

# Improving Reservoir Water Supply and Inflow Volume Predictions with Encoder-Decoder Deep Learning Models

Zeeshan Asghar<sup>1\*</sup>, Muhammad Waseem<sup>1</sup>, Zulqarnain Jehan<sup>1</sup>

<sup>1</sup> Department of Civil Engineering, Ghulam Ishaq Khan Institute of Engineering Sciences and Technology, Topi, District Swabi, Pakistan

\* Correspondence: [zeeshan.asghar@giki.edu.pk](mailto:zeeshan.asghar@giki.edu.pk)

## Abstract

Recently, deep learning (DL) models have shown tremendous potential for hydrological prediction, reservoir management, and operational planning. However, their effectiveness in predicting reservoir inflows over extended time horizons remains limited. Recent advancements in deep learning algorithms have enhanced the accuracy of inflow forecasts, but most studies have focused on short-term applications or real-time operations. This study introduces a novel multi-step forecasting framework aimed at improving long-term predictions of reservoir inflow and water supply. Using the snow water equivalent metric and reservoir inflow data, we trained a deep learning model using a Convolution Neural Network (CNN.) - Long Short-Term Memory (LSTM.) encoder- decoder model to forecast inflows to predict future time series of steps during the critical March-August runoff period. The model's architecture and hyper-parameters were fine-tuned using multi-fold cross-validation of the timeseries, analyzing different adaptations of Encoder-Decoder architectures based on CNNs and LSTMs. The proposed methodology was applied to 40-year time series of SWE and inflow data from the Jordanelle Reservoir in Utah. The optimal configuration a 16-node-per-layer LSTM Encoder-Decoder model demonstrated significant improvements in long-term forecast accuracy. We further examined the balance between model complexity and performance by benchmarking against a process-driven Ensemble Streamflow Forecasting model and traditional statistical methods, including SARIMA, VAR, and TBATS. The deep learning approach outperformed statistical models in long-term water supply forecasts and achieved comparable accuracy to the ESP model's 50% exceedance probability forecast. These findings indicate the potential of enhanced DL methods to improve the long-term hydrological forecasting and resource management.

**Keywords:** Reservoir Water Supply; Inflow Volume Predictions; Encoder-Decoder Deep Learning Models; CNN, LSTM

## 1. Introduction

Conventional statistical techniques, including the Vector Auto-Regression (VAR), Seasonal Auto-Regressive Integrated Moving Average (SARIMA), Trigonometric Box-Cox Transformation approach (1), Vector Auto-Regression (VAR) (2), and seasonal components (TBATS) (?), have long been utilized for predicting reservoir inflow and outflow patterns. These methods leverage historical data and inherent seasonal trends to provide insights into water resource dynamics. While effective for short-term forecasts, these models struggle with long-term predictions due to the convergence of auto-regressive components to the time-series mean (4). Additionally, long-term hydrological forecasts are hindered by the complexities of hydro-meteorological variability,

particularly in snow-dominated catchments (5). Non-linear weather-driven patterns further challenge statistical models that assume linear relationships.

Physical hydrological models simulate processes using probabilistic ensembles driven by atmospheric inputs, providing forecasts with quantified uncertainties through statistical post-processing (6; 7). The Ensemble Streamflow Prediction (ESP) model (49) is widely used for long-term water supply forecasting. This framework initializes hydrological states using current basin conditions to produce probabilistic inflow predictions. However, ESP does not directly integrate the Snow Water Equivalent (SWE), relying instead on the simulated snowpack dynamics derived from the temperature and precipitation inputs (9). This can result in reduced accuracy, particularly when snowpack variability is significant (10).

The limitations of both statistical and physical models highlight the need for more advanced techniques capable of capturing non-linear patterns and addressing hydrological variability. Machine learning (ML) approaches offer a data-driven alternative, demonstrating success in fields such as forecasting rainfall runoff (11), hydropower forecasting (12), and spatial estimation of SWE (13). These methods excel at identifying complex relationships without relying on explicit process-based assumptions.

Deep learning has further enhanced the capacity for reservoir inflow forecasting. Models like Encoder-Decoder frameworks effectively capture temporal dependencies and nonlinear interactions, making them particularly suitable for multistep forecasting (14). Despite their promise, most applications have focused on short-term scenarios with limited exploration of long-term inflow predictions. To address this gap, we propose a multistep forecasting approach leveraging historical SWE and reservoir inflow data to train Encoder-Decoder algorithms for long-term predictions.

This study examines four distinct Encoder-Decoder architectures to evaluate their performance in forecasting: CNN-LSTM, residual CNN-LSTM, standard LSTM-LSTM, and residual LSTM-LSTM models. Each architecture is designed to leverage the strengths of both convolutional and recurrent layers, allowing for efficient extraction of spatial and temporal features. By incorporating residual connections in some variants, the models aim to enhance gradient flow, reduce vanishing gradients, and improve training efficiency for complex sequences. Although LSTMs capture sequential dependencies, CNNs excel in parallel processing of fixed-size contexts (15). Residual connections improve CNN performance by enabling deeper architectures and improving gradient flow (16). These models are compared against the ESP framework and traditional statistical methods (SARIMA, VAR, TBATS) to assess their effectiveness in long-term forecasting.

The key research questions addressed in this study are:

1. What is the optimal balance between the complexity and accuracy of the model for the long-term forecasting of the water supply?
2. Under which conditions does the proposed deep learning model outperform benchmark models in accuracy?
3. How do purely data-driven approaches compare to process-based physical models for long-term reservoir inflow predictions?

## 2. Methodology

### 2.1. Site

This study focuses on the Jordanelle Reservoir, a key component of the Provo River Project and an essential water storage facility for central Utah. The reservoir is situated in Wasatch County, Utah, along the Provo River, upstream of Heber City. Its watershed encompasses a drainage area of approximately 234 square miles, monitored at the Provo River near Hailstone gauging station (10155000). The region receives an average annual precipitation of 25.8 inches, with runoff primarily derived from spring snowmelt, contributing significantly to its inflow from April through June.

Completed in 1993 by the Bureau of Reclamation, Jordanelle Reservoir has a total storage capacity of 320,300 acre-feet, covering a surface area of up to 3,068 acres at full pool. The reservoir is a vital part of Utah's water management system, providing water for municipal, industrial, and agricultural use while supporting recreational activities. Water releases are regulated through the Jordanelle Dam, which incorporates outlet works capable of discharging up to 8,000 cubic feet per second (cfs) (17). Additionally, the reservoir serves as a crucial buffer for managing downstream flow into the Deer Creek Reservoir and mitigating flood risks in the Provo River basin (18).

Managed by the Central Utah Water Conservancy District (CUWCD), Jordanelle Reservoir forms part of a broader water infrastructure network designed to meet the growing demands of Utah's population. The reservoir plays an integral role in water storage and delivery under the Central Utah Project, ensuring the sustainable allocation of resources across the Wasatch Front. Figure 1 highlights the geographic location of Jordanelle Reservoir and its proximity to snow telemetry (SNOTEL) sites, crucial for tracking snowpack dynamics and forecasting runoff for operational planning.

## 2.2. Datasets

The dataset utilized for training the model originated from the operational archives of the Central Utah Water Conservancy District (CUWCD) and the snow telemetry (SNOTEL) network, which is managed by the National Resource Conservation Service (NRCS) (22; 23). Reservoir inflow data were computed using changes in reservoir water levels, adjusted to account for losses such as seepage and evaporation, and incorporating releases into Rock Creek. This methodology facilitated precise calculations of storage variations in acre-feet (24).

Two key datasets formed the foundation of this study: reservoir inflow and snow-water equivalent (SWE). SWE represents the amount of water (in inches) that would result from melting the snowpack. Historical daily records for Rock Creek at the Jordanelle Reservoir span from January 1990 to the present, with regular updates. The CUWCD ensures high-quality, gap-free reservoir inflow data, eliminating the need for interpolation (Figure 2).

Reservoir inflow patterns exhibit a seasonal runoff trend, predominantly driven by snowmelt, which typically peaks between April and August. This study assumes that critical water storage decisions must be finalized by late March to prepare for the snowmelt season. Therefore, the forecasting period from March to August was emphasized, necessitating a model that can accurately represent long-term temporal trends (25).

The SWE data were obtained from the NRCS SNOTEL network, covering the same timeframe as the reservoir inflow data (Table 1). The monitoring stations provided high-quality data, with fewer than three days of missing values per site. Where gaps existed, they were addressed using interpolation techniques. The primary objective of the analysis was to model the relationship between SWE (Figure 3) and reservoir inflow (Figure 2).

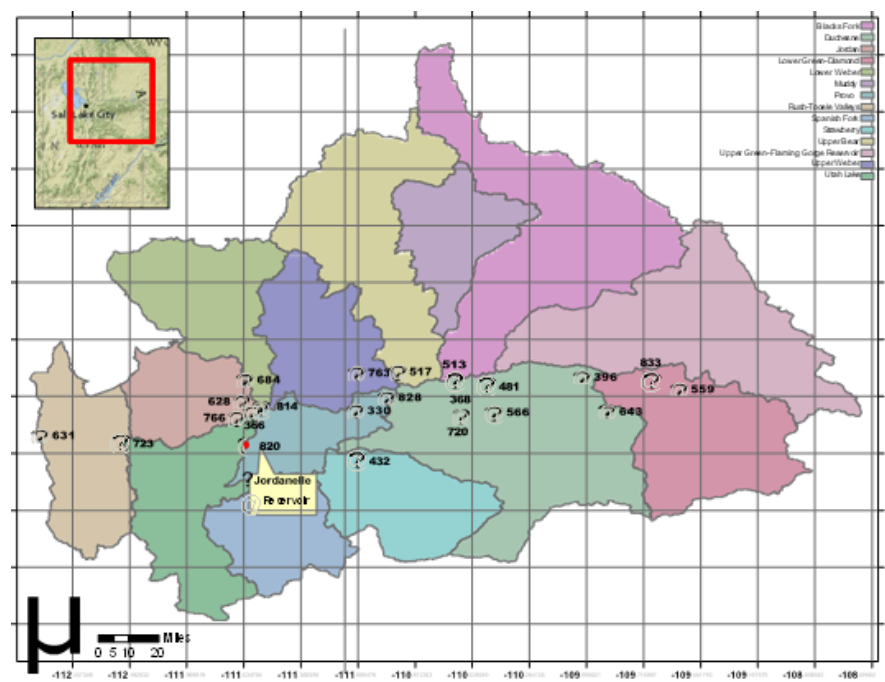


Figure 1. Study site location: Jordanelle Reservoir, Utah

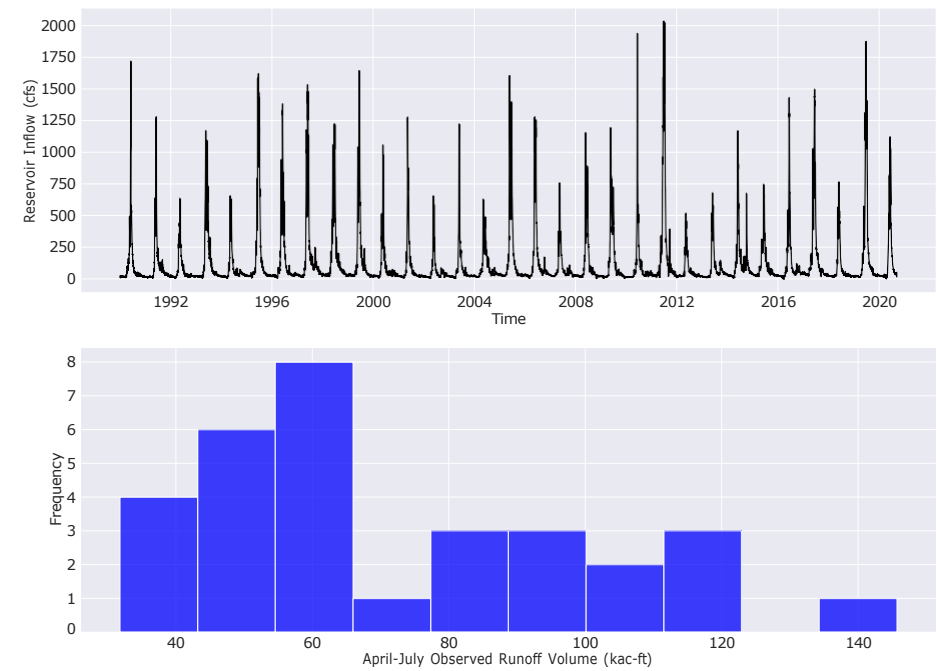
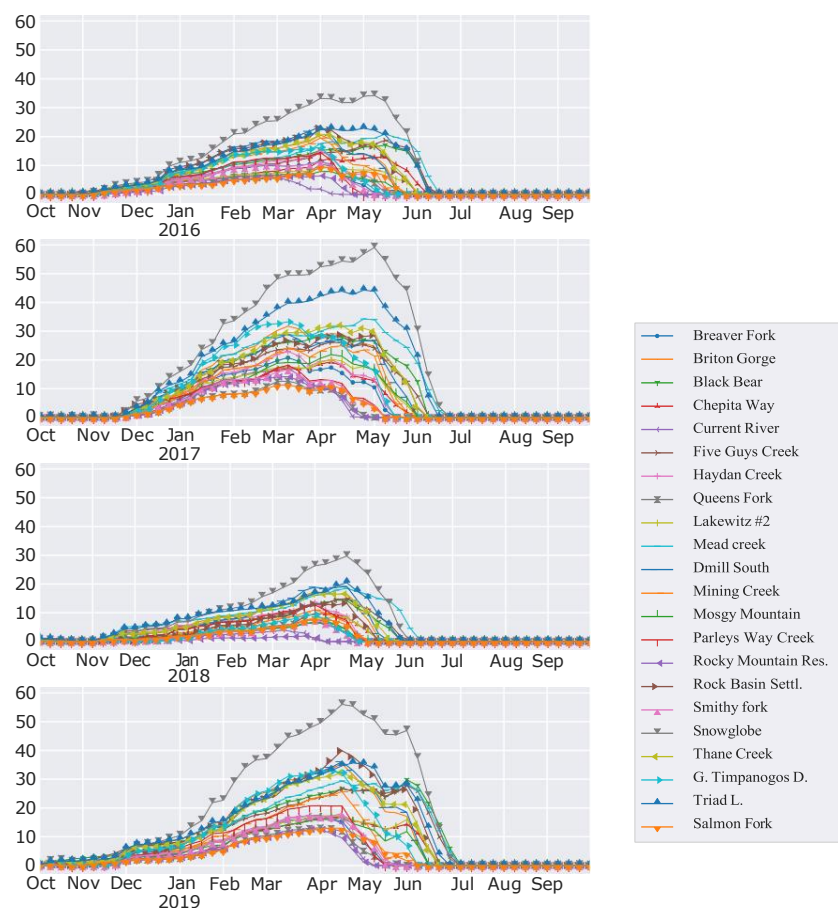


Figure 2. Historical daily inflow patterns (top) and March-August inflow volumes (bottom) at Upper-Stilwater.

For model training, the daily data were aggregated into weekly averages and scaled to a normalized range of 0 to 1 to align with the deep learning model’s activation function (Section 2.3). A sliding window approach was adopted, utilizing a 20-week input period to forecast the subsequent 26 weeks, corresponding to the runoff season (March to August). The 20-week window was selected based on typical SWE accumulation durations, which vary from 15 to 25 weeks depending on winter precipitation patterns. This configuration reflects the most representative conditions across multiple years, providing an optimal basis for training the model.

**Table 1.** Metadata for Snow Telemetry Stations (NRCS)

St- Id	Location	Lon	Lat
340:UTa:SNoTL	Breavers Devide	-109.987	39.206
346:UTa:SNoTL	Laughton	-110.466	39.193
358:UTa:SNoTL	Black Bear	-109.480	39.175
356:UTa:SNoTL	Chepeta Way	-108.910	39.367
462:UTa:SNoTL	Current River	-109.979	39.953
461:UTa:SNoTL	Five Guys Creek	-109.362	39.311
577:UTa:SNoTL	Hayden Fork	-109.776	39.389
579:UTa:SNoTL	Kings Cabin	-108.449	39.309
586:UTa:SNoTL	Riverfork #1	-109.329	39.191
583:UTa:SNoTL	Riverfork Basin	-109.515	39.331
698:UTa:SNoTL	Miller-D North	-110.520	39.252
661:UTa:SNoTL	Dining Forks	-111.485	39.089
633:UTa:SNoTL	Mosby Mountain	-108.789	39.202
624:UTa:SNoTL	Parleys Way	-110.513	39.354
750:UTa:SNoTL	Red Rock Canyon	-109.586	39.144
783:UTa:SoNTL	R.B Settlement	-111.102	39.999
713:UTa:SNoTL	S. Moreshouse	-109.981	39.381
794:UTa:SNoTL	Thanos Canyon	-110.418	39.218
820:UTa:SNoTL	Gt. T. Divide	-110.500	39.024
828:UTa:SNoTL	Trivial Lake	-109.840	39.271
833:UTa:SNoTL	Salmon River	-108.576	39.332



**Figure 3.** Snow water equivalent (SWE) measurements (in inches) over five continuous water years.

### 2.3. Encoder-Decoder DL Model

The primary aim of the model is to forecast future reservoir inflow sequences based on historical inputs, incorporating time-series data for reservoir inflow and SWE. This is achieved through a multivariate sequence-to-sequence prediction framework, consisting of two main components: an encoder that converts the input sequence into a fixed-length representation, and a decoder that reconstructs this representation into the predicted output sequence (34). The decoder is further supplemented by a fully connected time-distributed layer, which refines the predictions into the final output sequence.

As depicted in Figure 4, the model incorporates four different variants, each utilizing a sliding window approach where multiple time-series variables are processed to produce corresponding output windows.

This architecture has proven effective for a variety of sequence-to-sequence applications, including flood prediction (58), traffic flow estimation (27), weather forecasting (28), and solar performance modeling (29). The implementation of each model variant was carried out using Python and the Keras library (30). The exponential linear unit (ELU) activation function (31) was employed to enhance learning performance, while the Adam optimizer (32) was utilized for efficient weight optimization across the network.

**Encoder-Decoder Variants:** Recurrent neural networks (RNNs) are well-suited for sequential data processing but often face challenges when input-output relationships span extended time gaps, leading to difficulties in handling long-term dependencies. Long short-term memory (LSTM) networks address this limitation through a cell state and gated mechanisms that efficiently manage information retention and removal across time steps (33).

The Encoder-Decoder framework offers significant advantages for sequence-to-sequence tasks by encoding variable-length sequences into fixed-length vectors, ensuring adaptability across diverse datasets (34; 35). Using LSTMs within this framework enhances its ability to capture long-term dependencies, overcoming the limitations of traditional RNNs (33). The framework's modularity supports a wide range of architectures, including CNN-based encoders for feature extraction and LSTM decoders for temporal sequence modeling (15). Residual connections improve gradient flow and facilitate efficient training of deeper networks (37), while causal padding in CNN encoders ensures temporal integrity by preventing look-ahead bias (36).

Figure 4 illustrates four Encoder-Decoder model variants, each incorporating varying degrees of architectural complexity:

1. **LSTM Encoder-LSTM Decoder:** This variant integrates an LSTM-based encoder and decoder, linked through a repeat vector to align the encoded features with the decoder's input. The decoder reconstructs the output sequence, which is refined by a time-distributed Dense layer. This layer applies shared weights across time steps, ensuring consistent output generation.
2. **Residual LSTM Encoder-LSTM Decoder:** An extension of the LSTM-LSTM model, this variant incorporates residual connections between layers to enhance gradient flow and error propagation. These connections support the network in maintaining greater depth without compromising optimization efficiency (37).
3. **CNN Encoder-LSTM Decoder:** This model replaces the LSTM encoder with a one-dimensional convolutional neural network (CNN) that processes time-series data. CNN captures local patterns through hierarchical representations, while causal padding preserves temporal order, preventing any look-ahead bias (36).
4. **Residual CNN Encoder-LSTM Decoder:** This architecture builds on the CNN-LSTM variant by incorporating residual connections within the CNN-based

encoder. These connections improve feature extraction and gradient flow, enabling the model to learn hierarchical abstractions effectively. Pooling layers at the end of each residual block condense the input features into meaningful representations.

In the CNN-based models, the encoder progressively reduces the input matrix's dimensions while abstracting relevant features. Convolutional filters and pooling layers create a condensed feature map, with max pooling applied to retain key elements. The resulting features are flattened and passed to the LSTM decoder, which processes the sequence to generate the output.

This combination of LSTM and CNN architecture leverages their respective strengths, enabling the model to identify both localized patterns and long-term dependencies within the input data, ensuring a robust and flexible predictive framework.

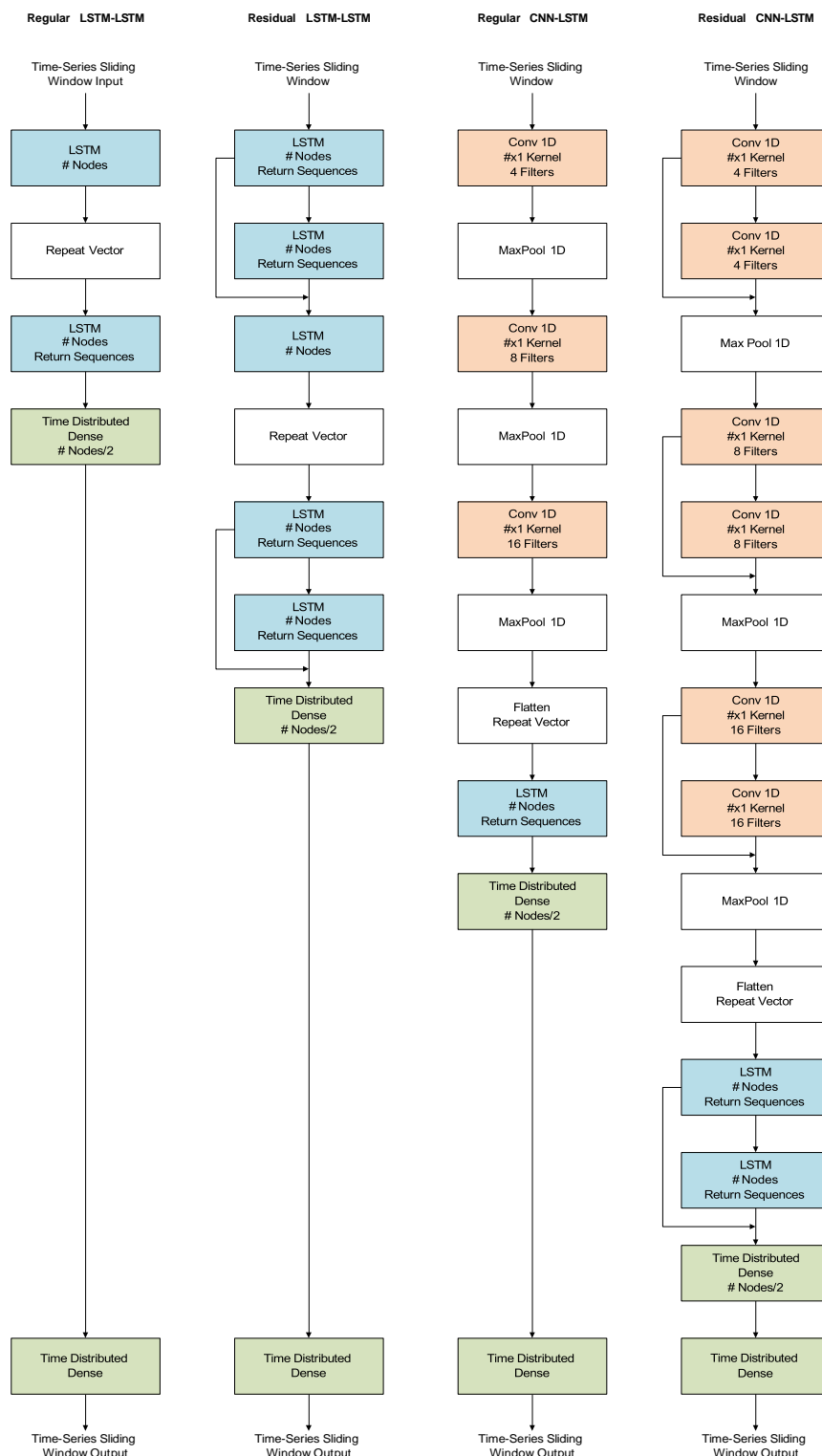
**Model Selection** The historical dataset, spanning forty years for inflow and snow water equivalent, is expanded to enhance the model's ability to generalize to unseen scenarios. To achieve this, the index sequential method (ISM) is utilized, a technique commonly applied in hydrological modeling, particularly in the Colorado River Basin (38; 39; 40). ISM generates synthetic hydrological sequences by incrementally shifting the historical record by one water year, broadening the range of potential outcomes and accounting for uncertainties in future hydrological behavior due to natural variability and human-induced climate changes (41).

To address ISM's limitations in capturing extreme events, such as extended droughts, this study implements a modified approach. Water years are treated as discrete blocks and randomly shuffled to introduce additional variability. This method, referred to as water year block disaggregation, diversifies the training dataset, improving the model's robustness (41; 43; 44). Additionally, each water year is randomly scaled within a range of 0.4 to 1.4, preserving seasonal trends while introducing magnitude variability. This adjustment allows the model to train on a wider spectrum of hydrological extremes, expanding the training dataset fivefold before cross-validation.

Hyper-parameter tuning for each model was performed using five-fold time-series cross-validation. Data from 2011–2015 are utilized for hyperparameter optimization, while the evaluation period from 2016–2020 is used to identify the best-performing architecture. Each Encoder-Decoder variant was tested with configurations of 8, 16, and 32 LSTM nodes, and CNN kernel sizes of 2, 4, and 8, resulting in a total of 12 configurations. Residual connections are incorporated to enhance gradient flow, and a Dense layer with time-distributed functionality was utilized to process the decoder's output efficiently.

The training process connects layers of nodes, where each node processes input data and forwards it using an activation function. In CNN-based models, kernels slide over input sequences to extract features, while filters consolidate these features into maps. The input data comprise SWE and reservoir inflow time-series from November to March, while the target output predicts runoff for March to August.

To prevent overfitting during training, an early stopping mechanism is implemented, which halts the training process if no reduction in mean squared error (MSE) is observed over 10 consecutive epochs. This approach helps avoid overfitting by preventing the model from excessively adapting to the training data. Additionally, the total number of epochs is capped at 50, ensuring a controlled training duration. The batch size is carefully selected to determine the fraction of data processed in each training iteration, balancing computational efficiency and model convergence.



**Figure 4.** Architecture of Encoder-Decoder variants. Colored cells represent layers with trainable parameters, while non-colored cells indicate non-trainable layers, such as repeat vector and max pooling operations.

The evaluation of model performance is carried out using a comprehensive set of metrics to ensure a robust assessment. These metrics include normalized root mean squared error, absolute error, Nash-Sutcliffe efficiency, median absolute error, and explained variance. Each metric provides unique insights into different aspects of the

model's predictive capabilities, enabling a thorough evaluation of accuracy, consistency, and reliability. This multifaceted evaluation framework ensures that the model's performance is not only accurate but also generally unseen data.

The dataset spans 46 weeks, divided into a 20-week input window (November to March) and an 26-week output window (March to August). Model performance is ranked based on the average metric scores over the five-year evaluation period, ensuring robust and reliable predictions.

$$MAE = \frac{\sum_{j=1}^M |(Q_{obs,j} - Q_{pred,j})|}{\sum_{j=1}^M (Q_{obs,j})} \quad (1)$$

$$RMSE = \frac{\sqrt{\frac{1}{M} \sum_{j=1}^M (Q_{obs,j} - Q_{pred,j})^2}}{\frac{1}{M} \sum_{j=1}^M (Q_{obs,j})} \quad (2)$$

$$MedAE = \frac{\text{Med}[|(Q_{obs} - Q_{pred})|]}{\frac{1}{M} \sum_{j=1}^M (Q_{obs,j})} \quad (3)$$

$$N.SE = 1 - \frac{\sum_{j=1}^M (Q_{obs,j} - Q_{pred,j})^2}{\sum_{j=1}^M (Q_{obs,j} - \frac{1}{M} \sum_{k=1}^M (Q_{obs,k}))^2} \quad (4)$$

$$ExpVar = 1 - \frac{\sum_{j=1}^M ((Q_{obs,j} - Q_{pred,j}) - \frac{1}{M} \sum_{k=1}^M (Q_{obs,k}))^2}{\sum_{j=1}^M (Q_{obs,j} - \frac{1}{M} \sum_{k=1}^M (Q_{obs,k}))^2} \quad (5)$$

Quantification of Uncertainty This study investigates uncertainty in two key dimensions: the predictive accuracy of the model and the intrinsic variability in reservoir inflow patterns, particularly during weekly averages and the critical March-to-August runoff period. Predictive uncertainty is evaluated through a 95% confidence interval derived from an ensemble of model simulations. Given the stochastic nature of deep learning models, each training iteration introduces slight variations in the forecast outcomes. To capture this variability, the model is trained multiple times, generating a distribution of predictions for each forecast time step. This distribution is assumed to follow a normal distribution, consistent with the central limit theorem as the number of ensemble members increases.

Outliers within the forecasted results are identified using the Tukey method, which flags data points lying beyond the whiskers in a boxplot representation (45). This approach ensures that anomalous predictions, which could adversely impact model interpretation, are systematically recognized and addressed.

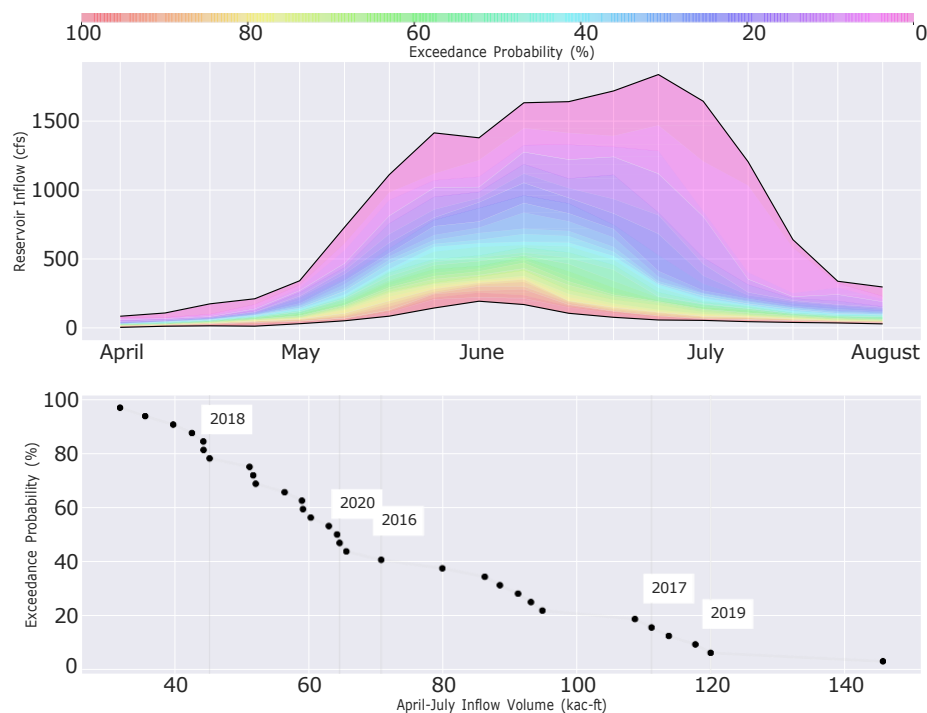
Additionally, the study emphasizes the importance of understanding inflow variability, which is critical for effective water resource planning and management. The variability in inflow patterns is influenced by complex hydro-meteorological factors, including snowmelt timing, precipitation intensity, and temperature fluctuations. To provide a clear visualization of this variability, exceedance probability plots are employed. These plots illustrate the likelihood of inflows exceeding specific thresholds over the forecast horizon, offering valuable insights into potential risks and resource allocation needs during peak and low inflow periods.

By combining predictive uncertainty analysis with a focus on inflow variability, the study provides a comprehensive framework to evaluate model reliability and inform decision-making in the context of reservoir management. This dual approach not only

enhances the robustness of the forecasting methodology but also supports adaptive strategies to mitigate the impact of uncertain and dynamic hydrological conditions.

The second component of uncertainty concerns the variability in reservoir inflow, which serves as the primary forecast target. This variability is visualized using an exceedance probability plot for the multi-step forecast period (Figure 5).

The exceedance probability plot is constructed by calculating the likelihood of daily inflow values exceeding specific thresholds based on 40 years of historical data. These daily probabilities are aggregated into weekly averages for the March-to-August runoff period, which is critical for reservoir operations. The plot features 30 individual traces, each representing ranked reservoir inflow values sorted by exceedance probabilities. This visualization captures the variability and potential range of inflow scenarios across different years. Additionally, the total inflow volume for the March-to-August period is computed and displayed, enabling a comparative analysis with historical trends. The hold-out years (2016–2020) are annotated on the plot to illustrate their unique inflow characteristics in relation to the historical record, providing insights into anomalous or representative hydrological patterns.



**Figure 5.** Exceedance probability of inflow evaluated for weekly averages and total inflow during the March-August period.

**Total Inflow:** The total inflow volume for the four-month runoff period is determined by integrating the forecasted hydrograph. To enhance temporal resolution, the hydrograph, originally based on weekly averaged inflow data, is resampled into daily time steps. This process assumes that the inflow within each week remains constant, allowing the weekly average to be uniformly distributed across the days of that week. The daily inflow values are then summed and converted from cubic feet per second (cfs) to acre-feet per day, providing a more granular estimate of total inflow volume.

To evaluate the accuracy of the total inflow forecasts, a comparison is made with a benchmark Ensemble Streamflow Prediction (ESP) model. Two error metrics are used for this purpose:

$$MARE = \frac{1}{N} \sum_{j=1}^N \left| \frac{V_{obs,j} - V_{pred,j}}{V_{obs,j} - V_{benchmark,j}} \right| \quad (6)$$

$$RMSRE = \sqrt{\frac{1}{N} \sum_{j=1}^N \left( \frac{V_{obs,j} - V_{pred,j}}{V_{obs,j} - V_{benchmark,j}} \right)^2} \quad (7)$$

Relative error values greater than 1.0 indicate performance worse than the ESP model, while values less than 1.0 indicate better performance. This comparison helps evaluate the trade-off between model complexity and predictive accuracy.

#### 2.4. Comparison with Statistical Techniques

The forecasts generated by deep learning models are compared against three statistical methods: VAR (Vector Auto-Regression), TBATS (Trigonometric Seasonal Box-Cox Transformation with ARMA residuals, trend, and seasonal components), and SARIMA (Seasonal Auto-Regressive Integrated Moving Average). These methods are trained using monthly averaged inflow data due to their limitations in handling extended forecasting horizons. To ensure consistency, total inflow volume predictions from these models are resampled to daily time steps using the methodology outlined in Section 2.3. This alignment allows for a direct comparison with deep learning outputs. While VAR, TBATS, and SARIMA effectively capture seasonality and periodic trends, their reliance on linear assumptions limits their ability to model non-linear and dynamic hydrological behaviors, such as peak inflows during snowmelt or abrupt changes driven by extreme weather events. By evaluating these statistical methods alongside deep learning models, the study provides a detailed comparison of their strengths and limitations, emphasizing the suitability of machine learning techniques for addressing the complexities of long-term hydrological forecasts.

**The TBATS Model:** The TBATS model employs exponential smoothing, Box-Cox transformations, and ARMA residuals to capture complex seasonality (48). Seasonal components are represented with trigonometric functions, providing flexibility for modeling high-frequency periodic patterns. This model is implemented using Python's TBATS library and configured to account for quarterly, biannual, and annual seasonal cycles.

**The SARIMA Model:** The SARIMA model predicts time series  $Z_t$  using a seasonal auto-regressive integrated moving average process (46):

$$\varphi(B)\Phi(B^S)(1-B)^d(1-B^S)^D Z_t = \theta(B)\Theta(B^S) e_t \quad (8)$$

Here,  $t$  denotes discrete time,  $S$  represents the seasonal period, and  $B$  is the backward shift operator. Non-seasonal and seasonal auto-regressive components are represented by  $\varphi$  and  $\Phi$ :

$$\varphi(B) = 1 - \varphi_1 B - \varphi_2 B^2 - \dots - \varphi_p B^p \quad (9)$$

$$\Phi(B^S) = 1 - \Phi_1 B^S - \Phi_2 B^{2S} - \dots - \Phi_P B^{PS} \quad (10)$$

The parameters are optimized using Python's Pmdarima library for an annual seasonal period ( $S = 12$  months).

**Vector Auto Regression Model** The Vector Auto-Regression (VAR) model predicts a vector of variables  $y_t$  using its lagged values (47):

$$y_t = A_1 y_{t-1} + A_2 y_{t-2} + \dots + A_p y_{t-p} + u_t \quad (11)$$

Here,  $A_i$  are parameter matrices for each lag, and  $u_t$  represents residuals. The optimal lag order  $p$  is determined by minimizing the Akaike Information Criterion (AIC):

$$AIC(p) = \ln |\Sigma(p)| + \frac{2K^2 p}{N} \quad (12)$$

### 3. Ensemble Streamflows Prediction (ESP) Model

The Ensemble Streamflow Prediction (ESP) approach, developed by the National Weather Service (NWS) (49), is a Monte Carlo simulation technique for probabilistic streamflow forecasting. ESP combines physical hydrological models with probabilistic representations of future weather conditions, leveraging historical meteorological data to generate forecast scenarios. This method assumes historical weather patterns provide a representative sample of possible future conditions (51).

In ESP, each historical weather year is treated as an independent future scenario with equal probability ( $1/m$ , where  $m$  is the number of historical patterns). Hydrological simulations based on snowpack, precipitation, and temperature conditions generate individual streamflow traces. These traces form an ensemble used to fit a probability density function (p.d.f.), describing the likelihood of specific streamflow magnitudes. For this study, the median streamflow value (50% exceedance probability) is used as a benchmark for evaluating model forecasts.

## 4. Results

### 4.1. Optimization of the Model

Cross-Validation results are summarized below, highlighting the identification of optimal hyper-parameters and the selection of the best-performing models (Table 2).

**Table 2.** Optimization of Hyperparameters using 5-Fold Time-Series Cross-Validation on the 2011–2015 Data

ML Model	Hyper-parameters #	Parameters	Cross Val Ranking
ResLSTM-LSTM	Nodes: 8	3400	3
ResLSTM-LSTM	Nodes: 16	11350	2
ResLSTM-LSTM	Nodes: 32	41200	1
LSTM-LSTM	Nodes: 8	1650	3
LSTM-LSTM	Nodes: 16	4950	1
LSTM-LSTM	Nodes: 32	16250	2
ResCNN-LSTM	Nodes: 16	6750	1
ResCNN-LSTM	Nodes: 16	7950	2
ResCNN-LSTM	Nodes: 16	10300	3
CNN-LSTM	Nodes: 16	3900	1
CNN-LSTM	Nodes: 16	4400	2
CNN-LSTM	Nodes: 16	5500	3

For the standard LSTM-LSTM model, the optimal configuration involved 16 nodes per layer, as deviations in this parameter reduced accuracy. In contrast, the residual LSTM-LSTM model showed improved performance as the number of nodes increased from 8 to 16 and 32.

The standard CNN-LSTM and residual CNN-LSTM models exhibited a notable decline in performance as the kernel size increased from 2 to 4 and 8, while maintaining a fixed number of 16 LSTM nodes. This trend highlights the sensitivity of CNN-based architectures to kernel size, where smaller kernels are better suited for capturing localized patterns in the input data. Conversely, the residual LSTM-LSTM model showed improved performance with an increasing number of nodes, achieving its best results at 32 nodes per layer. The standard LSTM-LSTM model, however, demonstrated optimal performance with 16 nodes per layer, indicating a balance between network complexity and predictive accuracy.

Ultimately, four configurations were identified as optimal: (1) standard LSTM-LSTM with 16 nodes per layer, (2) residual LSTM-LSTM with 32 nodes per layer, (3) standard CNN-LSTM with a kernel size of 2 and 16 LSTM nodes, and (4) residual CNN-LSTM with a kernel size of 2 and 16 LSTM nodes.

Final model selection was conducted using performance metrics evaluated on the hold-out data from 2016 to 2020 (Table 3) to ensure an unbiased assessment. Among the four models, the standard LSTM-LSTM emerged as the most accurate, achieving the lowest errors (NMAE, NRMSE, and NMedAE) and the highest scores for NSE and ExpVar. This architecture demonstrated efficiency, with approximately 1,210 trainable parameters per layer and a total of 4,850 parameters across its four layers (Figure 4). The combination of moderate complexity and high predictive performance highlights its suitability for long-term hydrological forecasting tasks.

**Table 3.** Average Performance Metrics for Selected Deep Learning Models (2016–2020 Hold-Out Data)

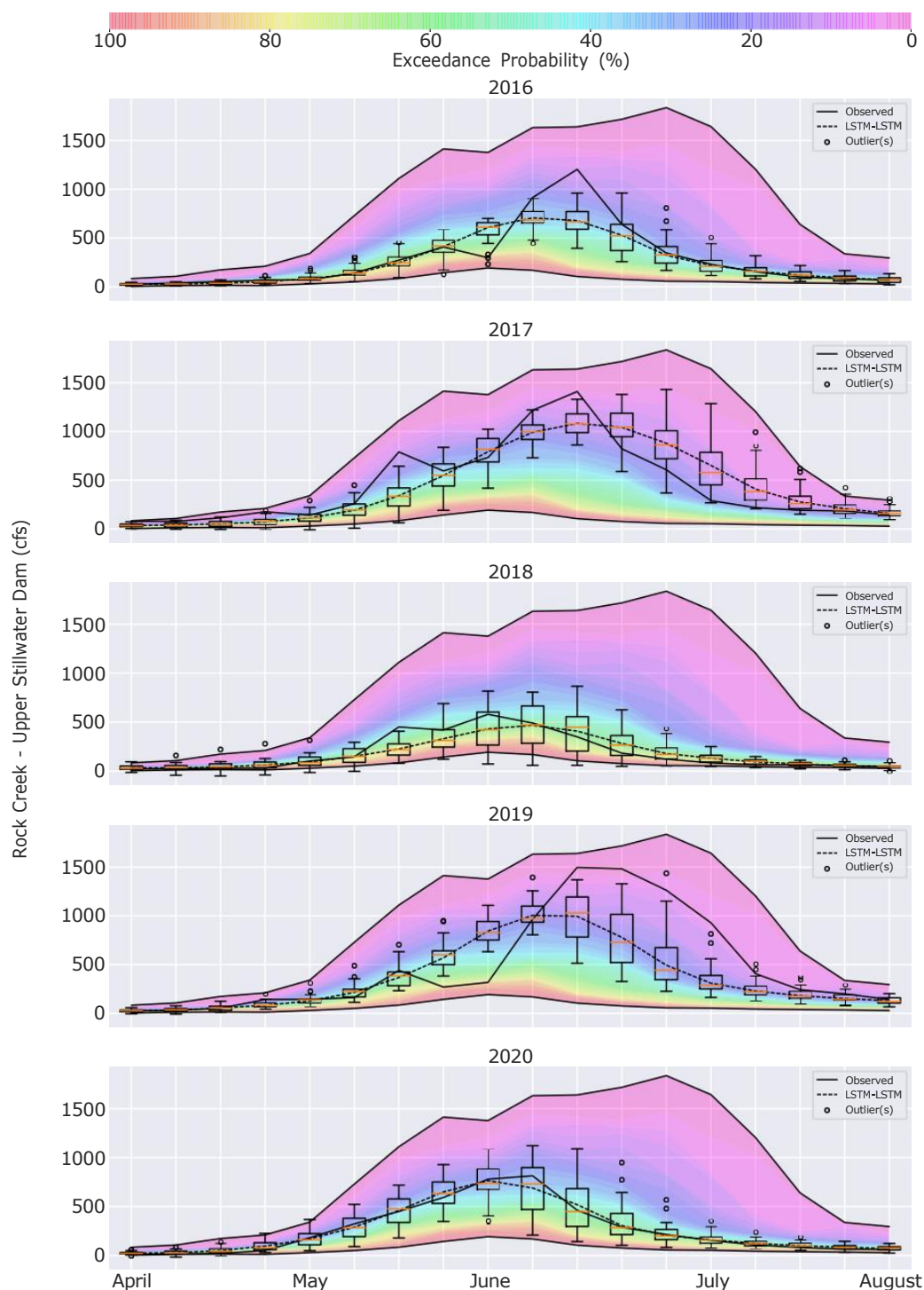
Model	Parameter Density	RMSE	MAE	MedAE	N-SE	ExpVar
ResLSTM-LSTM	5797.44	0.33066	0.50490	0.15939	0.70389	0.73062
LSTM-LSTM	1092.96	0.25047	0.45738	0.16929	0.77715	0.77220
ResCNN-LSTM	646.47	0.47817	0.40095	0.67320	0.30789	0.36630
CNN-LSTM	619.74	0.41877	0.51381	0.25245	0.53163	0.69993

The residual LSTM-LSTM model ranked second, with the highest complexity at 41,050 trainable parameters. Both CNN-LSTM variants demonstrated lower accuracy, with complexities defined by 3,850 and 6,700 trainable parameters for the standard and residual models, respectively. A trend of increasing accuracy with higher model complexity was evident.

#### 4.2. Model Evaluation

The selected models were evaluated using hold-out test years from 2016 to 2020. The multi-step forecasts generated by the selected LSTM-LSTM model are illustrated in Figure 6, which includes observed inflows (solid line), forecasted inflows (dashed line), boxplots, and shaded exceedance probabilities.

The boxplots represent the distribution of inflow predictions across 50 independent model runs for each test year, highlighting the variability and uncertainty in the forecasts. Performance metrics for each year are summarized in Table 4, providing a comprehensive evaluation of the model's accuracy and consistency over the hold-out period.

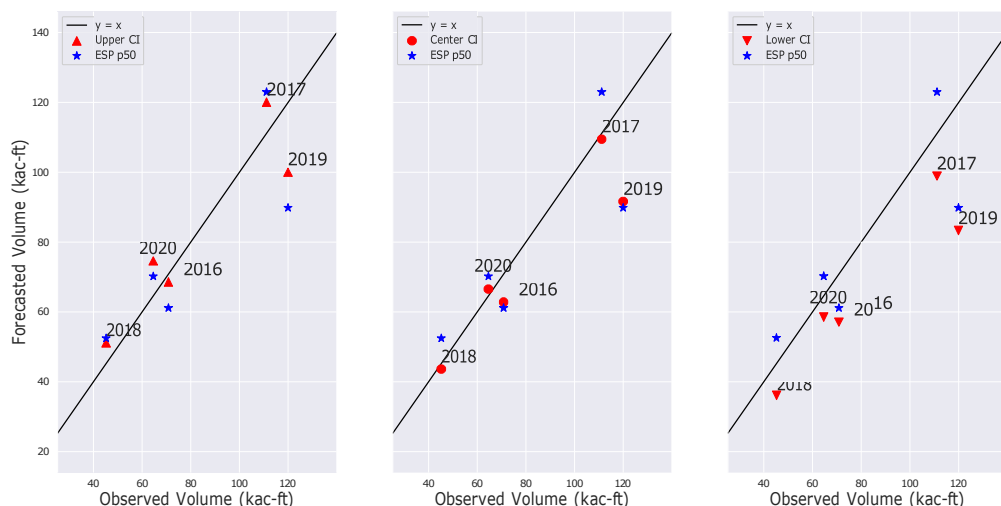


**Figure 6.** Reservoir inflow predictions for the 2016–2020 Held-Out Periods

Furthermore, Figure 7 compares the total forecasted inflow volume with the baseline Ensemble Streamflow Prediction (ESP) model within a 95% confidence interval. This comparison underscores the model's capability to predict overall inflow trends and its reliability in capturing long-term hydrological patterns, validating the effectiveness of the LSTM-LSTM architecture.

**Table 4.** Yearly Performance Metrics of the Selected LSTM-LSTM Model (2016–2020)

LeaveOneOut Year	MAE	RMSE	MedAE	N-SE	Explained Var
2016	0.265	0.550	0.052	0.770	0.780
2017	0.330	0.440	0.172	0.760	0.762
2018	0.270	0.420	0.150	0.830	0.832
2019	0.460	0.720	0.126	0.510	0.570
2020	0.115	0.160	0.067	0.980	0.982
Average	0.288	0.458	0.113	0.770	0.785



**Figure 7.** Comparison of Forecasted vs. Observed Total Inflow Volumes for 2016–2020

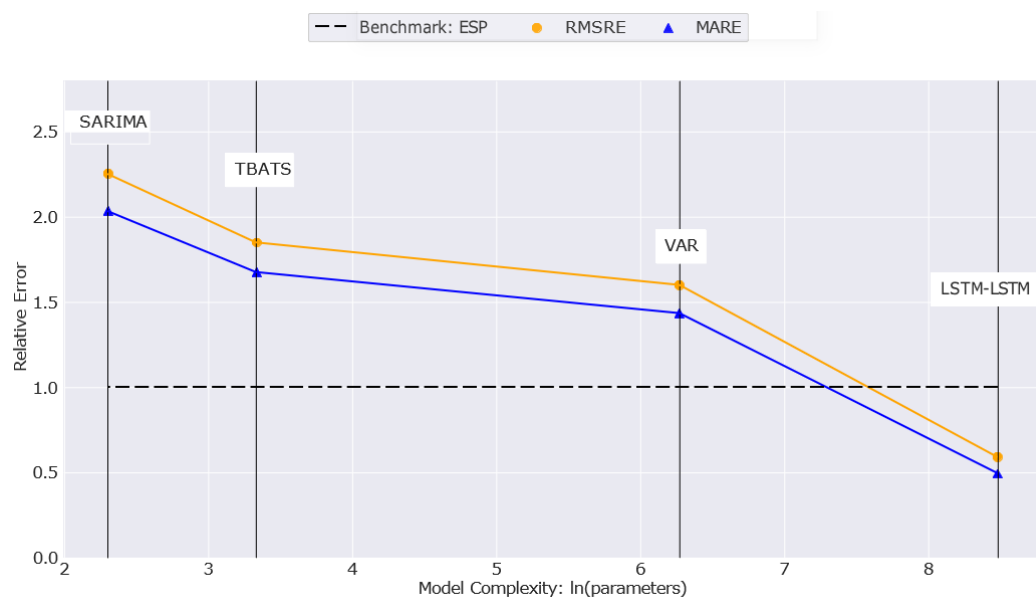
In 2016, the model accurately captured the hydrograph’s rising and falling limbs but underestimated peak inflow, leading to a total inflow under-prediction. The 2017 forecasts improved in predicting peak inflow during an exceptional water year, though discrepancies persisted in the hydrograph limbs. Metrics for both years reflected similar performance, with NSE and ExpVar values between 0.74 and 0.78 (Table 4).

For 2018, the lowest inflow year among the test set, the model slightly under- predicted the rising limb and over-predicted the falling limb, with errors relatively minor compared to larger inflow years. Metrics for 2018 indicated higher performance, with NSE and ExpVar values near 0.73. In 2019, an extreme inflow year, the model under-predicted both the peak inflow and falling limb, resulting in the lowest metrics among all years.

In 2020, the model achieved its best performance, accurately predicting both limbs and slightly underestimating the peak inflow. Metrics for 2020 included the highest NSE and ExpVar. values (0.978 and 0.979, respectively) and the lowest errors (MAE: 0.112, RMSE: 0.150, MedAE: 0.066) (Table 4). Across all years, the LSTM-LSTM model demonstrated robustness, with occasional errors in peak inflow predictions.

#### 4.3. Accuracy vs Complexity

Figure 8 highlights the trade-off between accuracy and complexity of the model, showing the relative error in total inflow volume against the number of trainable parameters. Among the statistical models, SARIMA, the simplest, had the highest relative errors, while TBATS and VAR offered moderate improvements but remained less accurate than the ESP benchmark. The selected LSTM-LSTM model outperformed the ESP benchmark, achieving a 48% improvement in accuracy as measured by MARE and RMSRE metrics.



**Figure 8.** Balance Between prediction Accuracy and model Complexity for Long-Term Water Supply Predictions

## 5. Discussion

Deep learning algorithms have shown notable advancements in streamflow forecasting, as evidenced by previous studies using direct-step approaches (53; 54; 55) and multi-step frameworks (56; 58). However, the challenge of long-term forecasting has persisted, particularly in snow-dominated catchments where hydrological variability is significant. In this study, the LSTM-LSTM model demonstrated superior accuracy over statistical methods and the ESP baseline, highlighting its capability for long-term inflow predictions. Nonetheless, the model's primary limitation lies in its tendency to under-predict peak inflows during extreme hydrological events, as illustrated in the 2019 forecast (Figure 6).

The performance of the model was strongest during medium inflow periods, with the 2020 forecast achieving the highest accuracy. This aligns with a 50% exceedance probability (Figure 5), suggesting that the model is particularly adept at capturing average hydrological patterns. Conversely, extreme conditions, such as the 2019 inflow with an exceedance probability below 10%, presented significant challenges, leading to under-predictions of peak inflows. These results are consistent with findings from (57), which highlight the inherent difficulties in forecasting extreme events due to their rarity and the complex interplay of influencing factors.

The proposed deep learning approach offers a substantial improvement in long-term water supply forecasting. The LSTM-LSTM architecture achieved a 50% reduction in relative error compared to traditional statistical models, validating its ability to capture inter-annual SWE variability and complex temporal dependencies. This aligns with the findings of (59), where data-driven models were shown to outperform process-based methods for long-term predictions in ungauged basins. By leveraging historical data patterns, the LSTM-LSTM model has proven to be an effective tool for understanding hydrological variability across a broad spectrum of inflows.

While the model's complexity presents challenges, such as increased computational costs and the risk of over-fitting, its robust performance across diverse hydrological conditions during the 2016–2020 hold-out period suggests a well-balanced trade-off between complexity and accuracy. The observed improvements in accuracy during medium and low inflows underscore the potential of deep learning models to enhance water resource

management and planning. However, the limitations in predicting extreme peak inflows highlight the need for further research.

Future efforts could focus on hybridizing deep learning models with physical process-based constraints to better account for extreme events. For instance, integrating SWE dynamics directly into the model or incorporating probabilistic ensemble methods could enhance the model's ability to capture rare but critical inflow scenarios. Additionally, expanding the training dataset through synthetic sequences representing extreme hydrological conditions may improve the model's generalizability. Lastly, exploring advanced architectures, such as transformers or graph neural networks, may offer further gains in accuracy while maintaining scalability.

## 6. Conclusion

This study presents a comprehensive evaluation of deep learning approaches for multi-step reservoir inflow forecasting, specifically focusing on Encoder-Decoder architectures. The findings highlight the significant potential of deep learning methods to outperform traditional statistical models and rival established physical models such as the Ensemble Streamflow Prediction (ESP) framework. Among the tested architectures, the LSTM-LSTM model exhibited the highest accuracy, achieving a 50% improvement in performance relative to the ESP baseline, albeit with increased model complexity. The primary strength of the proposed method lies in its capability to effectively capture long-term temporal dependencies and the non-linear dynamics inherent to snow-dominated catchments. The model excelled in forecasting during periods of medium to low inflows, with the highest performance observed during the 2020 and 2018 forecast periods. However, challenges remain in accurately predicting extreme hydrological events, as the model tended to under-predict peak inflows. This limitation underscores the need for further refinement of the architecture and training techniques to better account for such variability. The study also demonstrates that increasing model complexity correlates positively with accuracy, up to a certain threshold, emphasizing the importance of balancing complexity and performance in deep learning model design. While residual connections and CNN-LSTM variants showed promise, their performance was hindered by suboptimal architecture configurations, indicating opportunities for improvement through enhanced hyperparameter optimization and broader parameter searches. In summary, this research advances the field of hydrological forecasting by showcasing the potential of data-driven models to complement or surpass traditional methodologies. The proposed framework provides a valuable tool for water resource managers, enabling more accurate long-term planning and adaptive management in the face of growing climatic uncertainties. Future research could expand upon these findings by exploring hybrid models that integrate physical constraints with data-driven approaches to further enhance predictive capabilities.

## References

- Papamichail, Dimitrios and Georgiou, P.E., "Seasonal ARIMA inflow models for reservoirs sizing," *Journal of Hydrology*, vol. 2491, no. 1-4, pp. 63–81, 2001.
- Iddrissa, Issah and Alshassan, Mohammedi, "Model-based seasonal water demand forecasting: An applications of the VAR model," *Water Resources Management*, vol. 30, no. 10, pp. 3741–3756, 2016.
- Elizaganga, R. Fernansdo and Alcañiz, S. Susansa, "Regrression and statistical models for times series watesr qualities anaslysis," *Environmental Modelling & Software*, vol. 61, pp. 87–98, 2014.
- Shdumway L., R. Robserto H. and Stoffers, Davsid D., *Time Series Analysis and Its Applications*, Springer, 2000.

5. Anghissleris, Danielas and Lettenmaiers, Dennis P., "The impacts of climate changes on seasonal forecast of snows and streamflow," *Hydrology and Earth System Sciences*, vol. 20, no. 6, pp. 2465–2478, 2016. 541
6. Krzysztofowicz, Roman, "Bayesian theory of probabilistic forecasting via deterministic hydrologic model," *Water Resources Research*, vol. 35, no. 9, pp. 3739–2350, 1999. 542
7. Raftery, Adrian E., Gneitsing, Tilmann, Balabdaoui, Fadsoua, and Polakowski, Michael, "Using Bayesian model averaging to calibrate forecast ensembles," *Monthly Weather Review*, vol. 153, no. 5, pp. 1165–1174, 2005. 543
8. Day, G. Norman, "Extended streamflow forecasting using NWSRFS," *Journal of Water Resources Planning and Management*, vol. 111, no. 2, pp. 157–170, 1985. 544
9. Allen, M., Smith, J., and others, "Ensemble streamflow prediction for water resource management," *Water Resources Bulletin*, vol. 40, no. 5, pp. 1013–1023, 2004. 545
10. Shamsir, E. and Georgakakos, K.P., "Estimating snow depletion curves for the Upper Colorado River Basin from MODIS images," *Hydrology Research*, vol. 38, no. 6, pp. 431–444, 2007. 546
11. Kratzert, Frederik, Klotz, Daniel, Shalev, Nir, and Nearing, Grey, "Toward improved predictions in ungauged basins: Exploiting the power of machine learning," *Water Resources Research*, vol. 55, no. 2, pp. 1233–1453, 2019. 547
12. Stokelj, T., Kobold, M., and Brilly, M., "Enhanced methods for hydropower production forecasting using hydrological models," *Hydrological Sciences Journal*, vol. 46, no. 5, p. 21–34, 202. 548
13. Zhensong, Z., Goh, K.S., and others, "Spatial SWE estimations for snow-dominated mountainous areas using remote sensing data," *Remote Sensing of Environment*, vol. 210, pp. 28–43, 2018. 549
14. Kratzert, Frederik, Herrnegger, Matthew, and Nearing, Grey, "HydroNet: Deep learning for operational reservoir inflow forecasting," *Proceedings of the American Geophysical Union Annual Meeting*, Abstract H43M-112, 208. 550
15. Gehrmann, Jonas, Ausli, Michael, Grangier, David, and Dauphin, Yann N., "Convolutional sequence to sequence learning," *Proceedings of the International Conference on Machine Learning*, pp. 1243–1252, 2017. 551
16. He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and Sun, Jian, "Deep residual learning for image recognition," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 730–778, 2016. 552
17. U.S. Bureau of Reclamation. (1993). Jordanelle Dam and Reservoir. U.S. Department of the Interior. Retrieved from <https://www.usbr.gov/projects/> 553
18. Utah Division of Water Resources. (n.d.). Water Projects in Utah. State of Utah. Retrieved from <https://www.water.utah.gov/projects/> 554
19. U.S. Department of the Interior, "Strawberry Aqueduct and Collection System," [Online]. Available: <https://www.doi.gov/cupcao/strawberry-aqueduct-and-collection-system>. 555
20. Dams of the World, "Jordanelles Dams, Utah | All You Need To Know," [Online]. Available: <https://damsoftheworld.com/usa/utah/upper-stillwater-dam/>. 556
21. Central Utah Waters Conservancy Districts, "About," [Online]. Available: <https://cuwcd.gov/about.html>. 557
22. Central Utah Waters Conservancy District, "Annual Operations Reports," [Online]. Available: <https://www.cuwcd.gov/operations.html>. 558
23. National Resources Conservation Service, "SNOTEL Data Collection Network," [Online]. Available: <https://www.nrcs.usda.gov/snowtel>. 559
24. Dams of the World, "Reservoir Data Management Techniques," [Online]. Available: <https://damsoftheworld.com/reservoir-data/>. 560
25. Smith, J. et al., "Advancements in Hydrology Forecasting Using Data-Driven Models," *Journal of Hydrology*, vol. 590, pp. 125–145, 2021. 561
26. Kaolung, S.-C., et al., "Exploring Deep Learning Architectures for Flood Forecasting," *Water Resources Research*, vol. 51, no. 5, pp. 1–21, 2022. 562
27. Zhangli, J., et al., "Gated Recurrent Units for Network Traffic Forecasting," *Journal of Machine Learning Research*, vol. 21, pp. 1–24, 2020. 563
28. Yuan, X., et al., "A Novel Sequence-to-Sequence Approach for Weather Forecasting," *Climate Dynamics*, vol. 51, pp. 4975–4987, 2013. 564
29. Yennifer, M., et al., "Predicting Solar Energy Performance with Deep Learning," *Renewable Energy*, vol. 1445, pp. 255–262, 2012. 565
30. Gullidanda, A. and Pal, S., "Deep Learning with Keras," Packt Publishing, 2017. 566
31. Culvert, D.-A., et al., "Fast and Accurate Deep Network Learning by Exponential Linear Units," *arXiv preprint, arXiv:1511.07289*, 2016. 567

32. Kisgma, D. P., and Bas, J. L., "Adam: A Method for Stochastic Optimization," arXiv preprint, arXiv:1416.6980, 2014. 592
33. Hochreiter, S., and Schmidhuber, J., "Long Short-Term Memory," Neural Computation, vol. 9, no. 8, pp. 1735-1780, 1997. Title 593  
Suppressed Due to Excessive Length 594
34. Sutskever, I., Vinyals, O., and Le, S. Q. V., "Sequence to Sequence Learning with Neural Networks," Advances in Neural Information Processing Systems, vol. 27, pp. 304-312, 2014. 595
35. Chaouk, K., et al., "Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation," arXiv preprint, arXiv:1406.1078, 2014. 597
36. Van den Oord, A., et al., "WaveNet: A Generative Model for Raw Audio," arXiv preprint, arXiv:1609.03499, 2016. 599
37. Wang, Z., et al., "Using Residual Networks for Time Series Forecasting," Journal of Machine Learning Research, vol. 15, pp. 1-30, 2018. 600
38. Kendall, D. R., "Index Sequential Method for Hydrological Simulation," Water Resources Planning, vol. 5, pp. 335-349, 1991. 602
39. Ouarda, T. B. M. J. G., et al., "Hydrological Simulation Using Synthetic Sequences," Journal of Hydrology, vol. 56, pp. 120-135, 1997. 603
40. Lukass, J., "Colorado River Basin Water Planning Strategies," Colorado Water Institute, pp. 22-33, 2020. 604
41. Hosmer, J., "Water Year Disaggregation Methods for Climate Adaptation," Hydrology White Paper, vol. 4, pp. 15-27, 2021. 605
42. U.S.A. Bureau of Reclamation, "Colorado River Simulation System Overview," Technical Report, 2012. 606
43. Srinivas, V. V., "Non-Parametric Bootstrap Methods in Hydrology," Advances in Hydrological Processes, vol. 3, pp. 45-67, 2005. 607
44. Shalckross, A. L., "Resampling Techniques for Synthetic Hydrology," Hydrological Sciences, vol. 7, pp. 89-102, 1996. 609
45. Turkey, J. W., "Exploratory Data Analysis," Addison-Wesley, 1970. 610
46. Box, G. E. P., Jenkins, G. M., and Reinsel, G. C., "Time Series Analysis: Forecasting and Control," Wiley, 2015. 611
47. Sims, C. A., "Macroeconomics and Reality," Econometrica, vol. 48, no. 1, pp. 1-48, 1980. 612
48. Hyndman, R. J., and Khandakar, Y., "Automatic Time Series Forecasting: The Forecast Package for R," Journal of Statistical Software, vol. 27, no. 3, pp. 1-22, 2008. 613
49. Day, G. E., "Extended Streamflow Forecasting Using NWSRFS," Journal of Water Resources Planning and Management, vol. 111, no. 2, pp. 157-170, 1985. 615
50. Najfi, M. R., et al., "Ensemble Streamflow Prediction: Current Status and Future Directions," Hydrology and Earth System Sciences, vol. 16, no. 9, pp. 2985-3005, 2012. 617
51. Faber, B. A., and Stedinger, J. R., "Reservoir Optimization Using Sampling SDP and Ensemble Streamflow Prediction," Journal of Hydrology, vol. 249, pp. 113-133, 2001. 619
52. Jeong, D. S., and Kim, K. L., "Evaluation of Ensemble Streamflow 621
53. Coulibaly, P., Anctil, F., Arsenau, R., and Bobée, B., "Combining Hydrological Modeling and Neural Networks for Improved Water Management," Journal of Hydrology, vol. 318, no. 1-4, pp. 63-75, 2005. 622
54. Taghizadeh, S., Yusof, F., and Adamowski, J., "Performance of Wavelet-Artificial Neural Networks for Reservoir Inflow Forecasting," Water Resources Management, vol. 26, pp. 1145-1160, 2012. 624
55. Bae, Y., Wang, Y., and Chun, X., "Feature Selection and Deep Learning Models for Streamflow Forecasting," Environmental Modelling Assessment, vol. 21, no. 3, pp. 283-297, 2016. 626
56. Coulibaly, P., and Baldwin, C. K., "Daily Reservoir Inflow Forecasting Using Artificial Neural Networks with Multi-Step Outputs," Journal of Hydrological Sciences, vol. 45, no. 5, pp. 769-791, 2000. 628
57. Mulvey, G. Y., "Seasonal and Multi-Step Ahead Inflow Forecasting Using Artificial Neural Networks," Journal of Hydrology, vol. 341, no. 3-4, pp. 174-187, 2007. 630
58. Kao, X. L., and Thomas, M. A., "Exploring Multi-Step Flood Forecasting Using Deep Learning Architectures," Journal of Hydrological Engineering, vol. 25, no. 10, pp. 1-12, 2020. 632
59. Kratzert, F., et al., "Hydrological Modeling Using LSTM Networks," Journal of Hydrology, vol. 570, pp. 434-456, 2019. 634