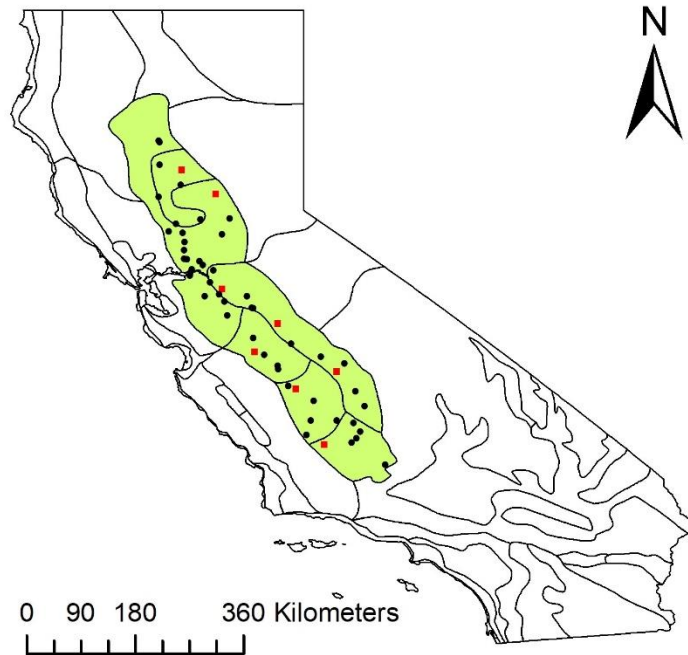


Figure 1. The study area (the Central Valley of California, highlighted in light green) and CIMIS stations. Black circles represent the training stations, and red squares represent the test stations.

As figure 1 shows, both training and test sets are widely spread over the Central Valley and uniformly distributed. To ensure the representativeness of the test set and its statistical similarity to the training set, we compared the ET_O distributions in these sets (figure 2). Figure 2 suggests that the ET_O distributions in training and test sets are reasonably similar.

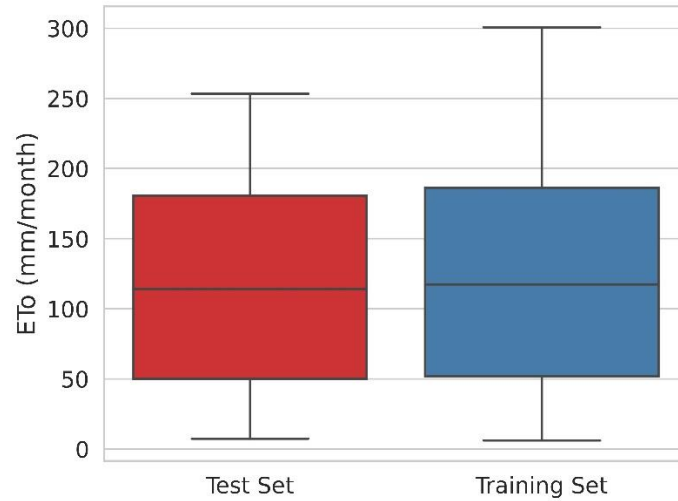


Figure 2. Distribution of reference evapotranspiration in the test and training sets

3. Methodology

3.1. Time series forecasting platform

We employed Darts, a Python library for time series manipulation and forecasting (Herzen et al., 2022). Darts encompasses various forecasting models, ranging from statistical methods to state-of-the-art deep neural network architectures specifically designed for time series forecasting. Readers are directed to Herzen et al. (2022) for further details on this library.

3.2. Local vs. global learning

Local learning refers to training a separate forecasting model for each individual time series. Therefore, in the local learning scheme of our study, a separate model is trained for each station

142 in the test set. All the stations in the test set have data available through the end of June 2023.
143 The last three years of data (i.e., July 2020 to June 2023) are used to test the models'
144 performance, and models are trained over the remaining time steps (figure 3). The forecasting
145 horizon for all models is 12 months (i.e., 12 time steps), and the stride is one month (figure 3).
146 Figure 3 demonstrates how the model is continually updated and retrained as new data becomes
147 available, always predicting 12 months ahead while moving forward in monthly increments
148 through the 36-month (3-year) testing period.

149 Contrary to local learning, a single forecasting model is trained over data pooled across multiple
150 time series in global learning. The global learning scheme of our study trains a single model
151 using the data from all 47 stations in the training set (figure 1). Data from the eight stations in the
152 test set are not introduced to this model in the training stage. Only after training is the model
153 used to forecast the last three years of data for each of the eight stations in the test set. After the
154 training stage, no further fine-tuning is performed on the globally learned models.

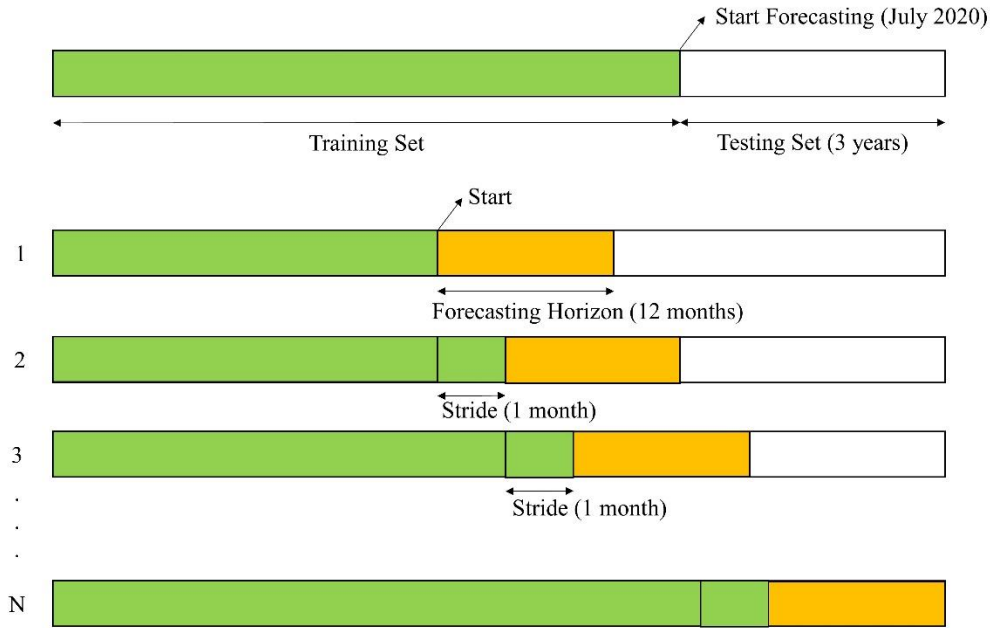


Figure 3. Data splitting and forecasting time steps for local learning

3.3. Statistical forecasting models

3.3.1. Seasonal autoregressive integrated moving average (SARIMA)

The ARIMA model is widely recognized in time series forecasting for its combination of autoregressive (AR), differencing (I), and moving average (MA) components (Contreras et al., 2003; Hyndman and Athanasopoulos, 2018). Differencing is employed to eliminate non-stationarity by computing the differences between consecutive observations. The AR component forecasts the variable of interest using a linear combination of its past values, while the MA component utilizes past forecast errors instead of the variable's past values in a regression-like manner. Although ARIMA is a non-seasonal model, it sets the foundation for the seasonal

166 ARIMA (SARIMA) model, which integrates additional seasonal terms into the ARIMA
167 framework.

168 In this study, we focus on the SARIMA model, which extends ARIMA by incorporating four
169 additional hyperparameters: P, D, and Q, representing the seasonal order for the AR,
170 differencing, and MA components, respectively, and m, which denotes the periodicity or number
171 of time steps in a complete seasonal period. For monthly data, m is set to 12 (Hyndman and
172 Athanasopoulos, 2018).

173 To optimize the hyperparameters of SARIMA, we used the *pmdarima* Python library. Each
174 station in the test set had its hyperparameters optimized individually, after which these station-
175 specific parameters were utilized to train the SARIMA models. The *AutoARIMA* function from
176 the *pmdarima* library was used to identify the optimal set of parameters for the SARIMA model,
177 resulting in a single fitted model. The maximum values for p, q, P, and Q were set to 5, and m
178 was fixed at 12, aligning with our monthly data.

179 The Akaike Information Criterion (AIC) was employed to select the best model, with an alpha
180 level of 0.05 for statistical significance. The Kwiatkowski–Phillips–Schmidt–Shin (KPSS) unit
181 root test was used to assess stationarity, and the Osborn-Chui-Smith-Birchenhall (OCSB)
182 seasonal unit root test was performed for seasonal stationarity. Model parameters were optimized
183 using the stepwise algorithm described by Hyndman and Khandakar (2008), and the limited-
184 memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS) algorithm with optional box constraints
185 was used as the optimization algorithm.

3.3.2. *Holt-Winters' exponential smoothing*

Exponential smoothing, a statistical forecasting technique introduced in the late 1950s (Holt, 2004; Winters, 1960), utilizes weighted averages of past data points to predict future values. This univariate method assigns exponentially decreasing weights to older observations, giving more significance to more recent data (Hyndman and Athanasopoulos, 2018). Additional details about this method can be found in Kalekar (2004) and Hyndman and Athanasopoulos (2018). Our study applies an additive approach to both trend and seasonal components, with a seasonal period set to 12. Moreover, the trend component is damped for better accuracy.

3.3.3. *Theta method*

The Theta model, introduced by Assimakopoulos and Nikolopoulos (2000), is a univariate forecasting approach that modifies the local curvature of time series data by applying a coefficient known as "Theta" (a real number) to the second differences in the data. This method decomposes the original series into multiple lines, extrapolates each line using suitable forecasting models, and then combines the predictions to arrive at a final forecast. This study employs the 4Theta model, an enhanced variation of the original Theta method (Spiliotis et al., 2020). After a manual search, we selected $\Theta = 2$ for this model. The seasonal period is specified as 12, with a multiplicative type of seasonality. The Theta lines are integrated using an additive model, and the trend is modeled linearly. For additional details on the Theta and 4Theta models, readers should consult Assimakopoulos and Nikolopoulos (2000) and Spiliotis et al. (2020).

206 **3.4. Machine learning model (LightGBM)**

207 In this study, we utilize LightGBM as a machine learning forecasting tool. LightGBM,
208 introduced by Ke et al. (2017) and developed by Microsoft as a free and open-source framework,
209 offers an efficient implementation of the gradient boosting algorithm and optimizes memory
210 usage. Gradient boosting is an ensemble technique combining decision tree models to form a
211 more robust predictive model. The fitting process involves a gradient descent optimization
212 algorithm that minimizes the loss gradient as the model parameters are adjusted. Within machine
213 learning literature and competitions, gradient boosting and decision tree-based models have been
214 shown to outperform other regression algorithms for tabular data (Shwartz-Ziv and Armon,
215 2022). Furthermore, various studies highlight that while gradient boosting algorithms perform
216 similarly in accuracy and runtime, LightGBM often excels (Al Daoud, 2019). Consequently, this
217 study adopts LightGBM as the primary machine learning model, utilizing 24 prior time steps.
218 LightGBM forecasts the next 12 time steps with a Multi-Input Multi-Output (MIMO) strategy.
219 For further details on the LightGBM methodology, readers can refer to Ke et al. (2017) and Al
220 Daoud (2019).

221 **3.5. Deep learning models**

222 All DL models use a Multi-Input Multi-Output (MIMO) strategy. We performed manual
223 hyperparameter tuning for all DL models. DL models take 24 previous time steps as input and
224 output the next 12 time steps. We used Google Colab GPU to train DL models.

3.5.1. N-BEATS

Neural basis expansion analysis for interpretable time series forecasting (N-BEATS), a deep neural network architecture introduced by Oreshkin et al. (2019), leverages backward and forward residual connections and a deeply stacked array of fully connected layers. Initially developed in 2019 to address univariate time series forecasting, N-BEATS is known for its rapid training capabilities and state-of-the-art performance across various datasets. For additional details on the N-BEATS architecture, readers are encouraged to consult Oreshkin et al. (2019). In this study, we adopted the generic architecture described by Oreshkin et al. (2019). Specifically, we employed four stacks, each containing four blocks. Each block included four fully connected layers preceding the final backcast-forecast forking layer, with each layer comprising 16 neurons. The expansion coefficient dimension was set to five, and the rectified linear unit (ReLU) function was used as the activation function for the encoder/decoder intermediate layer. Our grid search indicated that N-BEATS performed optimally without dropout; thus, the dropout probability was set to zero. The model was trained over 100 epochs with a batch size of 32.

3.5.2. Long short-term memory (LSTM)

Long Short-Term Memory (LSTM), introduced by Hochreiter and Schmidhuber (1997), is a type of recurrent neural network (RNN). RNN models are well-suited for tackling problems involving sequential input data, such as time series. However, traditional RNN models often struggle with long-term dependencies, failing to retain information over extended sequences. LSTM architecture is designed to address this issue, making it a significant advantage for this model.

For further details on LSTM, readers are referred to Hochreiter and Schmidhuber (1997) and Van Houdt et al. (2020). Our LSTM model features a single recurrent layer with 12 features in its hidden state. The dropout rate for this model is set to zero. The model is trained over 1,000 epochs with a batch size of 8.

3.5.3. Temporal convolutional network (TCN)

While convolutional neural networks (CNNs) are often associated with raster data, they can also be adapted for sequential data with appropriate modifications. The temporal convolutional network (TCN), introduced by Bai et al. (2018), is a convolutional architecture designed explicitly for sequence modeling. In this study, we utilize a dilated TCN for forecasting. For more details on this model, readers are referred to Bai et al. (2018). Our model features a kernel size of 6 and 18 filters. The base of the exponent that determines the dilation at each level is set to 2. We applied weight normalization to the model and employed a dropout rate of 0.1. The model was trained over 1,000 epochs with a batch size of 32.

3.5.4. Transformer model

The Transformer model, introduced by Vaswani et al. (2017), represents a state-of-the-art deep learning architecture. Unlike traditional architectures, the Transformer does not depend on recurrence or convolutions to produce its output; instead, it utilizes an encoder-decoder structure. Central to this architecture is the multi-head attention mechanism, which can simultaneously focus on different positions in the sequence, making the Transformer particularly suitable for time series forecasting. Additionally, its highly parallelizable nature makes it well-suited for

training on GPUs. For further details on the Transformer model, readers are encouraged to consult Vaswani et al. (2017).

In our implementation, the transformer model's encoder and decoder inputs have 16 features, with one encoder layer and one decoder layer. The multi-head attention mechanism utilizes four heads. The feedforward network dimension is set to 128, and we use the ReLU activation function in the encoder/decoder intermediate layers. Based on grid search results, the dropout rate is set to 0.1. The model is trained over 1,200 epochs with a batch size of 32.

3.5.5. Temporal fusion Transformer (TFT)

The temporal fusion transformer (TFT), introduced by Lim et al. (2021), is a cutting-edge deep learning architecture designed for interpretable multi-horizon time series forecasting. This novel attention-based model integrates recurrent layers for local processing and a self-attention mechanism for capturing long-term dependencies. The TFT can learn temporal relationships at varying scales and includes specialized components for selecting relevant features. Readers are referred to Lim et al. (2021) for an in-depth understanding of this architecture.

In our study, we configured the TFT with a hidden state size of 16 and a hidden size of 8 for processing continuous variables. The architecture includes one layer for each LSTM encoder and decoder. The model employs four attention heads, with the multi-head attention query applied exclusively to the future (decoder) part. A gated residual network serves as the feedforward component. The training was conducted using PyTorch's mean squared error (MSE) loss function. The model was trained over 700 epochs with a batch size of 32.

3.5.6. *N-HiTS*

Neural hierarchical interpolation for time series forecasting (N-HiTS) is a novel DL model introduced by Challu et al. (2023). N-HiTS is like N-BEATS but aims to enhance performance while reducing computational costs by incorporating multi-rate sampling of inputs and multi-scale interpolation of outputs. By assembling its predictions sequentially and emphasizing components with different frequencies and scales, N-HiTS addresses the challenges of prediction volatility and computational complexity. Extensive experiments have demonstrated that N-HiTS significantly improves accuracy and efficiency, outperforming state-of-the-art Transformer models by reducing computation time dramatically. Detailed information about N-HiTS can be found in Challu et al. (2023).

In our N-HiTS architecture, we employed four stacks, each containing four blocks. Each block included four fully connected layers, with each layer comprising 32 neurons. The ReLU function was the activation function for the encoder/decoder intermediate layer. Similar to our N-BEATS model, the dropout probability was set to zero. The model was trained over 100 epochs with a batch size of 32.

3.5.7. *N-Linear and D-Linear*

N-Linear and D-Linear models are introduced by Zeng et al. (2023) as simple yet effective alternatives to Transformer-based models for long-term time series forecasting. Unlike Transformers, these "embarrassingly" simple models avoid the temporal information loss associated with self-attention mechanisms by employing a one-layer linear architecture. Despite their simplicity, extensive experiments show that N-Linear and D-Linear models outperform

sophisticated Transformer models significantly in accuracy and efficiency. More details can be found in the original study introducing these models (Zeng et al., 2023). For our D-Linear model, the size of the kernel for the moving average is set to 24. We trained both N-Linear and D-Linear models over 100 epochs with a batch size of 32.

3.5.8. Time-series Dense Encoder (TiDE)

Time series dense encoder (TiDE) developed by Das et al. (2023) resembles Transformers but strives for enhanced performance and reduced computational cost by utilizing multilayer perceptron (MLP)-based encoder-decoders without the use of attention mechanisms. This approach maintains the simplicity and speed of linear models while effectively handling nonlinear dependencies, achieving near-optimal error rates for linear dynamical systems. Empirically, TiDE matches or surpasses previous methods on established long-term time-series forecasting benchmarks, operating faster than the best Transformer models. Detailed information about TiDE can be found in Das et al. (2023).

Our TiDE architecture consists of one residual block in the encoder layer and one in the decoder layer, while the width of the layers in the residual blocks is set to 128. The width of the layers in the past and the future covariate projection residual blocks is 4. The output of the decoder layer is 16-dimensional. The temporal decoder layers have a width of 32 units. The dropout probability is 0.1. The model is trained over 100 epochs with a batch size of 32.

3.6. Automatic hyperparameter optimization

We employed Optuna, a hyperparameter optimization framework for automatic hyperparameter tuning of the best-performing DL models in global learning (Akiba et al., 2019). Optuna allows users to construct the hyperparameter search space dynamically by offering a define-by-run API. Optuna has a user-friendly setup, and it provides efficient sampling and pruning algorithms for customization. Detailed information about Optuna and its algorithm can be found in Akiba et al. (2019). We used the Davis CIMIS station from the training set to optimize hyperparameters using Optuna, where the last three years of its data are used as the validation set and the remaining as the train set.

4. Results and discussion

4.1. Local learning

Figure 4 shows the performance of locally learned models on the test set. Holt-Winters is the most accurate forecasting model based on RMSE, MAE, and R2 scores. This finding aligns with Ahmadi et al. (2023) analysis of a similar dataset, where the Holt-Winters model's accuracy was comparable with the best-performing DL forecasting models.

As figure 4 suggests, DL models' performance in local learning falls behind simpler statistical models. This might be because DL models have more complex algorithms involving hundreds or even thousands of trainable parameters (readers are referred to Ahmadi et al. (2023) to learn more about the number of trainable parameters in DL and statistical models). We need extensive datasets to train these big models effectively, which is hard to provide in the case of local learning.

346 Another reason for the superiority of simpler models is the inherent simplicity of monthly ET_0
347 and, more broadly, monthly hydrometeorological time series. These time series are usually
348 dominated by clear seasonal patterns due to predictable climatic cycles with few nonlinear or
349 complex interactions (Wilks, 2011). The Holt-Winters model excels in this context as it is
350 explicitly designed to handle periodic data, incorporating specific components for modeling
351 seasonality, trend, and level (Holt, 2004; Winters, 1960). Holt-Winters is less prone to overfitting
352 as a simpler model, making it more effective for the straightforward patterns prevalent in
353 monthly ET_0 time series. In contrast, advanced DL models like N-BEATS, TiDE, and
354 Transformers have high capacity and can capture complex patterns, but they are prone to
355 overfitting when applied to these relatively simple patterns (Goodfellow et al., 2016).

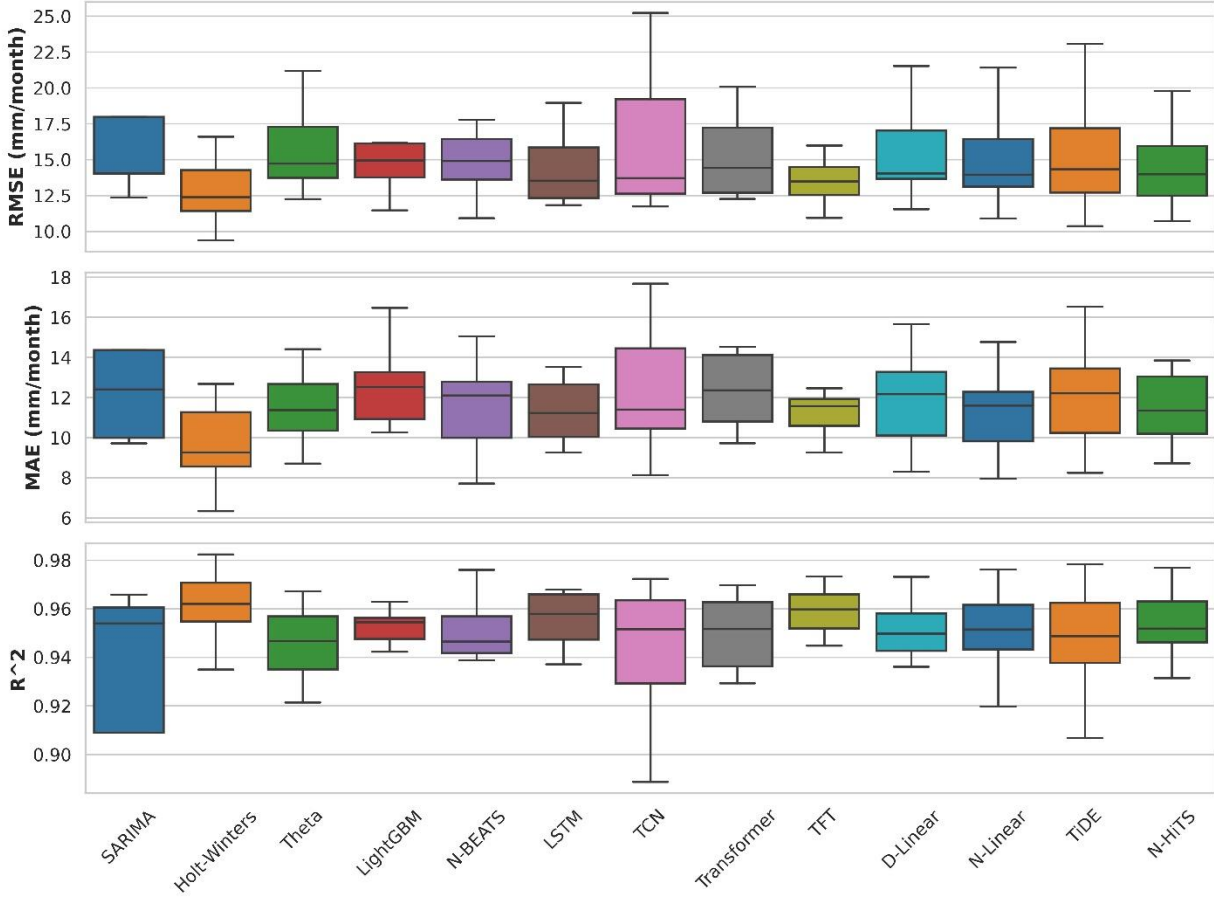


Figure 4. Results of local learning on the test set

4.2. Global learning

As figure 5 shows, global learning significantly enhances the performance of N-BEATS, TiDE, and N-HITS. The performance of LSTM, D-Linear, and N-Linear stays about the same for local and global learning schemes, while the performance of TCN, Transformer, and TFT improves meaningfully in global learning. As mentioned earlier, the DL models require vast amounts of data to be trained effectively. Contrary to local learning, where models are trained over a single time series, in global learning, DL models learn from multiple time series, which, as figure 5 suggests, results in higher forecasting accuracy.

366 N-BEATS, TiDE, and N-HiTS have high-capacity architectures capable of capturing complex
367 patterns in data (Oreshkin et al., 2019; Das et al., 2023; Challu et al., 2023). When employed in
368 global learning, these models can capture and generalize broader patterns across different time
369 series in the dataset, boosting performance. Conversely, LSTM is designed to capture long-range
370 dependencies in sequential data and might be sensitive to the length and type of sequences,
371 which can vary between local and global learning schemes (Hochreiter and Schmidhuber, 1997).
372 Furthermore, the linear architectures of D-Linear and N-Linear models may already effectively
373 capture linear relationships within the local context, offering limited gains from additional global
374 patterns (Zeng et al., 2023). These simpler models may thus experience a performance plateau,
375 as they have less room for improvement when introducing more training data and adopting
376 global schemes.

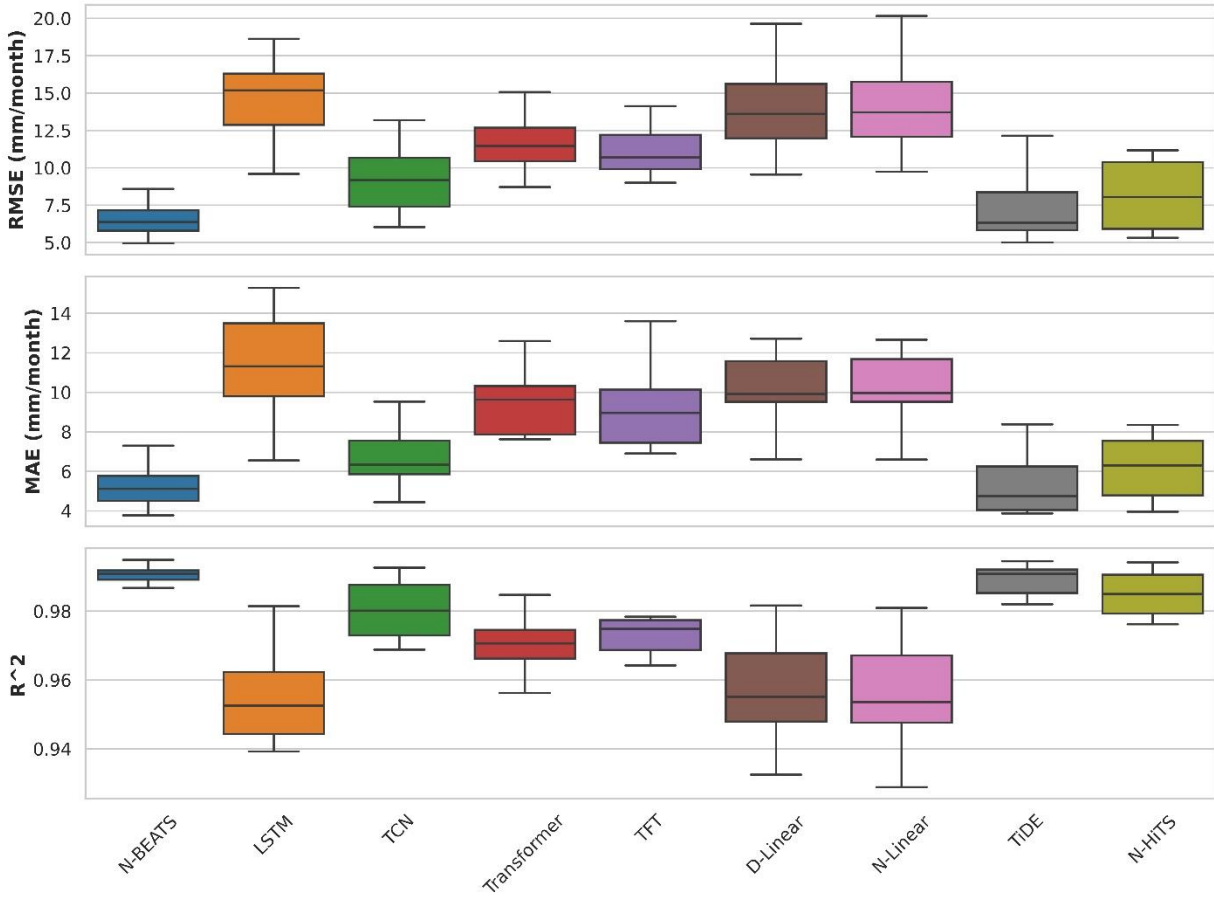


Figure 5. Results of global learning on the test set

4.3. Automatic hyperparameter optimization

We employed Optuna for automatic hyperparameter tuning of the best-performing DL models: N-BEATS, TiDE, and N-HITS. As figure 6 demonstrates, Optuna improved the performance of N-BEATS and N-HITS models. However, the TiDE model's accuracy is higher with manual hyperparameter, and automatic hyperparameter optimization did not enhance TiDE's performance (figure 6). The discrepancy of hyperparameter optimization effectiveness for these

386 models can be rooted in their characteristics. Different models have varying sensitivities to
387 hyperparameter settings. Additionally, optimization landscapes have various levels of
388 complexity for different models, with potential ruggedness and local minima, resulting in higher
389 performance in manual tuning.

390 It should be noted that figure 6 also compares the best-performing DL models in the global
391 learning scheme with Holt-Winters, which had the best accuracy in local learning. As can be
392 seen, when trained globally, DL models outperform Holt-Winters by a significant margin.
393 Automatic hyperparameter tuning enhances the performance of two of three models, making the
394 difference between their and Holt-Winter's overall accuracies even larger.

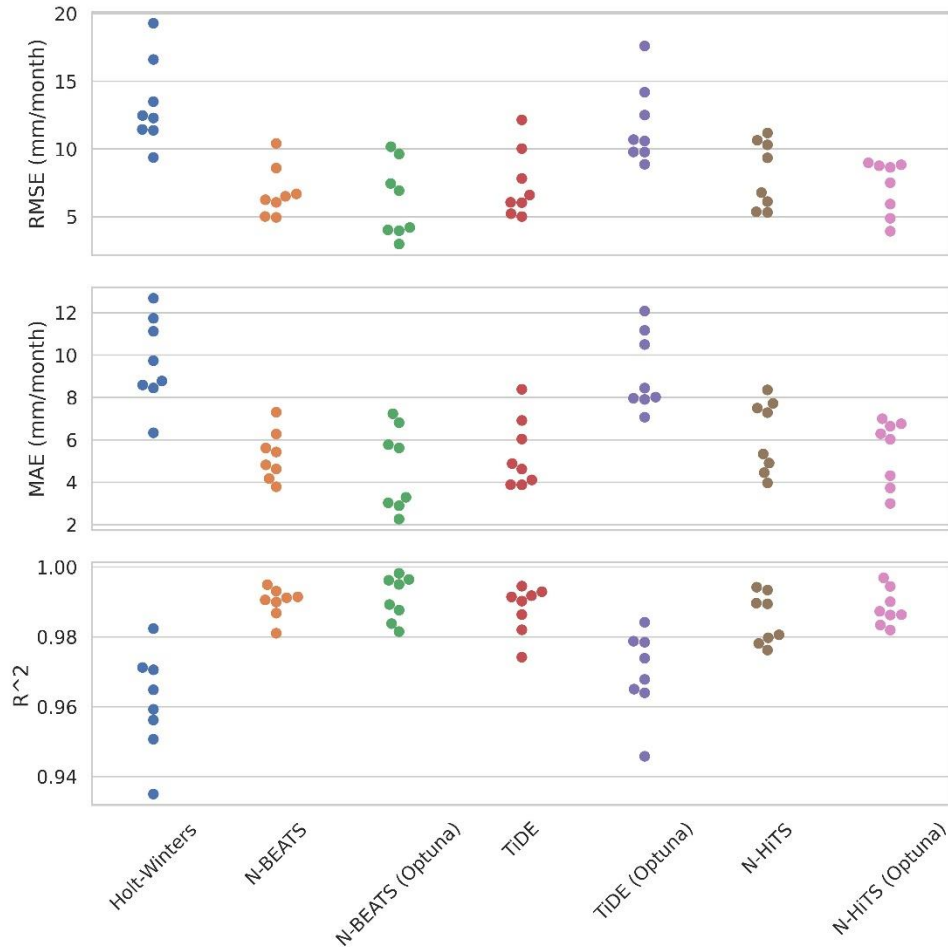


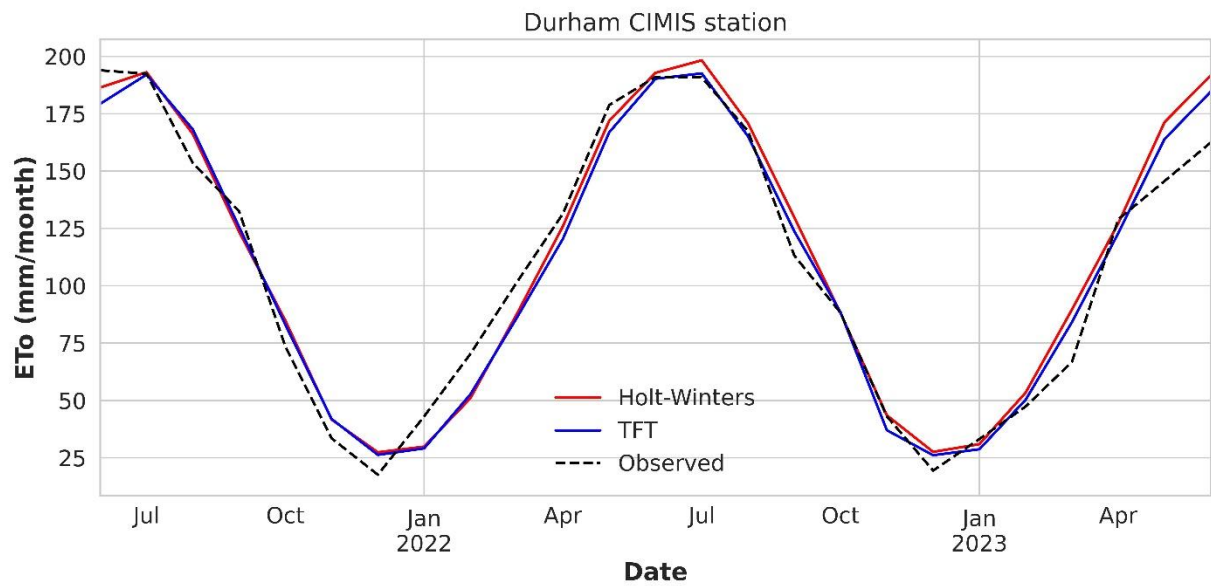
Figure 6. Results of best-performing globally learned deep learning models on the test set with and without automatic hyperparameter tuning by Optuna

4.4. Local versus global learning

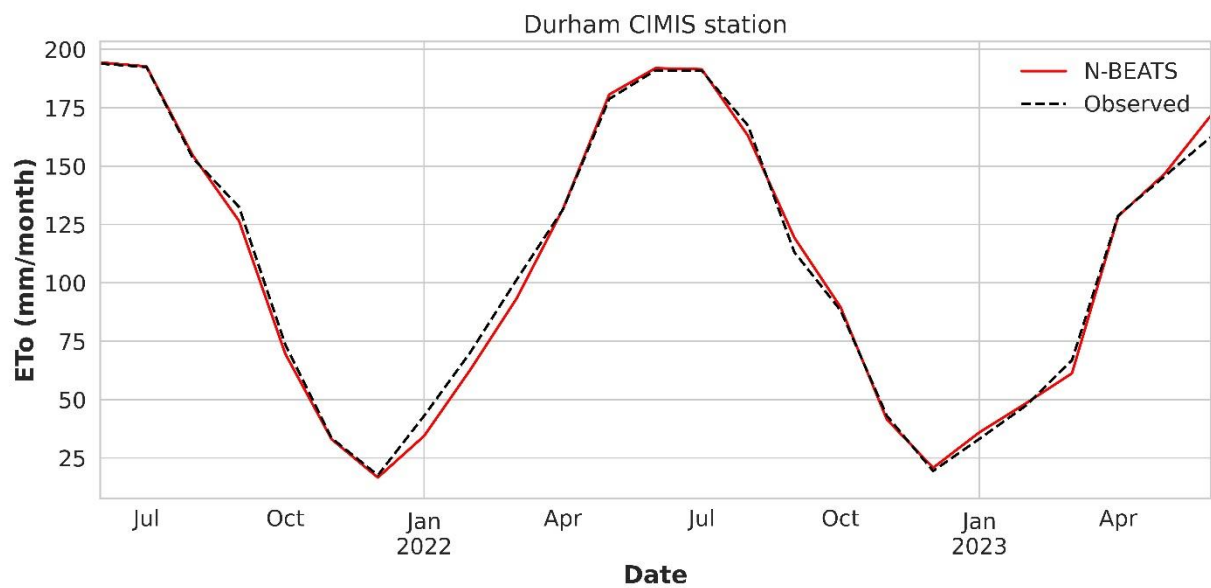
Figure 7 compares the forecasting accuracy in local versus global learning. We used the Durham CIMIS station to depict how these learning schemes affect the forecasted time series. Durham CIMIS station is one of the stations in the test set with data available from 1986 to 2023. As figure 7-a illustrates, TFT -the best-performing DL model in local learning- is very similar to Holt-Winters, which is a statistical forecasting model with a much simpler algorithm and lower

computational cost. Conversely, when we use the global learning scheme, DL forecasting models generate accurate forecasts far closer to the observed values (figures 7-b and 7-c). In other words, we can infer that global learning unleashes the true potential of advanced deep forecasting models such as N-BEATS and N-HiTS. As mentioned earlier, N-BEATS and N-HiTS have high-capacity architecture and can model complex nonlinear patterns effectively. However, as our findings suggest, the local learning scheme cannot provide them with enough data.

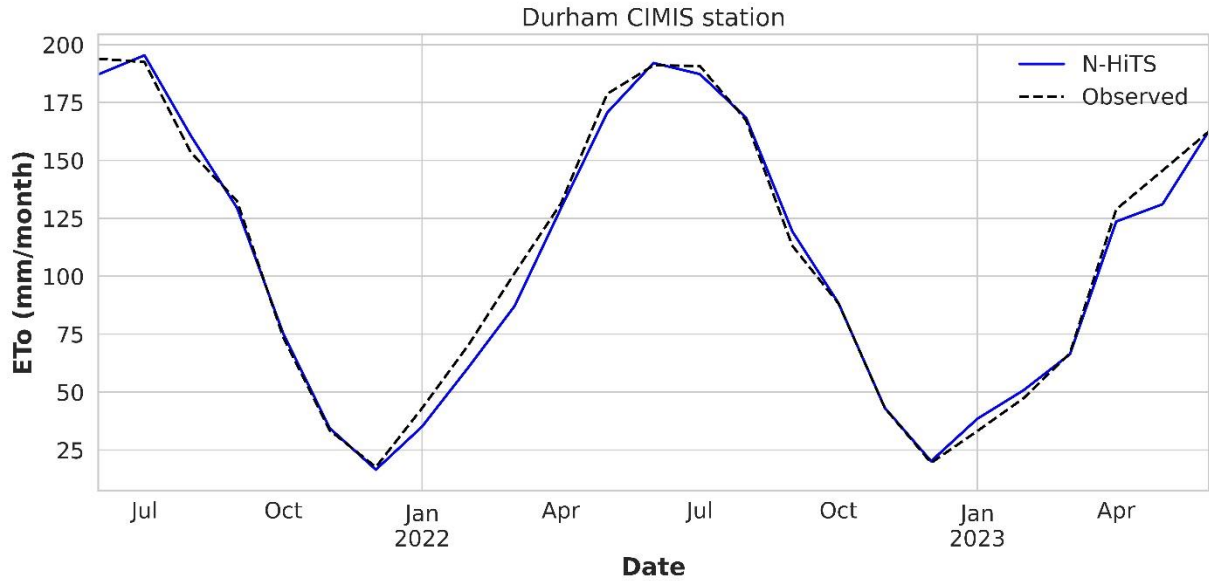
In addition to improving accuracy, using global learning for DL models can enhance their generalizability. Trained over a pooled dataset of hydrometeorological time series with characteristics like the target, these globally learned models can work reliably on unseen instances. This is especially important in hydrometeorological time series in ungauged locations and newly established measurement stations. Deep global forecasting models can leverage large data availabilities from adjacent locations or even stations from other regions with matching climatic characteristics. Monthly hydrometeorological time series are usually seasonality-driven; therefore, global learning can open the door for mixing different variables in the training set, making the models more nuanced to complex patterns. Future studies can apply global learning to various hydrometeorological time series and climates. As larger training datasets often result in more accurate DL models, we hypothesize that including more time series from different hydrometeorological variables leads to better performance and generalization.



(a)



(b)



(c)

Figure 7. Accuracy in global vs. local learning in forecasting reference evapotranspiration in Durham CIMIS station. Best statistical and deep learning local learning models (a), N-BEATS global learning with Optuna (b), N-HiTS global learning with Optuna (c)

423

424 5. Summary and conclusion

425 This study explored how global learning affects the forecasting accuracy and generalizability of
 426 deep learning models. Although we focused on reference evapotranspiration forecasting in the
 427 Central Valley of California, we hypothesize that our findings are applicable to other seasonality-
 428 driven hydrometeorological time series and different climates and regions. Our findings revealed
 429 that global learning significantly enhances the accuracy of deep learning models. Global learning
 430 could unleash the potential of high-capacity deep architectures by providing them with more
 431 training data. We also argued how using global learning can advance the generalizability of deep
 432 forecasting models. Deep global learning can leverage large data availabilities in gauged regions

and facilitate reliable hydrometeorological time series forecasting in ungauged areas and recently established stations with short data histories. We hypothesize that employing an array of hydrometeorological variables in the training set can further improve the accuracy and generalizability of deep forecasting models in future studies.

Acknowledgements

This research was funded by the California Department of Water Resources and the University of California, Davis grant number 4600014165-01.

References

- Ahmadi, A., Daccache, A., Sadegh, M. and Snyder, R.L., 2023. Statistical and deep learning models for reference evapotranspiration time series forecasting: A comparison of accuracy, complexity, and data efficiency. *Computers and Electronics in Agriculture*, 215, p.108424.
- Ahmadi, A., Daccache, A., Snyder, R.L. and Suvočarev, K., 2022. Meteorological driving forces of reference evapotranspiration and their trends in California. *Science of the Total Environment*, 849, p.157823.
- Ahmadi, A., Kazemi, M.H., Daccache, A. and Snyder, R.L., 2024. SolarET: A generalizable machine learning approach to estimate reference evapotranspiration from solar radiation. *Agricultural Water Management*, 295, p.108779.

450 Akiba, T., Sano, S., Yanase, T., Ohta, T. and Koyama, M., 2019, July. Optuna: A next-
 451 generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD*
 452 *international conference on knowledge discovery & data mining* (pp. 2623-2631).

453 Al Daoud, E., 2019. Comparison between XGBoost, LightGBM and CatBoost using a home
 454 credit dataset. *International Journal of Computer and Information Engineering*, 13(1), pp.6-10.

455 Allen, R.G., Pereira, L.S., Raes, D. and Smith, M., 1998. Crop evapotranspiration-Guidelines for
 456 computing crop water requirements-FAO Irrigation and drainage paper 56. *Fao, Rome*, 300(9),
 457 p.D05109.

458 Allen, R.G., Pruitt, W.O., Wright, J.L., Howell, T.A., Ventura, F., Snyder, R.L., Itenfisu, D.,
 459 Steduto, P., Berengena, J., Baselga Yrisarry, J., Smith, M., Pereira, L.S., Raes, D., Perrier, A.,
 460 Alves, I., Walter, I. and Elliott, R. 2006. A recommendation on standardized surface resistance
 461 for hourly calculation of reference ETo by the FAO56 Penman-Monteith method. *Agricultural*
 462 *Water Manual*, 81: 1-22.

463 Allen, R.G., Walter, I.A., Elliott, R.L., Howell, T.A., Itenfisu, D., Jensen, M.E. and Snyder, R.L.
 464 2005. The ASCE Standardized Reference Evapotranspiration Equation. *Amer. Soc. of Civil Eng.*
 465 Reston, Virginia, 192p.

466 Assimakopoulos, V. and Nikolopoulos, K., 2000. The theta model: a decomposition approach to
 467 forecasting. *International journal of forecasting*, 16(4), pp.521-530.

468 Bai, S., Kolter, J.Z. and Koltun, V., 2018. An empirical evaluation of generic convolutional and
 469 recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*.

470 Bandara, K., Bergmeir, C. and Smyl, S., 2020. Forecasting across time series databases using
471 recurrent neural networks on groups of similar series: A clustering approach. *Expert systems with*
472 *applications*, 140, p.112896.

473 Challu, C., Olivares, K.G., Oreshkin, B.N., Ramirez, F.G., Canseco, M.M. and Dubrawski, A.,
474 2023, June. Nhits: Neural hierarchical interpolation for time series forecasting. In *Proceedings of*
475 *the AAAI conference on artificial intelligence* (Vol. 37, No. 6, pp. 6989-6997).

476 Chen, S., Feng, Y., Mao, Q., Li, H., Zhao, Y., Liu, J., Wang, H. and Ma, D., 2024. Improving the
477 accuracy of flood forecasting for Northeast China by the correction of global forecast rainfall
478 based on deep learning. *Journal of Hydrology*, p.131733.

479 Chen, Z., Zhu, Z., Jiang, H. and Sun, S., 2020. Estimating daily reference evapotranspiration
480 based on limited meteorological data using deep learning and classical machine learning
481 methods. *Journal of Hydrology*, 591, p.125286.

482 Chia, M.Y., Huang, Y.F., Koo, C.H., Ng, J.L., Ahmed, A.N. and El-Shafie, A., 2022. Long-term
483 forecasting of monthly mean reference evapotranspiration using deep neural network: A
484 comparison of training strategies and approaches. *Applied Soft Computing*, 126, p.109221.

485 Contreras, J., Espinola, R., Nogales, F.J. and Conejo, A.J., 2003. ARIMA models to predict next-
486 day electricity prices. *IEEE transactions on power systems*, 18(3), pp.1014-1020.

487 Das, A., Kong, W., Leach, A., Mathur, S., Sen, R. and Yu, R., 2023. Long-term forecasting with
488 tide: Time-series dense encoder. *arXiv preprint arXiv:2304.08424*.

489 Dong, J., Zhu, Y., Jia, X., Han, X., Qiao, J., Bai, C. and Tang, X., 2022. Nation-scale reference
 490 evapotranspiration estimation by using deep learning and classical machine learning models in
 491 China. *Journal of Hydrology*, 604, p.127207.

492 Ferreira, L.B. and da Cunha, F.F., 2020. Multi-step ahead forecasting of daily reference
 493 evapotranspiration using deep learning. *Computers and electronics in agriculture*, 178,
 494 p.105728.

495 Gleick, P.H., 2003. Water use. *Annual review of environment and resources*, 28(1), pp.275-314.

496 Goodfellow, I., Bengio, Y. and Courville, A., 2016. *Deep learning*. MIT press.

497 Gocić, M., Motamedi, S., Shamshirband, S., Petković, D., Ch, S., Hashim, R. and Arif, M., 2015.
 498 Soft computing approaches for forecasting reference evapotranspiration. *Computers and*
 499 *Electronics in Agriculture*, 113, pp.164-173.

500 Hanak, E., 2011. *Managing California's water: From conflict to reconciliation*. Public Policy
 501 Instit. of CA.

502 Herzen, J., Lässig, F., Piazzetta, S.G., Neuer, T., Tafti, L., Raille, G., Van Pottelbergh, T.,
 503 Pasička, M., Skrodzki, A., Huguenin, N. and Dumonal, M., 2022. Darts: User-friendly modern
 504 machine learning for time series. *Journal of Machine Learning Research*, 23(124), pp.1-6.

505 Holt, C.C., 2004. Forecasting seasonals and trends by exponentially weighted moving
 506 averages. *International journal of forecasting*, 20(1), pp.5-10.

507 Hochreiter, S. and Schmidhuber, J., 1997. Long short-term memory. *Neural computation*, 9(8),
 508 pp.1735-1780.

509 Hyndman, R.J. and Athanasopoulos, G., 2018. *Forecasting: principles and practice*. OTexts.

510 Hyndman, R.J. and Khandakar, Y., 2008. Automatic time series forecasting: the forecast package
511 for R. *Journal of statistical software*, 27, pp.1-22.

512 Kalekar, P.S., 2004. Time series forecasting using holt-winters exponential smoothing. *Kanwal*
513 *Rekhi school of information Technology*, 4329008(13), pp.1-13.

514 Karbasi, M., Jamei, M., Ali, M., Malik, A. and Yaseen, Z.M., 2022. Forecasting weekly
515 reference evapotranspiration using Auto Encoder Decoder Bidirectional LSTM model hybridized
516 with a Boruta-CatBoost input optimizer. *Computers and Electronics in Agriculture*, 198,
517 p.107121.

518 Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q. and Liu, T.Y., 2017.
519 Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information*
520 *processing systems*, 30.

521 Le, M.H., Kim, H., Do, H.X., Beling, P.A. and Lakshmi, V., 2024. A framework on utilizing of
522 publicly availability stream gauges datasets and deep learning in estimating monthly basin-scale
523 runoff in ungauged regions. *Advances in Water Resources*, 188, p.104694.

524 Li, W., Duan, Q., Miao, C., Ye, A., Gong, W. and Di, Z., 2017. A review on statistical
525 postprocessing methods for hydrometeorological ensemble forecasting. *Wiley Interdisciplinary*
526 *Reviews: Water*, 4(6), p.e1246.

527 Lim, B., Arık, S.Ö., Loeff, N. and Pfister, T., 2021. Temporal fusion transformers for
528 interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4),
529 pp.1748-1764.

530 Mehdi-zadeh, S., Fathian, F., Safari, M.J.S. and Adamowski, J.F., 2019. Comparative assessment
531 of time series and artificial intelligence models to estimate monthly streamflow: a local and
532 external data analysis approach. *Journal of Hydrology*, 579, p.124225.

533 Milly, P.C., Betancourt, J., Falkenmark, M., Hirsch, R.M., Kundzewicz, Z.W., Lettenmaier, D.P.
534 and Stouffer, R.J., 2008. Stationarity is dead: Whither water management?. *Science*, 319(5863),
535 pp.573-574.

536 Oreshkin, B.N., Car-pov, D., Chapados, N. and Bengio, Y., 2019. N-BEATS: Neural basis
537 expansion analysis for interpretable time series forecasting. *arXiv preprint arXiv:1905.10437*.

538 Quinn, J.D., Reed, P.M., Giuliani, M. and Castelletti, A., 2024. Average domination: A new
539 multi-objective value metric applied to assess the benefits of forecasts in reservoir operations
540 under different flood design levels. *Advances in Water Resources*, 185, p.104638.

541 Salinas, D., Flunkert, V., Gasthaus, J. and Januschowski, T., 2020. DeepAR: Probabilistic
542 forecasting with autoregressive recurrent networks. *International journal of forecasting*, 36(3),
543 pp.1181-1191.

544 Shwartz-Ziv, R. and Armon, A., 2022. Tabular data: Deep learning is not all you
545 need. *Information Fusion*, 81, pp.84-90.

546 Smyl, S., 2020. A hybrid method of exponential smoothing and recurrent neural networks for
547 time series forecasting. *International journal of forecasting*, 36(1), pp.75-85.

548 Spiliotis, E., Assimakopoulos, V. and Makridakis, S., 2020. Generalizing the theta method for
549 automatic forecasting. *European Journal of Operational Research*, 284(2), pp.550-558.

550 Van Houdt, G., Mosquera, C. and Nápoles, G., 2020. A review on the long short-term memory
551 model. *Artificial Intelligence Review*, 53(8), pp.5929-5955.

552 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł. and
553 Polosukhin, I., 2017. Attention is all you need. *Advances in neural information processing*
554 *systems*, 30.

555 Wilks, D.S., 2011. *Statistical methods in the atmospheric sciences*. Academic press.

556 Winters, P.R., 1960. Forecasting sales by exponentially weighted moving averages. *Management*
557 *science*, 6(3), pp.324-342.

558 Zeng, A., Chen, M., Zhang, L. and Xu, Q., 2023, June. Are transformers effective for time series
559 forecasting?. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 37, No. 9, pp.
560 11121-11128).

561 Zhang, Y., Zhao, Z. and Zheng, J., 2020. CatBoost: A new approach for estimating daily
562 reference crop evapotranspiration in arid and semi-arid regions of Northern China. *Journal of*
563 *Hydrology*, 588, p.125087.

Declaration of interests

☒The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Reference evapotranspiration (ET_0) indicates the atmospheric water demand and is decisive in regional to global water cycles. Like other hydrometeorological time series at monthly scales, monthly ET_0 time series are primarily driven by seasonality. A reliable forecast of these time series is crucial for sustainable water resources planning and management. Although the current and previous research on hydrometeorological time series forecasting focuses on local learning (i.e., training a forecasting model on a single time series and using the trained model for future time steps), our study points to the superiority of the global learning scheme. In global learning, the forecasting model is trained over a pool of multiple time series and tested on new instances. To quantify how deep learning (DL) models can benefit from global learning in hydrometeorological forecasting, our study uses monthly ET_0 time series from 55 standardized weather stations in the Central Valley of California. This study 1) compares the performance of statistical and deep forecasting models in local learning, 2) quantifies the performance improvement of DL models in global versus local learning, and 3) automatically optimizes hyperparameters of the best-performing DL models to achieve state-of-the-art forecasting accuracy. Our findings reveal that while statistical models such as Holt-Winters outperform DL models in local learning, global learning can unleash the true potential of high-capacity DL models such as N-BEATS and N-HiTS. This approach results in RMSE values below 10 mm/month for one-year-ahead forecasts on unseen stations. In addition to superior accuracies, global learning enhances the generalizability of DL models, making them applicable to ungauged locations and recently established weather stations. Our findings also point to the benefits of automatic hyperparameter optimization in deep global forecasting.

- Global learning enhances the accuracy of deep learning models for ET_O forecasting.
- High-capacity DL models outperform traditional methods with global learning.
- Achieves RMSE below 10 mm/month for one-year-ahead ET_O forecasts on unseen stations.
- Models are generalizable to ungauged regions and new weather stations.
- Automatic hyperparameter tuning boosts the performance of deep global forecasting.

August 22, 2024

To

Editorial Board, Journal of Hydrology

Dear Editor-in-Chief,

I am pleased to submit our manuscript, "*Enhancing the Accuracy and Generalizability of Reference Evapotranspiration Forecasting in California Using Deep Global Learning*," for consideration in the Journal of Hydrology. Our research addresses the critical question of how global learning can unleash the potential of deep learning models for hydrometeorological time series forecasting, an area of growing importance in hydrology and water resources.

Our study harnesses the power of state-of-the-art deep learning architectures and a novel forecasting strategy. We tested the application of deep global learning in the Central Valley of California, one of the world's most hydrologically altered and agriculturally productive regions.

Key findings of our research include:

- Global learning enhances the forecasting accuracy of deep learning models.

While locally learned deep forecasting models fall behind statistical models in terms of accuracy, global learning unleashes the true power of these advanced architectures by introducing them to vast data.

- Global learning elevates the generalizability of deep learning models.

Globally learned models take advantage of data availability in gauged regions and can be applied to newly established stations and data-scarce regions.

- Automatic hyperparameter tuning further advances the accuracy of deep learning models.

Our results underscore the importance of deep global learning in hydrological time series forecasting. This manuscript not only contributes new insights into the field of hydrology but also offers practical implications for leveraging advances in time series forecasting for more sustainable water resources management.

This original manuscript has not been published or submitted for publication elsewhere. All co-authors have read and approved the final manuscript, and there are no conflicts of interest to disclose.

We believe our findings will interest the readers of the Journal of Hydrology, and we look forward to your feedback.

I appreciate your consideration.

Sincerely,
Arman Ahmadi, Ph.D.
University of California, Berkeley
Email: a.ahmadi@berkeley.edu