# Probabilistic Hierarchical Interpolation and Interpretable Configuration for Flood Prediction

Mostafa Saberian[1], Vidya Samadi[2*,] Ioana Popescu[3]

1. The Glenn Department of Civil Engineering, Clemson University, Clemson, SC
2. Department of Agricultural Sciences, Clemson University, Clemson, SC.
3. Department of Hydroinformatics and Socio-Technical Innovation, IHE Delft Institute for Water Education, Delft, the Netherlands

*Corresponding author: samadi@clemson.edu

## Abstract

The last few years have witnessed the rise of neural networks (NNs) applications for hydrological time series modeling. By virtue of their capabilities, NN models can achieve unprecedented levels of performance when learn how to solve increasingly complex rainfall-runoff processes via data, making them pivotal for the development of computational hydrologic tasks such as flood predictions. The NN models should, in order to be considered practical, provide a probabilistic understanding of the model mechanisms and predictions and hints on what could perturb the model. In this paper, we developed two probabilistic NN models, i.e., Neural Hierarchical Interpolation for Time Series Forecasting (N-HiTS) and Network-Based Expansion Analysis for Interpretable Time Series Forecasting (N-BEATS) and benchmarked them with long short-term memory (LSTM) for flood prediction across two headwater streams in Georgia and North Carolina, USA. To generate a probabilistic prediction, a Multi-Quantile Loss was used to assess the 95th percentile prediction uncertainty (95PPU) of multiple flooding events. We conducted extensive flood prediction experiments demonstrating the advantages of hierarchical interpolation and interpretable architecture, where both N-HiTS and N-BEATS provided an average accuracy improvement of almost 5% (NSE) over the LSTM benchmarking model. On a variety of flooding events, both N-HiTS and N-BEATS demonstrated significant performance improvements over the LSTM benchmark and showcased their probabilistic predictions by specifying a likelihood parameter.

**Keywords:** Probabilistic Flood Prediction; Neural Networks; N-HiTS; N-BEATS; LSTM; Headwater Stream.

## Key Points

- N-HiTS and N-BEATS predictions reflect interpretability and hierarchical representations of data to reduce neural network complexities.
- Both N-HiTS and N-BEATS models outperformed the LSTM in mathematically defining uncertainty bands.
- Predicting the magnitude of the recession curve of flood hydrographs was particularly challenging for all models.

**Plain Language Summary**

Recent progress in NN accelerated improvements in the performance of catchment modeling. Yet flood modeling remains a very difficult task. Focusing on two headwater streams, we developed N-HiTS and N-BEATS models and benchmarked them with LSTM to predict flooding. N-HiTS and N-BEATS outperformed LSTM for flood predictions. We demonstrated how the proposed models can be augmented with an uncertainty approach to predict flooding that is interpretable without considerable loss in accuracy.

**1. Introduction**

The last few years have been characterized by an upsurge in the neural networks (NN) applications in hydrology. As opaque NN models are increasingly being employed to make important hydrological predictions, the demand for creating legitimate NN models is increasing in the hydrology community. However, maintaining coherence while producing accurate predictions can be a challenging problem (Olivares et al., 2024). There is a general agreement on the importance of providing probabilistic NN prediction (Samadi et al., 2020), especially in the case of flood prediction (Martinaitis et al., 2023).

Flood occurrences have witnessed an alarming surge in frequency and severity globally. Jonkman (2005) studied a natural disaster database (EM-DAT, 2023) and reported that over 27 years, more than 175000 people died, and close to 2.2 billion were affected directly by floods worldwide. These numbers are likely an underestimation due to unreported events (Nevo et al., 2022). In addition, the United Nations Office for Disaster Risk Reduction reported that flooding has been the most frequent, widespread weather-related natural disaster since 1995, claiming over 600,000 lives, affecting around 4 billion people globally, and causing annual economic damage of more than 100 billion USD (UNISDR, 2015). This escalating trend has necessitated the need for better flood prediction and management strategies. Scholars have successfully implemented different flood models such as deterministic (e.g., Roelvink et al., 2009, Thompson and Frazier, 2014; Barnard et al., 2014; Erikson et al., 2018) and physically based flood models (e.g., Basso et al., 2016; Chen et al., 2016; Pourreza-Bilondi et al., 2017; Saksena et al., 2019; Refsgaard et al., 2021) in various environmental systems over the past several decades. These studies have heightened the need for precise flood prediction, they have also unveiled limitations inherent in existing deterministic and physics-based models. While evidence suggests that both deterministic and physics-based approaches are meaningful and useful (Sukovich et al., 2014; Zafarmomen et al., 2024), their forecasts rest heavily on imprecise and subjective expert opinion; there is a challenge for setting robust evidence-based thresholds to issue flood warnings and alerts (Palmer, 2012). Moreover, many of these traditional flood models particularly physically explicit models rely heavily on a particular choice of numerical approximation and describe multiple process parameterizations only within a fixed spatial architecture (e.g., Clark et al., 2015). Recent NN models have shown promising results across a large variety of flood modeling applications (e.g.,

68  Nevo et al., 2022; Pally and Samadi, 2022; Dasgupta et al., 2023; Zhang et al., 2023) and encourage the

69  use of such methodologies as core drivers for neural flood prediction (Windheuser et al., 2023).

70  Earlier adaptations of these intelligent techniques showed promising for flood prediction (e.g., Hsu et al.,

71  1995; Tiwari and Chatterjee, 2010). However, recent efforts have taken NN application to the next level,

72  providing uncertainty assessment (Sadeghi Tabas and Samadi, 2022) and improvements over various

73  spatio-temporal scales, regions, and processes (e.g., Kratzert et al., 2018; Park and Lee, 2023; Zhang et al.,

74  2023). Nevo et al., (2022) were the first scholars who employed long short-term memory (LSTM) for flood

75  stage prediction and inundation mapping, achieving notable success during the 2021 monsoon season. Soon

76  after, Russo et al. (2023) evaluated various NN models for predicting flood depth in urban systems,

77  highlighting the potential of data-driven models for urban flood prediction. Similarly, Defontaine et al.

78  (2023) emphasized the role of NN algorithms in enhancing the reliability of flood predictions, particularly

79  in the context of limited data availability. Windheuser et al., (2023) studied flood gauge height forecasting

80  using images and time series data for two gauging stations in Georgia, USA. They used multiple NN models

81  such as Convolutional Neural Network (ConvNet/CNN) and LSTM to forecast floods in near real-time (up

82  to 72 hours). In a sequence, Wee et al., (2023) used Impact-Based Forecasting (IBF) to propose a Flood

83  Impact-Based Forecasting system (FIBF) using flexible fuzzy inference techniques, aiding decision-makers

84  in a timely response. Zou et al. (2023) proposed a Residual LSTM (ResLSTM) model to enhance and

85  address flood prediction gradient issues. They integrated Deep Autoregressive Recurrent (DeepAR) with

86  four recurrent neural networks (RNNs), including ResLSTM, LSTM, Gated Recurrent Unit (GRU), and

87  Time Feedforward Connections Single Gate Recurrent Unit (TFC-SGRU). They showed that ResLSTM

88  achieved superior accuracy. While these studies reported the superiority of NN models for flood modeling,

89  they highlighted a number of challenges, notably (i) the limited capability of proposed NN models to

90  capture the spatial variability and magnitudes of extreme data over time, (ii) the lack of a sophisticated

91  mechanism to capture different flood magnitudes and synthesize the prediction, and (iii) inability of the NN

92  models to process data in parallel and capture the relationships between all elements in a sequential manner.

93  Recent advances in neural time series forecasting showed promising results that can be used to address the

94  above challenges for flood prediction. Recent techniques include the adoption of the attention mechanism

95  and Transformer-inspired approaches (Fan et al. 2019; Alaa and van der Schaar 2019; Lim et al. 2021)

96  along with attention-free architectures composed of deep stacks of fully connected layers (Oreshkin et al.

97  2020).  All of these approaches are relatively easy to scale up in terms of flood magnitudes (small to major

98  flood predictions), compared to LSTM and have proven to be capable of capturing spatiotemporal

99  dependencies (Challu et al., 2022). In addition, these architectures can capture input-output relationships

100  implicitly while they tend to be more computationally efficient. Many state-of-the-art NN approaches for

101  flood forecasting have been established based on LSTM. There are cell states in the LSTM networks that

102 can be interpreted as storage capacity often used in flood generation schemes. In LSTM, the updating of

103 internal cell states (or storages) is regulated through a number of gates: the first gate regulates the storage

104 depletion, the second one regulates storage fluctuations, and the third gate regulates the storages outflow

105 (Tabas and Samadi, 2022). The elaborate gated design of the LSTM partly solves the long-term dependency

106 problem in flood time series prediction (Fang et al., 2020), although, the structure of LSTMs is designed in

107 a sequential manner that cannot directly connect two nonadjacent portions (positions) of a time series.

108 In this paper, we developed attention-free architecture, i.e. Neural Hierarchical Interpolation for Time

109 Series Forecasting (N-HiTS; Challu et al., 2022) and Network-Based Expansion Analysis for Interpretable

110 Time Series Forecasting (N-BEATS; Oreshkin et al., 2020) and benchmarked these models with LSTM for

111 flood prediction. We developed fully connected N-BEATS and N-HiTS architectures using multi-rate data

112 sampling, synthesizing the flood prediction outputs via multi-scale interpolation.

113 We implemented all algorithms for flood prediction on two headwater streams i.e., the Lower Dog River,

114 Georgia, and the Upper Dutchmans Creek, North Carolina, USA to ensure that the results are reliable and

115 comparable. The results of N-BEATS and N-HiTS techniques were compared with the benchmarking

116 LSTM to understand how these techniques can improve the representations of rainfall and runoff

117 dispensing over a recurrence process. Notably, this study represents a pioneering effort, as to the best of

118 our knowledge, this is the first instance in which the application of N-BEATS and N-HiTS algorithms in

119 the field of flood prediction has been explored. The scope of this research will focus on:

120

121 **(i)Flood prediction in a hierarchical fashion with interpretable outputs:** We built N-BEATS and N-

122 HiTS for flood prediction with a very deep stack of fully connected layers to implicitly capture input-output

123 relationships with hierarchical interpolation capabilities. The predictions also involve programming the

124 algorithms with decreasing complexity and aligning their time scale with the final output through multi-

125 scale hierarchical interpolation and interpretable architecture. Predictions were aggregated in a hierarchical

126 fashion that enabled the building of a very deep neural network with interpretable configurations.

127 **(ii)     Uncertainty quantification of the models by employing probabilistic approaches:** a Multi-

128 Quantile Loss (MQL) was used to assess the 95th percentile prediction uncertainty (95PPU) of multiple

129 flooding events. MQL was integrated as the loss function to account for probabilistic prediction. MQL

130 trains the model to produce probabilistic forecasts by predicting multiple quantiles of the distribution of

131 future values.

132 **(iii)     Exploring headwater stream response to flooding:** Understanding the dynamic response of

133 headwater streams to flooding is essential for managing downstream flood risks. Headwater streams

134 constitute the uppermost sections of stream networks, usually comprising 60% to 80% of a catchment area.

135 Given this substantial coverage and the tendency for precipitation to increase with elevation, headwater

136 streams are responsible for generating and controlling the majority of runoff in downstream portions
137 (MacDonald and Coe, 2007).
138 The remainder of this paper is structured as follows. Section 2 presents the case study and data, NN models,
139 performance metrics, and sensitivity and uncertainty approaches. Section 3 focuses on the results of flood
140 predictions including sensitivity and uncertainty assessment and computation efficiency. Finally, Section 4
141 concludes the paper.
142

143 **2. Methodology**

144 **2.1. Case Study and Data**

145 This research used two headwater gauging stations located at the Lower Dog River watershed, Georgia
146 (GA; USGS02337410, Dog River gauging station), and the Upper Dutchmans Creek watershed, North
147 Carolina (NC; USGS0214269560, Killian Creek gauging station). As depicted in Figures 1 and 2, the Lower
148 Dog River and the Upper Dutchmans Creek watersheds are located in the west and north parts of two
149 metropolitan cities, Atlanta and Charlotte. As shown in Figure 1, the Lower Dog River stream gauge is
150 established southeast of Villa Rica in Carroll County, where the USGS has regularly monitored discharge
151 data since 2007 in 15-minute increments. The Lower Dog River is a stream with a length of 15.7 miles
152 (25.3 km; obtained from the U.S. Geological Survey [USGS] National Hydrography Dataset high-
153 resolution flowline data), an average elevation of 851.94 meters, and the watershed area above this gauging
154 station is 66.5 square miles (172 km2; obtained from the Georgia Department of Natural Resources). This
155 watershed is covered by 15.2% residential area, 14.6% agricultural land, and ~70% forest (Munn et al.,
156 2020). Killian Creek gauging station at the Upper Dutchmans Creek watershed is established
157 in Montgomery County, NC, where the USGS has regularly monitored discharge data since 1995 in 15-
158 minute increments. The Upper Dutchmans Creek is a stream with a length of 4.9 miles (7.9 km), an average
159 elevation of 642.2 meters (see Table 1), and the watershed area above this gauging station is 4 square miles
160 (10.3 km2) with less than 3% residential area and about 93% forested land use (the United States
161 Environmental Protection Agency).
162  The Lower Dog River has experienced significant flooding in the last decades. For example, in September
163 2009, the creek, along with most of northern GA, experienced heavy rainfall (5 inches, equal to 94 mm).
164 The Lower Dog River, overwhelmed by large amounts of overland flow from saturated ground in the
165 watershed, experienced massive flooding in September 2009 (Gotvald, 2010). The river crested at 33.8 feet
166 (10.3 m) with a peak discharge of 59,900 cfs (1,700 m3/s), nearly six times the 100-year flood level
167 (McCallum and Gotvald, 2010). In addition, Dutchmans Creek has experienced significant flooding in
168 February 2020. According to local news (WCCB Charlotte, 2020), the flood in Gaston County caused

169     significant infrastructure damage and community disruption. Key impacts included the threatened collapse

170     of the Dutchman's Creek bridge in Mt. Holly and the closure of Highway 7 in McAdenville, GA.

171

172                     Table 1.  The Lower Dog River and Upper Dutchmans Creek's physical characteristics.

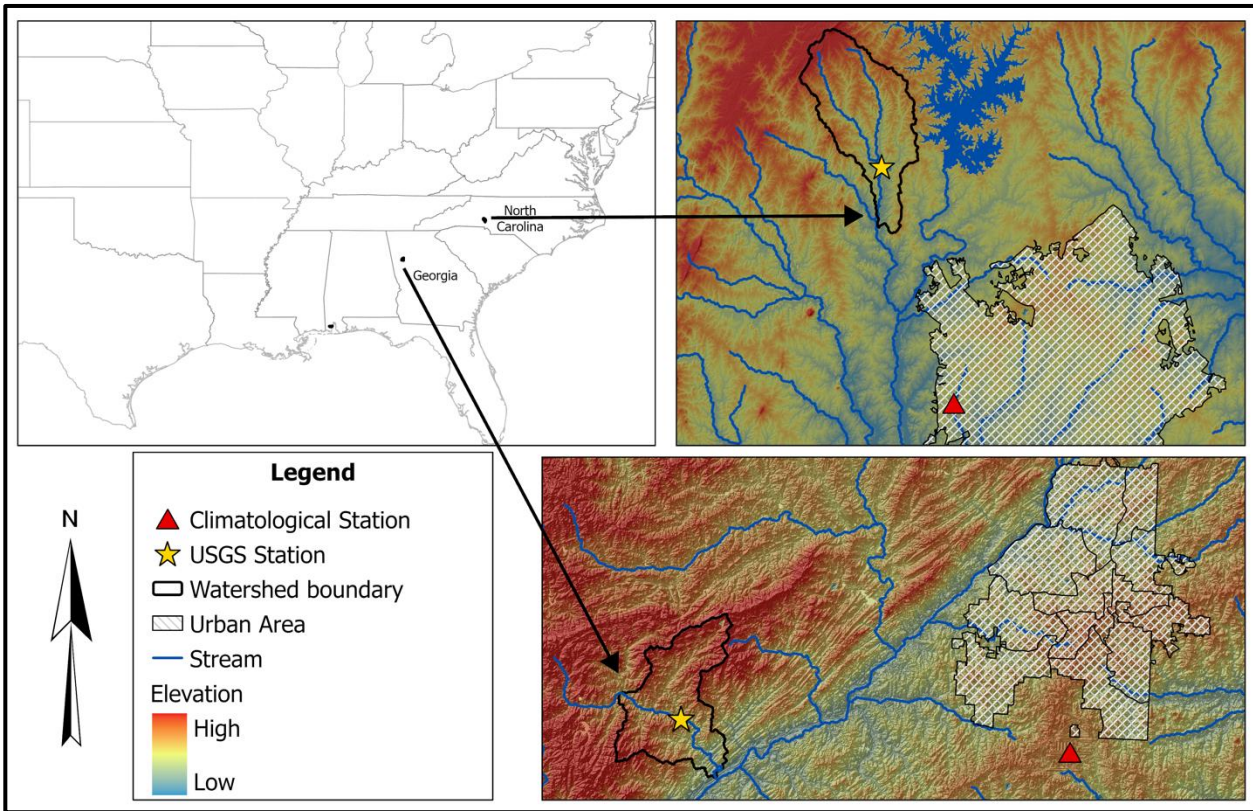| Watershed | USGS Station ID Number | Average Elevation (m) | Stream Length (km) | Watershed area (km2) |
|---|---|---|---|---|
| Lower Dog River watershed, GA | USGS02337410 | 851.9 | 25.3 | 172 |
| Upper Dutchmans Creek watershed, NC | USGS0214269560 | 642.2 | 7.9 | 10.3 |

173



175    Figure 1.  The Lower Dog River and The Upper Dutchmans Creek watersheds are located in GA and NC.
176       The proximity of the watersheds to Atlanta and Charlotte (urban area) are also displayed on the map.

177

178     To provide the meteorological forcing data, i.e., precipitation, temperature, and humidity, were extracted

179     from the National Oceanic and Atmospheric Administration's (NOAA) Local Climatological Data

180     (LCD). We used the NOAA precipitation, temperature, and humidity data of Atlanta Hartsfield Jackson

181     International Airport and Charlotte Douglas Airport stations as an input variable for neural network

182     algorithms. The data has been monitored since January 1, 1948, and July 22, 1941, with an hourly interval
183     which was used as an input variable for constructing neural networks.

184     To fill in the missing values in the data, we used the spline interpolation method. We applied this method
185     to fill the gaps in time series data, although the missing values were insignificant (less than 1%). In addition,
186     we employed the Minimum Inter-Event Time (MIT) approach to precisely identify and separate individual
187     storm events. The MIT-based event delineation is pivotal for accurately defining storm events. This method
188     allowed us to isolate discrete rainfall episodes, aiding a comprehensive analysis of storm events. Moreover,
189     it provided a basis for event-specific examination of flood responses, such as initial condition and cessation
190     (loss), runoff generation, and runoff dynamics.

191     The hourly rainfall dataset consists of distinct rainfall occurrences, some consecutive and others clustered
192     with brief intervals of zero rainfall. As these zero intervals extend, we aim to categorize them into distinct
193     events. It's worth noting that even within a single storm event, we often encounter short periods of no
194     rainfall, known as intra-storm zero values. In the MIT method, we defined a storm event as a discrete rainfall
195     episode surrounded by dry periods both preceding and following it, determined by an MIT (Asquith et al.,
196     2005; Safaei-Moghadam et al., 2023). There are many means to determine an MIT value. One practical
197     approximation is using serial autocorrelation between rainfall occurrences. MIT approach uses
198     autocorrelation that measures the statistical dependency of rainfall data at one point in time with data at
199     earlier, or lagged times within the time series. The lag time represents the gap between data points being
200     correlated. When the lag time is zero, the autocorrelation coefficient is unity, indicating a one-to-one
201     correlation. As the lag time increases, the statistical correlation diminishes, converging to a minimum value.
202     This signifies the fact that rainfall events become progressively less statistically dependent or, in other
203     words, temporally unrelated. To pinpoint the optimal MIT, we analyzed the autocorrelation coefficients for
204     various lag times, observing the point at which the coefficient approaches zero. This lag time signifies the
205     minimum interval of no rainfall, effectively delineating distinct rainfall events.

206     **2.2. NN Algorithms**

207     **2.2.1. LSTM**

208     LSTM is an RNN architecture widely used as a benchmark model for flood neural time series
209     modeling. LSTM networks are capable of selectively learning order dependence in sequence prediction
210     problems (Sadeghi Tabas and Samadi, 2022). These networks are powerful because they can capture the
211     temporal features, especially the long-term dependencies (Hochreiter et al., 2001), and are independent of
212     the length of the input data sequences meaning that each sample is independent from another one.

213    The memory cell state within LSTM plays a crucial role in capturing extended patterns in data, making it
214    well-suited for dynamic time series modeling such as flood prediction. An LSTM cell uses the following
215    functions to compute flood prediction.

$$i_t = \sigma(A_i x_t + B_i h_{t-1} + c_i) \qquad \text{(Equation 1)}$$

$$f_t = \sigma(A_f x_t + B_f h_{t-1} + c_f) \qquad \text{(Equation 2)}$$

$$o_t = \sigma(A_o x_t + B_o h_{t-1} + c_o) \qquad \text{(Equation 3)}$$

$$m_t = f_t \odot m_{t-1} + i_t \odot tanh(A_g x_t + B_g h_{t-1} + c_g) \qquad \text{(Equation 4)}$$

$$h_t = o_t \odot tanh(m_t) \qquad \text{(Equation 5)}$$

216    Where $x_t$ and $h_t$ represent the input and the hidden state at time step $t$, respectively. $\odot$ denotes element-
217    wise multiplication, $tanh$ stands for the hyperbolic tangent activation function, and $\sigma$ represents the
218    sigmoid activation function. $A$, $B$, and $c$ are trainable weights and biases that undergo optimization during
219    the training process. $m_t$ and $h_t$ are cell states at time step $t$ that are employed in the input processing for
220    the next time step. $m_t$ represents the memory state responsible for preserving long-term information, while
221    $h_t$ represents the memory state preserving short-term information. The LSTM cell consists of a forget gate
222    $f_t$, an input gate $i_t$ and an output gate $o_t$ and has a cell state $m_t$. At every time step $t$, the cell gets the data
223    point $x_t$ with the output of the previous cell $h_{t-1}$ (Windheuser et al., 2023). The forget gate then defines if
224    the information is removed from the cell state, while the input gate evaluates if the information should be
225    added to the cell state and the output gate specifies which information from the cell state can be used for
226    the next cells.

227    We used two LSTM layers with 128 cells in the first two hidden layers as encoder layers, which were then
228    connected to two multilayer perceptron (MLP) layers with 128 neurons as decoder layers. The LSTM
229    simulation was performed with these input layers along with the *Adam* optimizer (Kingma and Ba,
230    2014), *tanh* activation function, and a single lagged dependent-variable value to train with a learning rate
231    of 0.001. The architecture of the proposed LSTM model is illustrated in Figure 2.
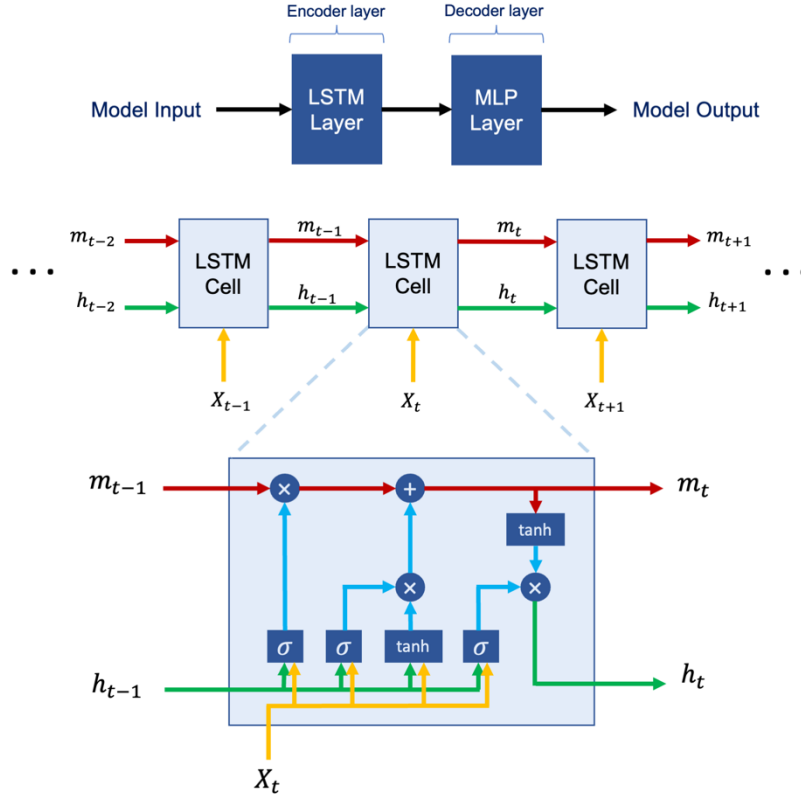
Figure 2. The structure of LSTM programmed in this research. We used *tanh* and *sigmoid* as activation functions along with 2 layers of LSTM, 2 layers of MLP, and 128 cells in each layer.

### 2.2.2. N-BEATS

N-BEATS is a deep learning architecture based on backward and forward residual links and the very deep stack of fully connected layers specifically designed for sequential data forecasting tasks (Oreshkin et al., 2020). This architecture has a number of desirable properties including interpretability. The N-BEATS architecture distinguishes itself from existing architectures in several ways. First, the algorithm approaches forecasting as a non-linear multivariate regression problem instead of a sequence-to-sequence challenge. Indeed, the core component of this architecture (as depicted in Figure 3) is a fully connected non-linear regressor, which takes the historical data from a time series as input and generates multiple data points for the forecasting horizon. Second, the majority of existing time series architectures are quite limited in depth, typically consisting of one to five LSTM layers. N-BEATS employs the residual principle to stack a substantial number of layers together, as illustrated in Figure 3. In this configuration, the basic block not only predicts the next output but also assesses its contribution to decomposing the input, a concept that is referred to as "backcast" (see Oreshkin et al. 2020).

9

249    The basic building block in the architecture features a fork-like structure, as illustrated in Figure 3 (bottom).
250    The $l$-th block (for the sake of brevity, the block index $l$ is omitted from Figure 3) takes its respective input,
251    $x_l$, and produces two output vectors: $\hat{x}_l$ and $\hat{y}_l$. In the initial block of the model, $x_l$ corresponds to the
252    overall model input, which is a historical lookback window of a specific length, culminating with the most
253    recent observed data point. For the subsequent blocks, $x_l$ is derived from the residual outputs of the
254    preceding blocks. Each block generates two distinct outputs: 1. $\hat{y}_l$: This represents the forward forecast of
255    the block, spanning a duration of $H$ time units. 2. $\hat{x}_l$: This signifies the block's optimal estimation of $x_l$,
256    which is referred to "backcast." This estimation is made within the constraints of the functional space
257    available to the block for approximating signals (Oreshkin et al., 2020).

258    Internally, the fundamental building block is composed of two elements. The initial element involves a
259    fully connected network, which generates forward expansion coefficient predictors, $\theta_l^f$, and a backward
260    expansion coefficient predictor, $\theta_l^b$. The second element encompasses both backward basis layers, $g_l^b$, and
261    forward basis layers, $g_l^f$. These layers take the corresponding forward $\theta_l^f$ and backward $\theta_l^b$ expansion
262    coefficients as input, conduct internal transformations using a set of basis functions, and ultimately yield
263    the backcast, $\hat{x}_l$, and the forecast outputs, $\hat{y}_l$, as previously described by Oreshkin et al. (2020). The
264    following equations describe the first element:

$$h_{l,1} = FC_{l,1}(x_l), \quad h_{l,2} = FC_{l,2}(h_{l,1}), \quad h_{l,3} = FC_{l,3}(h_{l,2}), \quad h_{l,4} = FC_{l,4}(h_{l,3}). \quad \text{(Equation 6)}$$

$$\theta_l^b = \text{LINEAR}_l^b(h_{l,4}), \qquad \theta_l^b = \text{LINEAR}_l^b(h_{l,4}) \qquad \text{(Equation 7)}$$

265    The LINEAR layer, in essence, functions as a straightforward linear projection, meaning $\theta_l^f = W_l^f h_{l,4}$. As
266    for the fully connected (FC) layer, it takes on the role of a conventional FC layer, incorporating RELU non-
267    linearity as an activation function.

268    The second element performs the mapping of expansion coefficients $\theta_l^f$ and $\theta_l^b$ to produce outputs using
269    basis layers, resulting in $\hat{y}_l = g_l^f(\theta_l^f)$ and $\hat{x}_l = g_l^b(\theta_l^b)$. This process is defined by the following equation:

$$\hat{y}_l = \sum_{i=1}^{\dim(\theta_l^f)} \theta_{l,i}^f v_i^f, \qquad \hat{x}_l = \sum_{i=1}^{\dim(\theta_l^b)} \theta_{l,i}^b v_i^b \qquad \text{(Equation 8)}$$

270    Within this context, $v_i^f$ and $v_i^b$ represent the basis vectors for forecasting and backcasting, respectively,
271    while $\theta_{l,i}^f$ corresponds to the $i$-th element of $\theta_l^f$.

272     The N-BEATS uses a novel hierarchical doubly residual architecture which is illustrated in Figure 3 (top
273     and middle). This framework incorporates two residual branches, one traversing the backcast predictions
274     of each layer, while the other traverses the forecast branch of each layer. The following equation describes
275     this process:

$$x_l = x_{l-1} - \hat{x}_{l-1} \quad , \quad \hat{y} = \sum_l \hat{y}_l \qquad \text{(Equation 9)}$$

276     As mentioned earlier, in the specific scenario of the initial block, its input corresponds to the model-level
277     input $x$. In contrast, for all subsequent blocks, the backcast residual branch $x_l$ can be conceptualized as
278     conducting a sequential analysis of the input signal. The preceding block eliminates the portion of the signal
279     $\hat{x}_{l-1}$ that it can effectively approximate, thereby simplifying the prediction task for downstream blocks.
280     Significantly, each block produces a partial forecast $\hat{y}_l$ , which is initially aggregated at the stack level and
281     subsequently at the overall network level, establishing a hierarchical decomposition. The ultimate forecast
282     $\hat{y}$ is the summation of all partial forecasts (Oreshkin et al., 2020).

283     The N-BEATS model has two primary configurations: generic and interpretable. These configurations
284     determine how the model structures its blocks and how it processes time series data. In the generic
285     configuration, the model uses a stack of generic blocks that are designed to be flexible and adaptable to
286     various patterns in the time series data. Each generic block consists of fully connected layers with ReLU
287     activation functions. The key characteristic of the generic configuration is its flexibility. Since the blocks
288     are not specialized for any specific pattern (like trend or seasonality), they can learn a wide range of patterns
289     directly from the data (Oreshkin et al., 2020). In the interpretable configuration, the model architecture
290     integrates distinct trend and seasonality components. This involves structuring the basis layers at the stack
291     level specifically to model these elements, allowing the stack outputs to be more easily understood.

292     **Trend Model:** In this stack $g_{s,l}^b$ and $g_{s,l}^f$ are polynomials of a small degree $p$, functions that vary slowly
293     across the forecast window, to replicate monotonic or slowly varying nature of trends:

$$\hat{y}_{s,l} = \sum_{i=0}^{p} \theta_{s,l,i}^f t^i \qquad \text{(Equation 10)}$$

294     The time vector $t = [0, 1, 2, \dots, H-2, H-1]^T/H$ is specified on a discrete grid ranging from 0 to
295     *(H−1)/H*, projecting $H$ steps into the future. Consequently, the trend forecast represented in matrix form is:

$$\hat{y}_{s,l}^{tr} = T\theta_{s,l}^f \qquad \text{(Equation 11)}$$

296 Where the polynomial coefficients, $\theta_{s,l}^f$, predicted by an FC network at layer $l$ of stack $s$, are described by

297 Equations (6) and (7). The matrix $T$, consisting of powers of $t$, is represented as $[1, t, \ldots, t^p]$. When $p$ is

298 small, such as 2 or 3, it compels $\hat{y}_{s,l}^{tr}$ to emulate a trend (Oreshkin et al., 2020).

299 Seasonality model: In this stack $g_{s,l}^b$ and $g_{s,l}^f$ are periodic functions, to capture the cyclical and recurring

300 characteristics of seasonality, such that $y_t = y_{t-\Delta}$, where $\Delta$ is the seasonality period. The Fourier series

301 serves as a natural foundation for modeling periodic functions:

$$\hat{y}_{s,l} = \sum_{i=0}^{\frac{H}{2}-1} \theta_{s,l,i}^f \cos(2\pi it) + \theta_{s,l,i+[H/2]}^f \sin(2\pi it) \qquad \text{(Equation 12)}$$

302

303 Consequently, the seasonality forecast is represented in the following matrix form:

$$\hat{y}_{s,l}^{seas} = S\theta_{s,l}^f \qquad \text{(Equation 13)}$$

$$S = \left[1, \cos(2\pi t), \ldots, \cos\left(2\pi\left[\frac{H}{2} - 1\right]t\right), \sin(2\pi t), \ldots, \sin\left(2\pi\left[\frac{H}{2} - 1\right]t\right)\right] \qquad \text{(Equation 14)}$$

304

305 Where the Fourier coefficients $\theta_{s,l}^f$, that predicted by an FC network at layer l of stack $s$, are described by

306 Equations (6) and (7). The matrix $S$ represents sinusoidal waveforms. As a result, the forecast $\hat{y}_{s,l}^{seas}$

307 becomes a periodic function that imitates typical seasonal patterns (Oreshkin et al., 2020).
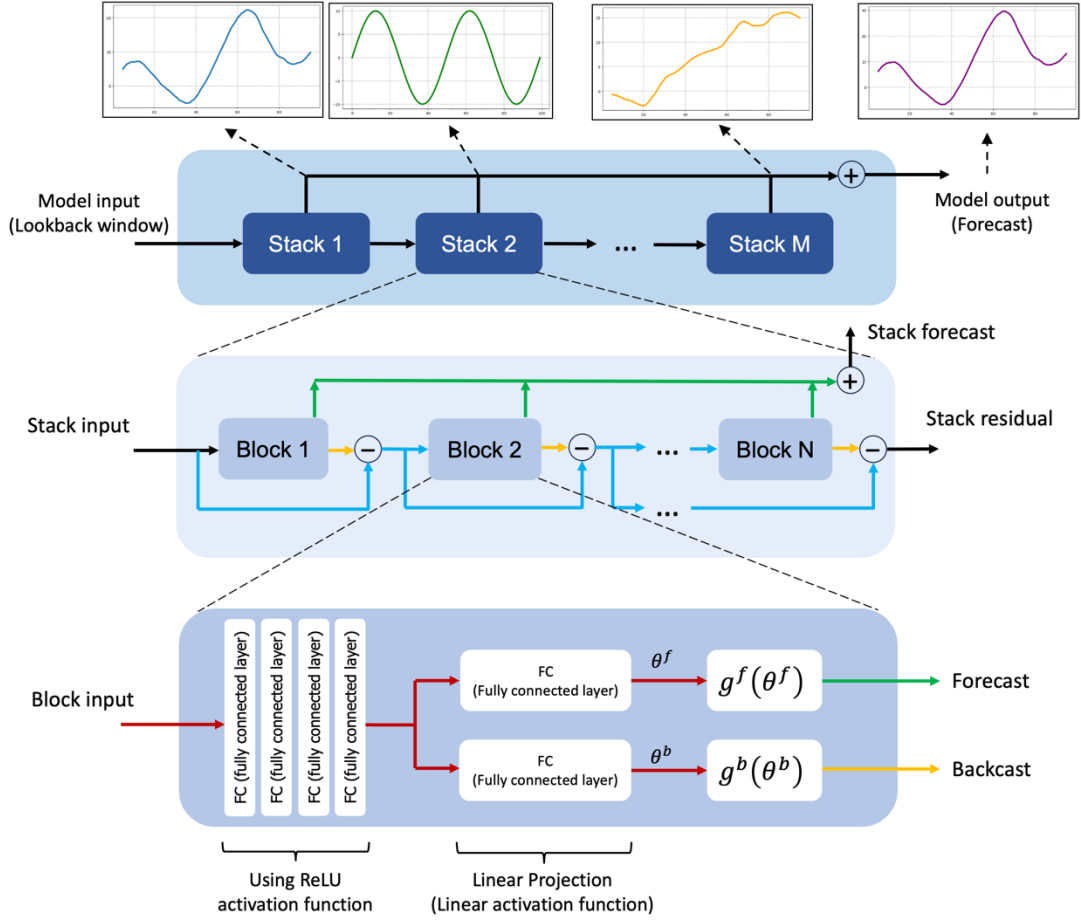
308
309                           Figure 3. The N-BEATS modeling structure used in this research.

310    **2.2.3. N-HiTS**

311    N-HiTS builds upon the N-BEATS architecture but with improved accuracy and computational efficiency

312    for long-horizon forecasting. N-HiTS utilizes multi-rate sampling and multi-scale synthesis of forecasts,

313    leading to a hierarchical forecast structure that lowers computational demands and improves prediction

314    accuracy (Challu et al., 2022).

315    Like N-BEATS, N-HiTS employs local nonlinear mappings onto foundational functions within numerous

316    blocks. Each block includes an MLP that generates backcast and forecast output coefficients. The backcast

317    output refines the input data for the following blocks, and the forecast outputs are combined to generate the

318    final prediction. Blocks are organized into stacks, with each stack dedicated to grasping specific data

319    attributes using its own distinct set of functions. The network's input is a sequence of $L$ lags (look-back

320    period), with $S$ stacks, each containing $B$ blocks (Challu et al., 2022).

321     In each block, a *MaxPool* layer with varying kernel sizes ($k_l$) is employed at the input, enabling the block

322     to focus on specific input components of different scales. Larger kernel sizes emphasize the analysis of

323     larger-scale, low-frequency data, aiding in improving long-term forecasting accuracy. This approach,

324     known as multi-rate signal sampling, alters the effective input signal sampling rate for each block's MLP

325     (Challu et al., 2022).

326     Additionally, multi-rate processing has several advantages. It reduces memory usage, computational

327     demands, the number of learnable parameters, and helps prevent overfitting, while preserving the original

328     receptive field. The following operation is applicable to the input $y_{t-L:t,l}$ of each block, with the first block

329     ($l = 1$) using the network-wide input, where $y_{t-L:t,1} \equiv y_{t-L:t}$.

$$y_{t-L:t,l} = MaxPool\,(y_{t-L:t,l}, k_l) \qquad \text{(Equation 15)}$$

330     In many multi-horizon forecasting models, the number of neural network predictions matches the horizon's

331     dimensionality, denoted as $H$. For instance, in N-BEATS, the number of predictions $\left|\theta_l^f\right| = H$. This results

332     in a significant increase in computational demands and an unnecessary surge in model complexity as the

333     horizon $H$ becomes larger (Challu et al., 2022).

334     To address these challenges, N-HiTS proposes the use of temporal interpolation. This model manages the

335     parameter counts per unit of output time ($\left|\theta_l^f\right| = \lceil r_l\, H \rceil$) by defining the dimensionality of the interpolation

336     coefficients with respect to the expressiveness ratio $r_l$. To revert to the original sampling rate and predict

337     all horizon points, this model employs temporal interpolation through the function $g$:

$$\hat{y}_{\tau,l} = g(\tau, \theta_l^f), \qquad \forall \tau \in \{t + 1,\dots,t + H\}, \qquad \text{(Equation 16)}$$

$$\tilde{y}_{\tau,l} = g(\tau, \theta_l^b), \qquad \forall \tau \in \{t - L,\dots,t\}, \qquad \text{(Equation 17)}$$

$$g(\tau,\theta) = \theta[t_1] + \left(\frac{\theta[t_2] - \theta[t_1]}{t_2 - t_1}\right)(\tau - t_1) \qquad \text{(Equation 18)}$$

$$t_1 = \arg\min_{t \in \tau:t \leq \tau} \tau - t, \qquad t_2 = t_1 + 1/r_l \qquad \text{(Equation 19)}$$

338     The hierarchical interpolation approach involves distributing expressiveness ratios over blocks, integrated

339     with multi-rate sampling. Blocks closer to the input employ more aggressive interpolation, generating lower

340     granularity signals. These blocks specialize in analyzing more aggressively subsampled signals. The final

341     hierarchical prediction, $\hat{y}_{t+1:t+H}$, is constructed by combining outputs from all blocks, creating

342  interpolations at various time-scale hierarchy levels. This approach maintains a structured hierarchy of
343  interpolation granularity, with each block focusing on its own input and output scales (Challu et al., 2022).

344  To manage a diverse set of frequency bands while maintaining control over the number of parameters,
345  exponentially increasing expressiveness ratios are recommended. As an alternative, each stack can be
346  dedicated to modeling various recognizable cycles within the time series (e.g., weekly, or daily) employing
347  matching $r_l$. Ultimately, the residual obtained from backcasting in the preceding hierarchy level is
348  subtracted from the input of the subsequent level, intensifying the next-level block's attention on signals
349  outside the previously addressed band (Challu et al., 2022).

$$\hat{y}_{t+1:t+H} = \sum_{l=1}^{L} \hat{y}_{t+1:t+H,l} \qquad \text{(Equation 20)}$$

$$y_{t-L:t,l+1} = y_{t-L:t,l} - \tilde{y}_{t-L:t,l} \qquad \text{(Equation 21)}$$
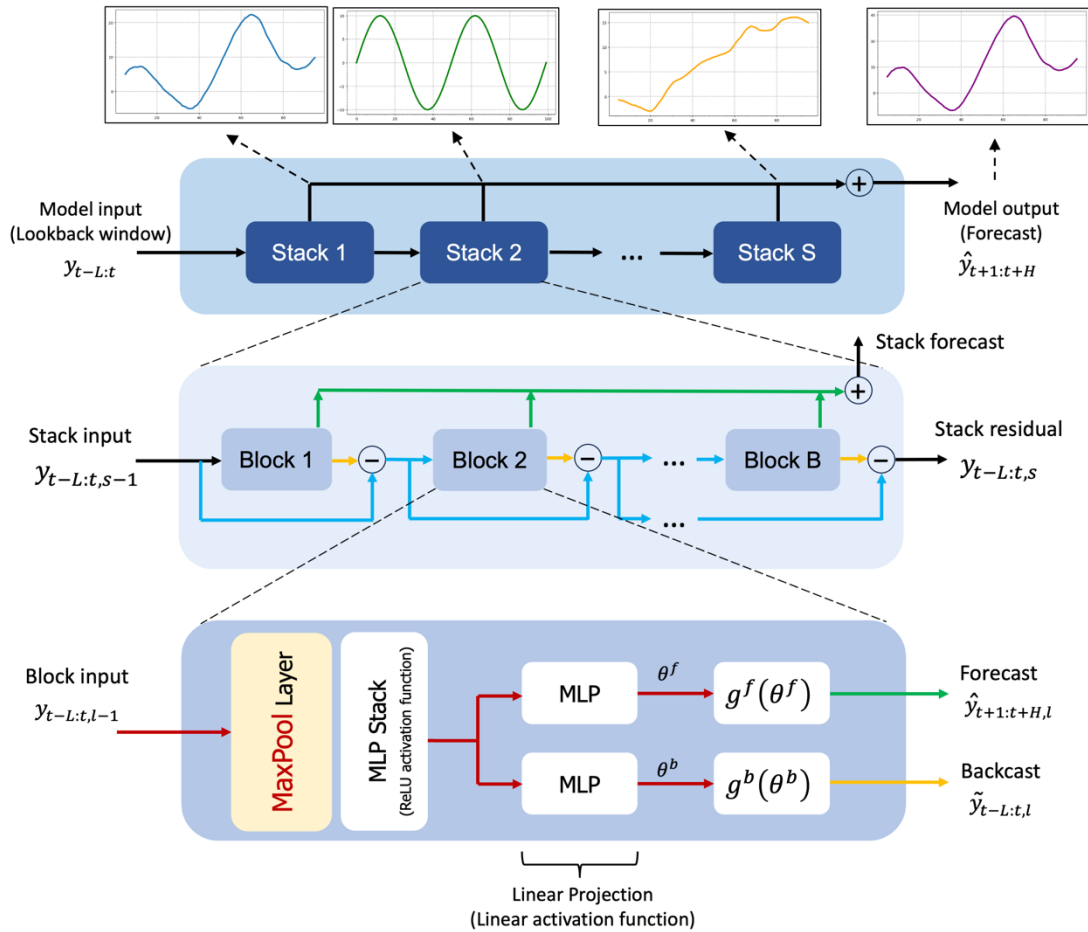


350

351 Figure 4. The structure of N-HiTS model programmed in this study. The architecture includes several
352 Stacks, each Stack includes several Block, where each block consists of a MaxPool layer and a multi-
353 layer which learn to produce coefficients for the backcast and forecast outputs of its basis.

354 **2.3. Performance Metrics**

355 To comprehensively evaluate the accuracy of flood predictions, we utilized a suite of metrics, including

356 Nash-Sutcliffe Efficiency (NSE; Nash and Sutcliffe, 1970), persistent Nash-Sutcliffe Efficiency (persistent-

357 NSE), Kling–Gupta efficiency (KGE; Gupta et al. 2009), Root Mean Square Error (RMSE), Mean

358 Absolute Error (MAE), Peak Flow Error (PFE), and Time to Peak Error (TPE; Evin et al., 2023; Lobligeois

359 et al., 2014). These metrics collectively facilitate a rigorous assessment of the model's performance in

360 reproducing the magnitude of observed peak flows and the shape of the hydrograph.

361 NSE measures the model's ability to explain the variance in observed data and assesses the goodness-of-fit

362 by comparing the observed and simulated hydrographs. In hydrological studies, the NSE index is a widely

363 accepted measure for evaluating the fitting quality of models (McCuen et al., 2006). It is calculated as:

$$NSE = 1 - \frac{\sum_{i=1}^{n}\left(Q_{s_i} - Q_{o_i}\right)^2}{\sum_{i=1}^{n}\left(Q_{o_i} - \overline{Q_o}\right)^2} \qquad \text{(Equation 22)}$$

364 Where $Q_{o_i}$ represents observed value at time $i$, $Q_{s_i}$ represents simulated value at time $i$, $\overline{Q_o}$ is the mean

365 observed values and $n$ is the number of data points. An NSE value of 1 indicates a perfect match between

366 the observed and modeled data, while lower values represent the degree of departure from a perfect fit.

367 As the models are designed to predict one hour ahead, the persistent-NSE is essential for evaluating their

368 performance. The standard NSE measures the model's sum of squared errors relative to the sum of squared

369 errors when the mean observation is used as the forecast value. In contrast, persistent-NSE uses the most

370 recent observed data as the forecast value for comparison (Nevo et al., 2022). The persistent-NSE is

371 calculated as:

$$persistent - NSE = 1 - \frac{\sum_{i=1}^{n}\left(Q_{s_i} - Q_{o_i}\right)^2}{\sum_{i=1}^{n}\left(Q_{o_i} - Q_{o_{i-1}}\right)^2} \qquad \text{(Equation 23)}$$

372 Where $Q_{o_i}$ represents the observed value at time $i$, $Q_{s_i}$ represents the simulated value at time $i$, $Q_{o_{i-1}}$ is the

373 observed value at the last time step $(i - 1)$ and $n$ is the number of data points.

16

374 The KGE is a widely used performance metric in hydrological modeling and combines multiple aspects of

375 model performance, including correlation, variability bias, and mean bias. The KGE metric is calculated

376 using the following equation:

$$KGE = 1 - \sqrt{(r-1)^2 + (\alpha-1)^2 + (\beta-1)^2}$$ (Equation 24)

377 Where $r$ represents Pearson correlation coefficient between observed $Q_o$ and simulated $Q_s$ values.

378 $\alpha$ represents bias ratio, calculated as $\alpha = \frac{\mu_s}{\mu_o}$ where $\mu_s$ and $\mu_o$ are the means of simulated and observed data,

379 respectively. $\beta$ represents variability ratio, calculated as $\beta = \frac{\sigma_s/\mu_s}{\sigma_o/\mu_o}$ where $\sigma_s$ and $\sigma_o$ are the standard

380 deviations of simulated and observed data, respectively.

381 RMSE quantifies the average magnitude of errors between observed and modeled values, offering insights

382 into the absolute goodness-of-fit, while MAE is a measure of the average absolute difference between the

383 modeled values and the observed values and provides a measure of the average magnitude of errors. RMSE

384 is calculated as:

$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(Q_{o_i} - Q_{s_i})^2}$$ (Equation 25)

385 and MAE is calculated as:

$$MAE = \frac{1}{n}\sum_{i=1}^{n}|Q_{o_i} - Q_{s_i}|$$ (Equation 26)

386 Where $Q_{o_i}$ represents observed value at time $i$, $Q_{s_i}$ represents simulated value at time $i$, and $n$ is the number

387 of data points. RMSE and MAE provide information about the magnitude of modeling errors, with smaller

388 values indicating a better model fit.

389 PFE quantifies the magnitude disparity between observed and modeled peak flow values. The PFE metric

390 is defined as:

$$PFE = \frac{|Q_{o\,max} - Q_{s\,max}|}{Q_{o\,max}}$$ (Equation 27)

391    Where $Q_{o\,max}$ represents the observed peak flow value, and $Q_{s\,max}$ signifies the simulated peak flow value.

392    The PFE metric, expressed as a dimensionless value, provides a quantitative measure of the relative error

393    in predicting peak flow magnitudes concerning the observed values. A smaller PFE denotes more accurate

394    modeling of peak flow magnitudes, with a value of zero indicating a perfect match.

395    TPE assesses the temporal alignment of peak flows in the observed and modeled hydrographs. The TPE

396    metric is computed as:

$$TPE = \left| T_{o\,max} - T_{s\,max} \right| \qquad \text{(Equation 28)}$$

397    Where $T_{o\,max}$ signifies the time at which the peak flow occurs in the observed hydrograph, and $T_{s\,max}$

398    represents the time at which the peak flow occurs in the simulated hydrograph. TPE that is measured in

399    units of time (hours), provides insight into the precision of peak flow timing. Smaller TPE values indicate

400    a superior alignment between the observed and modeled peak flow timing, while larger TPE values indicate

401    discrepancies in the temporal occurrence of peak flows.

402    The utilization of these five metrics, PFE, persistent-NSE, TPE, NSE, and RMSE, collectively provides a

403    robust and multifaceted assessment of flood prediction performance. This approach ensures that both the

404    magnitude and timing of peak flows, as well as the overall hydrograph shape, are accurately calibrated and

405    validated.

406    **2.4. Sensitivity and Uncertainty Analysis**

407    When implementing NN models, it's crucial to understand how each parameter affects the model's

408    performance or outputs. To achieve this, we systematically excluded each parameter from the model one

409    by one (the Leave-One-Out method). For each exclusion, we retrained the model without that specific

410    parameter and then tested its performance against a test dataset. This method helps in understanding which

411    parameters are most critical to the model's performance and which ones have a lesser impact. It also allows

412    us to identify any parameters that may be redundant or have little effect on the overall outcome, thus

413    potentially simplifying the model without sacrificing accuracy.

414    In this study, we utilized probabilistic approaches to quantify the uncertainty in flood prediction. This

415    method is rooted in statistical techniques employed for the estimation of unknown probability distributions,

416    with a foundation in observed data. More specifically, we leveraged the Maximum Likelihood Estimation

417    (MLE) approach, which entails the determination of parameter values that optimize the likelihood function.

18

418    The likelihood function quantifies the probability of parameters taking particular values, given the observed
419    realizations.

420    We incorporated the MQL as a probabilistic error metric into the algorithmic architecture. MQL performs
421    an evaluation by computing the average loss for a predefined set of quantiles. This computation is grounded
422    in the absolute disparities between predicted quantiles and their corresponding observed values. By
423    considering multiple quantile levels, MQL provides a comprehensive assessment of the model's ability to
424    capture the distribution of the target variable, rather than focusing solely on point estimates.

425    The MQL metric also aligns closely with the Continuous Ranked Probability Score (CRPS), a standard tool
426    for evaluating predictive distributions. CRPS measures the difference between the predicted cumulative
427    distribution function and the observed values by integrating over all possible quantiles. The computation of
428    CRPS involves a numerical integration technique that discretizes quantiles and applies a left Riemann
429    approximation for CRPS integral computation. This process culminates in the averaging of these
430    computations over uniformly spaced quantiles, providing a robust evaluation of the predictive distribution
431    $\hat{F}_t$.

432    To calculate the 95th PPU, we utilized the 0.95 quantile level within the MQL. This quantile level directly
433    corresponds to the 95th percentile of the predicted distribution, providing an estimate of the 95% confidence
434    interval. By examining the model's performance at this specific quantile, we effectively assessed its ability
435    to accurately capture the predicted values with 95% confidence.

436    Incorporating MQL as a central metric in our study underscores its suitability for probabilistic forecasting,
437    particularly in the context of uncertainty quantification. Unlike traditional error metrics that focus on point
438    predictions, MQL captures both central tendencies and variability by penalizing errors symmetrically across
439    quantiles. This property ensures balanced and reliable assessments of the predictive distribution, ultimately
440    enhancing the robustness and interpretability of flood prediction models.

$$\text{MQL}\left(Q_\tau, \left[\hat{Q}_\tau^{q_1}, \dots, \hat{Q}_\tau^{q_i}\right]\right) = \frac{1}{n} \sum_{q_i} \text{QL}\left(Q_\tau, \hat{Q}_\tau^{q_i}\right) \qquad \text{(Equation 29)}$$

$$\text{CRPS}\left(Q_\tau, \hat{F}_\tau\right) = \int_0^1 \text{QL}\left(Q_\tau, \hat{Q}_\tau^{q_i}\right) dq \qquad \text{(Equation 30)}$$

$$\text{QL}\left(Q_\tau, \hat{Q}_\tau^q\right) = \frac{1}{H} \sum_{\tau=t+1}^{t+H} \left((1-q)\left(\hat{Q}_\tau^q - Q_\tau\right) + q\left(Q_\tau - \hat{Q}_\tau^q\right)\right) \qquad \text{(Equation 31)}$$

441    Where $Q_\tau$ represents observed value at time $\tau$, $\hat{Q}_\tau^q$ represents simulated value at time $\tau$, $q$ is the slope of the

442    quantile loss, and $H$ is the horizon of forecasting.
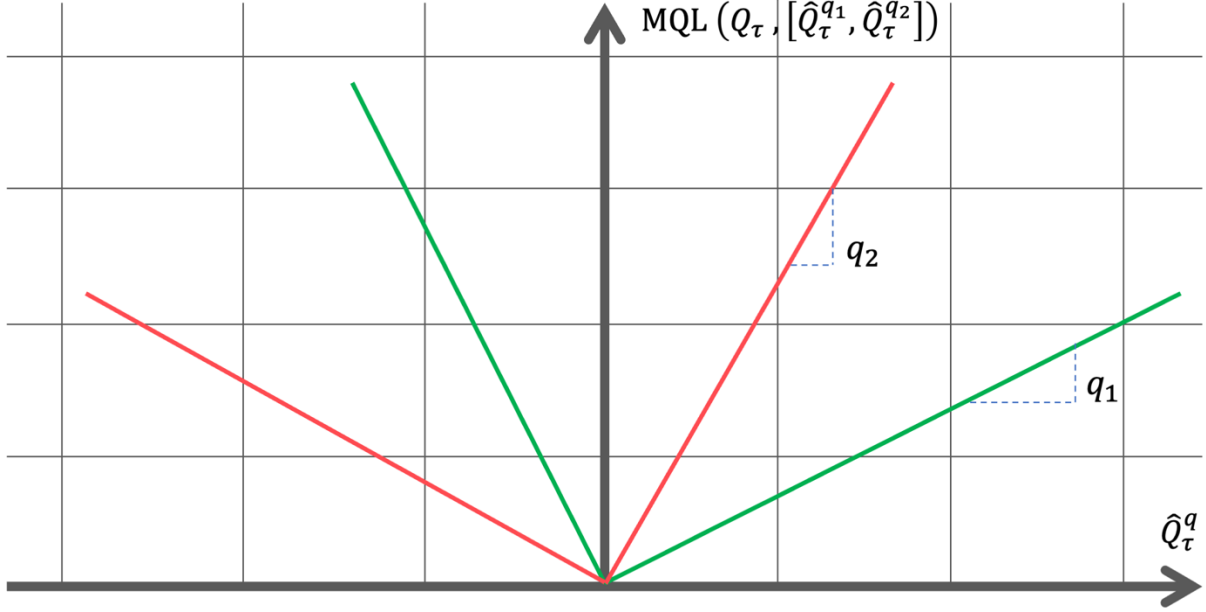


443

444    Figure 5. The MQL function which shows loss values for different parameters of $q$ when the true value is
445    $Q_\tau$.

446    Furthermore, we employed two key indices, the R-Factor and the P-Ffactor, to rigorously assess the quality

447    of uncertainty performance in our hydrological modeling. These metrics are instrumental in quantifying the

448    extent to which the model's predictions encompass the observed data, thereby providing valuable insights

449    into the model's predictive accuracy and reliability.

450    The P-Factor, or percentage of data within a 95PPU, is the first index used in this assessment. The P-Factor

451    quantifies the percentage of observed data that falls within the 95PPU, providing a measure of the model's

452    predictive accuracy. The P-Factor can theoretically vary from 0% to a maximum of 100%. A P-Factor of

453    100% signifies a perfect alignment between the model's predictions and the observed data within the

454    uncertainty band. In contrast, a lower P-Factor indicates a reduced ability of the model to predict data within

455    the specified uncertainty range.

$$P - Factor = \frac{Observations\ braketed\ by\ 95PPU}{Number\ of\ observations} \times 100 \qquad \text{(Equation 32)}$$
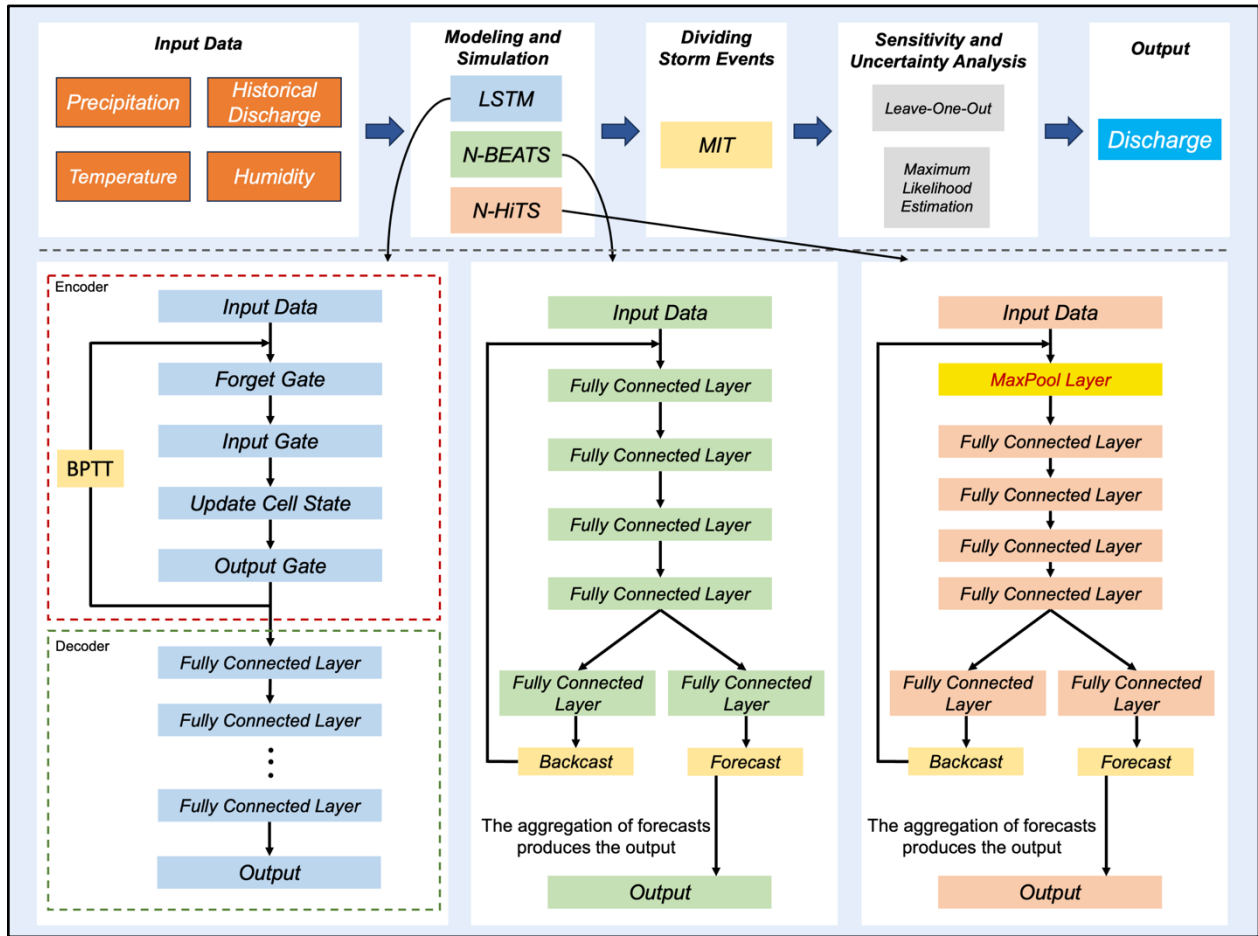
456    The R-Factor can be computed by dividing the average width of the uncertainty band by the standard

457    deviation of the measured variable. The R-Factor, with a minimum possible value of zero, provides a

458    measure of the spread of the uncertainty relative to the variability of the observed data. Theoretically, the

459　R-Factor spans from 0 to infinity, and a value of zero implies that the model's predictions precisely match

460　the measured data, with the uncertainty band being very narrow in relation to the variability of the observed

461　data.

$$R - Factor = \frac{Average\ width\ of\ 95PPU\ band}{Standard\ deviation\ of\ measured\ variables} \times 100 \quad \text{(Equation 33)}$$

462　In practice, the quality of the model is assessed by considering the 95% prediction band with the highest P-

463　Factor and the lowest R-Factor. This specific band encompasses the majority of observed records,

464　signifying the model's ability to provide accurate and reliable predictions while effectively quantifying

465　uncertainty. A simulation with a P-Factor of 1 and an R-Factor of 0 signifies an ideal scenario where the

466　model precisely matches the measured data within the uncertainty band (Abbaspour et al., 2007).

467　Figure 6 shows the workflow of programming N-BEATS, N-HiTS, and LSTM for flood prediction. As

468　illustrated, the initial step involved cleaning and preparing the input data, which was then used to feed the

469　models. The workflow for each model and their output generation processes are depicted in Figure 6. We

470　segmented the storm events using the MIT approach, as previously described. Following this, we conducted

471　a sensitivity analysis using the Leave-One-Out method and performed uncertainty analysis using the MLE

472　approach to construct the 95PPU band. This rigorous methodology ensures a robust evaluation of model

473　performance under varying conditions and highlights the models' predictive reliability and resilience. We

474　employed the "NeuralForecast" Python package to develop the N-BEATS, N-HiTS, and LSTM models.

475　This package provides a diverse array of NN models with an emphasis on usability and robustness.

476

477



478

Figure 6. The workflow of N-BEATS, N-HiTS, and LSTM implementation. The upper section of the figure illustrates multiple steps from data preprocessing to model evaluation. The lower section provides a detailed view of the workflow and implementation for each model, highlighting the specific processes and methodologies employed in generating the outputs. Backpropagation Through Time (BPTT) trains LSTM by unrolling the model through time, computing gradients for each time step, and updating weights based on temporal dependencies.

## 3. Results and Discussion

### 3.1. Independent Storms Delineation

MIT's contextual delineation of storm events laid the groundwork for in-depth evaluation of rainfall events, enabling isolation and separation of rainfall events that led to significant flooding events. The nuanced outcomes of the MIT assessment contributed significantly to the understanding of rainfall variability and distribution as the dominant contributor to flood generation.

491  During modeling implementation, the initial imperative was the precise distinction of storm events within
492  the precipitation time series data of each case study. Our findings demonstrate that on average a dry period
493  of 7 hours serves as the optimal MIT time for both of our case studies. This outcome signifies that when a
494  dry interval of more than 7 hours transpires between two successive rainfall events, these subsequent
495  rainfalls should be considered two distinct storm events. This determination underlines the temporal
496  threshold necessary for distinguishing between individual meteorological phenomena in two case studies.

497  **3.2. Hyperparameter Optimization**

498  In the context of hyperparameter optimization, we systematically considered and tuned various
499  hyperparameters for the N-HiTS, N-BEATS, and LSTM. Following extensive exploration and fine-tuning
500  of these hyperparameters, the optimal configurations were identified (see Table 2). For the N-HiTS model,
501  the most favorable outcomes were achieved with the following hyperparameter settings: 2000 epochs,
502  "identity" for scaler type, a learning rate of 0.001, a batch size of 32, input size of 24 hours, "identity" for
503  stack type, 512 units for hidden layers of each stack, step size of 1, MQLoss as loss function, and "ReLU"
504  for the activation function. As shown in Table 2, the N-HiTS model demonstrated superior performance
505  with 4 stacks, containing 2 blocks each, and corresponding coefficients of 48, 24, 12, and 1, showcasing
506  the significance of these settings for flood prediction.

507  This hyperparameter optimization was also conducted for the N-BEATS model. In this model, we
508  considered 2000 epochs, 3 stacks with 2 blocks, "identity" for scaler type, a learning rate of 0.001, a batch
509  size of 32, input size of 24 hours, "identity" for stack type, 512 units for hidden layers of each stack, step
510  size of 1, MQLoss as loss function, and "ReLU" for the activation function.

511  Moreover, the LSTM as a benchmark model yielded its best results with 5000 epochs, an input size of 24
512  hours, "identity" as the scaler type, a learning rate of 0.001, a batch size of 32, and "tanh" as the activation
513  function. Furthermore, the LSTM's hidden state was most effective with two layers containing 128 units,
514  and the MLP decoder thrived with two layers encompassing 128 units. These meticulously optimized
515  hyperparameter settings represent the culmination of efforts to ensure that each model operates at its peak
516  potential, facilitating accurate flood prediction.

517  Table 2. Optimized values for the hyperparameters.

| Hyperparameter | N-HiTS | N-BEATS | LSTM |
|---|---|---|---|
| Epoch | 2000 | 2000 | 5000 |
| Scaler type | identity | identity | standard |

| | | | |
|---|---|---|---|
| Learning rate | 0.001 | 0.001 | 0.001 |
| Batch size | 32 | 32 | 32 |
| Input size | 24 hours | 24 hours | 24 hours |
| Stack type | Seasonality, trend, identity, identity | Seasonality, trend, identity | * |
| Number of units in each hidden layer | 512 | 512 | 128 |
| Loss function | MQLoss | MQLoss | MQLoss |
| Activation function | ReLU | ReLU | tanh |
| Number of stacks | 4 | 3 | * |
| Number of blocks in each stack | 2 | 2 | * |
| Stacks' coefficients | 48,24,12,1 | * | * |

*Not applicable

In Table 2, "epoch" refers to the number of training steps, and "scaler type" indicates the type of scaler used for normalizing temporal inputs. The "learning rate" specifies the step size at each iteration while optimizing the model, and the "batch size" represents the number of samples processed in one forward and backward pass. The "loss function" quantifies the difference between the predicted outputs and the actual target values, while the "activation function" determines whether a neuron should be activated. The "stacks' coefficients" in the N-HITS model control the frequency specialization for each stack, enabling effective handling of different frequency components in the time series data.

Another hyperparameter for all three models is input size, which is a parameter that determines the maximum sequence length for truncated backpropagation during training and the number of autoregressive inputs (lags) that the models considered for prediction. Essentially, input size represents the length of the historical series data used as input to the model. This parameter offers flexibility in the models, allowing them to learn from a defined window of past observations, which can range from the entire historical dataset to a subset, tailored to the specific requirements of the prediction task. In the context of flood prediction, determining the appropriate input size is crucial to adequately capture the meteorological data preceding the flood event. To address this, we calculated the time of concentration ($TC$) of the watershed system and set the input size to exceed this duration. According to the Natural Resources Conservation Service (NRCS), for typical natural watershed conditions, the TC can be calculated from lag time, the time between peak rainfall and peak discharge, using the formula: $Lag\ time = TC \times 0.6$ (NRCS, 2009). Specifically, the

537      average *TC* in the Lower Dog River watershed and Upper Dutchmans Creek watershed was found to be 19

538      and 22 hours, respectively. As these represent the average *TC* for our case studies, we selected the 24 hours

539      for input data, slightly longer than the calculated avaerage *TC*, ensuring sufficient coverage of relevant

540      meteorological data preceding all flood events.

**3.3. Flood Prediction and Performance Assessment**

542      In this study, we conducted a comprehensive performance evaluation of N-HiTS, N-BEATS, and

543      benchmarked these models with LSTM, utilizing two case studies: the Lower Dog River and the Upper

544      Dutchmans Creek watersheds. Within these case studies, we trained and validated the models separately

545      for each watershed across a diverse set of storm events from 01/10/2007 to 01/10/2022 (15 years) in the

546      Lower Dog River and from 21/12/1994 to 01/10/2022 (27 years) in the Upper Dutchmans Creek. The

547      decision to train separate models for each catchment was made to account for the unique hydrological

548      characteristics and local features specific to each watershed. By training models individually, we aimed to

549      optimize performance by tailoring each model to the distinct rainfall-runoff relationship inherent in each

550      catchment. All algorithms were tested using unseen flooding events that occurred between 14/12/2022 and

551      28/03/2023. In the Dog River gauging station, two winter storms i.e., January 3rd to January 5th, 2023

552      (Event 1) and February 17th to February 18th, 2023 (Event 2), as well as a spring flood event that occurred

553      during March 26th to March 28th, 2023 (Event 3) were selected for testing. Additionally, three winter

554      flooding events, i.e., December 14th to December 16th, 2022 (Event 4), January 25th and January 26th,

555      2023 (Event 5), and February 11th to February 13th, 2023 (Event 6), were chosen to test the algorithms

556      across the Killian Creek gauging station in the Upper Dutchmans Creek. The rainfall events corresponding

557      to these flooding events were delineated using the MIT technique discussed in Section 3.1.

558      Our results for the Lower Dog River case study, explicitly demonstrated the accuracy of both N-HiTS and

559      N-BEATS in generating the winter and spring flood hydrographs compared to the LSTM model across all

560      selected storm events. Although, N-HiTS prediction slightly outperformed N-BEATS during winter

561      prediction (January 3rd to January 5th, 2023). In this event, N-HiTS outperformed N-BEATS with a

562      difference of 11.6% in MAE and 20% in RMSE. The N-HiTS slight outperformance (see Tables 3 and 4)

563      is attributed to its unique structure that allows the model to discern and capture intricate patterns within the

564      data. Specifically, N-HiTS predicted flooding events hierarchically using blocks specialized in different

565      rainfall frequencies based on controlled signal projections, through expressiveness ratios, and interpolation

566      of each block. The coefficients are then used to synthesize backcast through

567      $\tilde{y}_t - L: t, l$ and forecast $(\tilde{y}_{t+1}: t + H, l)$ outputs of the block as a flood value. The coefficients were locally

568      determined along the horizon, allowing N-HiTS to reconstruct nonstationary signals over time.

569    While the N-HiTS emerged as the most accurate in predicting flood hydrograph among the three models,
570    its performance was somehow comparable with N-BEATS. The N-BEATS model exhibited good
571    performance in two case studies. It consistently provided competitive results, demonstrating its capacity to
572    effectively handle diverse storm events and deliver reliable predictions. N-BEATS has a generic and
573    interpretable architecture depending on the blocks it uses. Interpretable configuration sequentially projects
574    the signal into polynomials and harmonic basis to learn trend and seasonality components while generic
575    configuration substitutes the polynomial and harmonic basis for identity basis and larger network's depth.
576    In this study, we used interpretable architecture, as it regularizes its predictions through projections into
577    harmonic and trend basis that is well-suited for flood prediction tasks. Using interpretable architecture,
578    flood prediction was aggregated in a hierarchical fashion. This enabled the building of a very deep neural
579    network with interpretable flood prediction outputs.

580    It is essential to underscore that, despite its strong performance, the N-BEATS model did not surpass the
581    N-HiTS model in terms of NSE, Persistent-NSE, MAE, and RMSE for the Lower Dog River case study.
582    Although both models showed almost the same KGE values. Notably, the N-BEATS model showcased
583    superior results based on the PFE metric, signifying its exceptional capability in accurately predicting flood
584    peaks. However, both N-HiTS and N-BEATS models overestimated the flood peak rate of Event 2 for the
585    Lower Dog River watershed. This event, which occurred from February 17$^{th}$ to February 18$^{th}$, 2023, was
586    flashy, short, and intense proceeded by a prior small rainfall event (from February 12$^{th}$ until February 13$^{th}$)
587    that minimized the rate of infiltration. This flash flood event caused by excessive rainfall in a short period
588    of time (<8 hours) was challenging to predict for N-BEATS and N-HiTS models. In addition, predicting
589    the magnitude of changes in the recession curve of the third event seems to be a challenge for both models.
590    The specific part of the flood hydrograph after the precipitation event, where flood diminishes during a
591    rainless is dominated by the release of runoff from shallow aquifer systems or natural storages. It seems
592    both models showed a slight deficiency in capturing this portion of the hydrograph when the rainfall amount
593    decreases over time in the Dog River gauging station.

594    Conversely, in the Killian Creek gauging station, the N-BEATS model almost emerged as the top performer
595    in predicting the flood hydrograph based on NSE, Persistent-NSE, RMSE, and PFE performance metrics
596    (see Tables 3 and 4).  KGE values remained almost the same for both models. In addition, both N-BEATS
597    and N-HiTS slightly overpredicted time to peak values for Event 5. This reflects the fact that when rainfall
598    value varies randomly around zero, it provides less to no information for the algorithms to learn the
599    fluctuations and patterns in time series data. Both N-HiTS and N-BEATS provided comparable results for
600    all events predicted in this study. N-HiTS builds upon N-BEATS by adding a MaxPool layer at each block.
601    Each block consists of an MLP layer that learns to produce coefficients for the backcast and forecast

602    outputs. This subsamples the time series and allows each stack to focus on either short-term or long-term

603    effects, depending on the pooling kernel size. Then, the partial predictions of each stack are combined using

604    hierarchical interpolation. This ability enhances N-HiTS capabilities to produce drastically improved,

605    interpretable, and computationally efficient long-horizon flood predictions.


606    In contrast, the performance of LSTM as a benchmark model lagged behind both N-HiTS and N-BEATS

607    models for all events across two case studies. Despite its extensive applications in various hydrology

608    domains, the LSTM model exhibited comparatively lower accuracy when tasked with predicting flood

609    responses during different storm events. Focusing on NSE, Persistent-NSE. KGE, MAE, RMSE, and PFE

610    metrics, it is noteworthy that all three models, across both case studies, consistently succeeded in capturing

611    peak flow rates at the appropriate timing. All models demonstrated commendable results with respect to

612    the TPE metric. In most scenarios, TPE revealed a value of 0, signifying that the models accurately

613    pinpointed the peak flow rate precisely at the expected time. In some instances, TPE reached a value of 1,

614    showing a deviation of one hour in predicting the peak flow time. This deviation is deemed acceptable,

615    particularly considering the utilization of short, intense rainfall for our analysis.


616    Our investigation into the performance of the three distinct forecasting models yielded compelling results

617    pertaining to their ability to generate 95PPU, as quantified by the P-Factor and R-Factor. These factors

618    serve as critical indicators for assessing the reliability and precision of the uncertainty bands produced by

619    the MLE. Our findings demonstrated that the N-HiTS and N-BEATS models outperformed the LSTM

620    model in mathematically defining uncertainty bands, in terms of R-Factor metric. The R-Factor, a crucial

621    metric for evaluating the average width of the uncertainty band, consistently favored the N-HiTS and N-

622    BEATS models over their counterparts. This finding was consistent across a diverse range of storm events.

623    In addition, coupling MLE with the N-HiTS and N-BEATS models demonstrated superior performance in

624    generating 95PPU when assessed through the P-Factor metric. The P-Factor represents another vital aspect

625    of uncertainty quantification, focusing on the precision of the uncertainty bands.


626    Figures 8 and 9 present graphical depictions of the predicted flood with uncertainty assessment for each

627    model as well as Flow Duration Curve (FDC) across two gauging stations.  As illustrated, the uncertainty

628    bands skillfully bracketed most of the observational data, reflecting the fact that MLE was successful in

629    reducing errors in flood prediction. FDC analysis also revealed that N-HiTS and N-BEATS models

630    skillfully predicted the flood hydrograph, however, both models were particularly successful in predicting

631    moderate to high flood events (1800-6000 and >6000 cfs). In the FDC plots, the x-axis denotes the

632    exceedance probability, expressed as a percentage, while the y-axis signifies flood in cubic feet per second.

633    Notably, these plots reveal distinctive patterns in the performance of the N-HiTS, N-BEATS, and LSTM

634    models. Within the lower exceedance probability range, particularly around the peak flow, the N-HiTS and
635    N-BEATS models demonstrated a clear superiority over the LSTM model, closely aligning with the
636    observed data. This observed trend is consistent when examining the corresponding hydrographs. Across
637    all events, the flood hydrographs generated by N-HiTS and N-BEATS exhibited a closer resemblance to
638    the observed data, particularly in the vicinity of the peak timing and rate, compared to the hydrographs
639    produced by the LSTM model. These findings underscore the enhanced predictive accuracy and reliability
640    of the N-HiTS and N-BEATS models, particularly in predicting moderate to high flood events as well as
641    critical hydrograph features such as peak flow rate and timing. The alignment of model-generated FDCs
642    and hydrographs with observed data in the proximity of peak flow further establishes the efficiency of N-
643    HiTS and N-BEATS in accurately reproducing the dynamics of flood generation mechanisms across two
644    headwater streams.

645

646                    Table 3. The performance metrics for the Lower Dog River flood predictions.

| Model | Performance Metric | Event 1 | Event 2 | Event 3 |
|---|---|---|---|---|
| N-HiTS | NSE | 0.995 | 0.991 | 0.992 |
| | Persistent-NSE | 0.947 | 0.931 | 0.948 |
| | KGE | 0.977 | 0.989 | 0.976 |
| | RMSE | 123.2 | 27.6 | 68.5 |
| | MAE | 64.1 | 12.0 | 37.8 |
| | PFE | 0.018 | 0.051 | 0.015 |
| | TPE (hours) | 0 | 1 | 0 |
| | P-Factor | 96.9 % | 100 % | 93.5 % |
| | R-Factor | 0.27 | 0.40 | 0.33 |
| N-BEATS | NSE | 0.991 | 0.989 | 0.993 |
| | Persistent-NSE | 0.917 | 0.916 | 0.956 |
| | KGE | 0.984 | 0.984 | 0.98 |
| | RMSE | 154.1 | 30.5 | 62.5 |
| | MAE | 72.6 | 13.6 | 35.9 |
| | PFE | 0.0005 | 0.031 | 0.0002 |
| | TPE (hours) | 0 | 1 | 0 |
| | P-Factor | 87.8 % | 100 % | 90.3 % |
| | R-Factor | 0.17 | 0.23 | 0.24 |
| LSTM | NSE | 0.756 | 0.983 | 0.988 |

| | | | |
|---|---|---|---|
| Persistent-NSE | -1.44 | 0.871 | 0.929 |
| KGE | 0.765 | 0.978 | 0.971 |
| RMSE | 841.1 | 37.9 | 79.5 |
| MAE | 369.4 | 18.6 | 42 |
| PFE | 0.258 | 0.036 | 0.016 |
| TPE (hours) | 1 | 0 | 0 |
| P-Factor | 81.8 % | 93.1 % | 96.7 % |
| R-Factor | 0.37 | 0.51 | 0.6 |

647

648 Table 4. The performance metrics for the Killian Creek flood predictions.

| Model | Performance Metric | Event 4 | Event 5 | Event 6 |
|---|---|---|---|---|
| | NSE | 0.991 | 0.971 | 0.991 |
| | Persistent-NSE | 0.885 | 0.806 | 0.844 |
| | KGE | 0.982 | 0.967 | 0.991 |
| | RMSE | 28.8 | 46.0 | 19.0 |
| N-HiTS | MAE | 17.9 | 23.8 | 11.5 |
| | PFE | 0.017 | 0.008 | 0.020 |
| | TPE (hours) | 0 | 0 | 0 |
| | P-Factor | 92.6 % | 90.9 % | 100 % |
| | R-Factor | 0.39 | 0.48 | 0.45 |
| | NSE | 0.992 | 0.973 | 0.989 |
| | Persistent-NSE | 0.908 | 0.821 | 0.823 |
| | KGE | 0.972 | 0.951 | 0.973 |
| | RMSE | 25.7 | 44.2 | 20.2 |
| N-BEATS | MAE | 18.3 | 25.9 | 14.0 |
| | PFE | 0.006 | 0.008 | 0.019 |
| | TPE (hours) | 0 | 0 | 0 |
| | P-Factor | 96.3 % | 86.3 % | 96.9 % |
| | R-Factor | 0.43 | 0.53 | 0.43 |
| | NSE | 0.952 | 0.892 | 0.935 |
| LSTM | Persistent-NSE | 0.4 | 0.27 | 0.087 |
| | KGE | 0.92 | 0.899 | 0.901 |

| | | | |
|---|---|---|---|
| **RMSE** | 65.7 | 89.2 | 50.3 |
| **MAE** | 41.1 | 45 | 35.9 |
| **PFE** | 0.031 | 0.058 | 0.098 |
| **TPE (hours)** | 1 | 0 | 0 |
| **P-Factor** | 70.4 % | 72.73 % | 81.82 % |
| **R-Factor** | 0.66 | 0.7 | 0.65 |

649



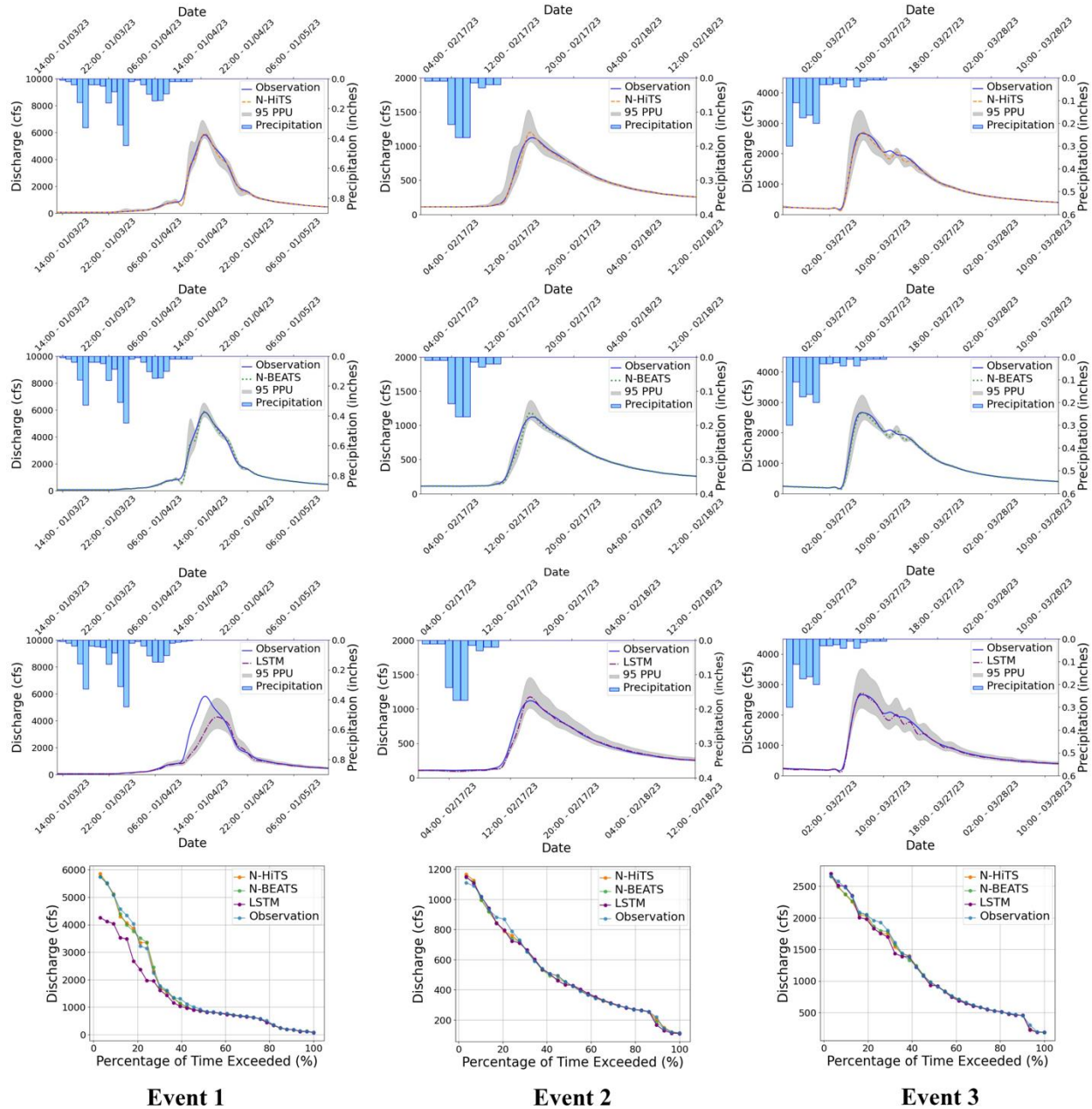Event 1                          Event 2                          Event 3

650

651      Figure 7. 95 PPU band and FDC plots of N-HiTS, N-BEATS, and LSTM models for the three selected
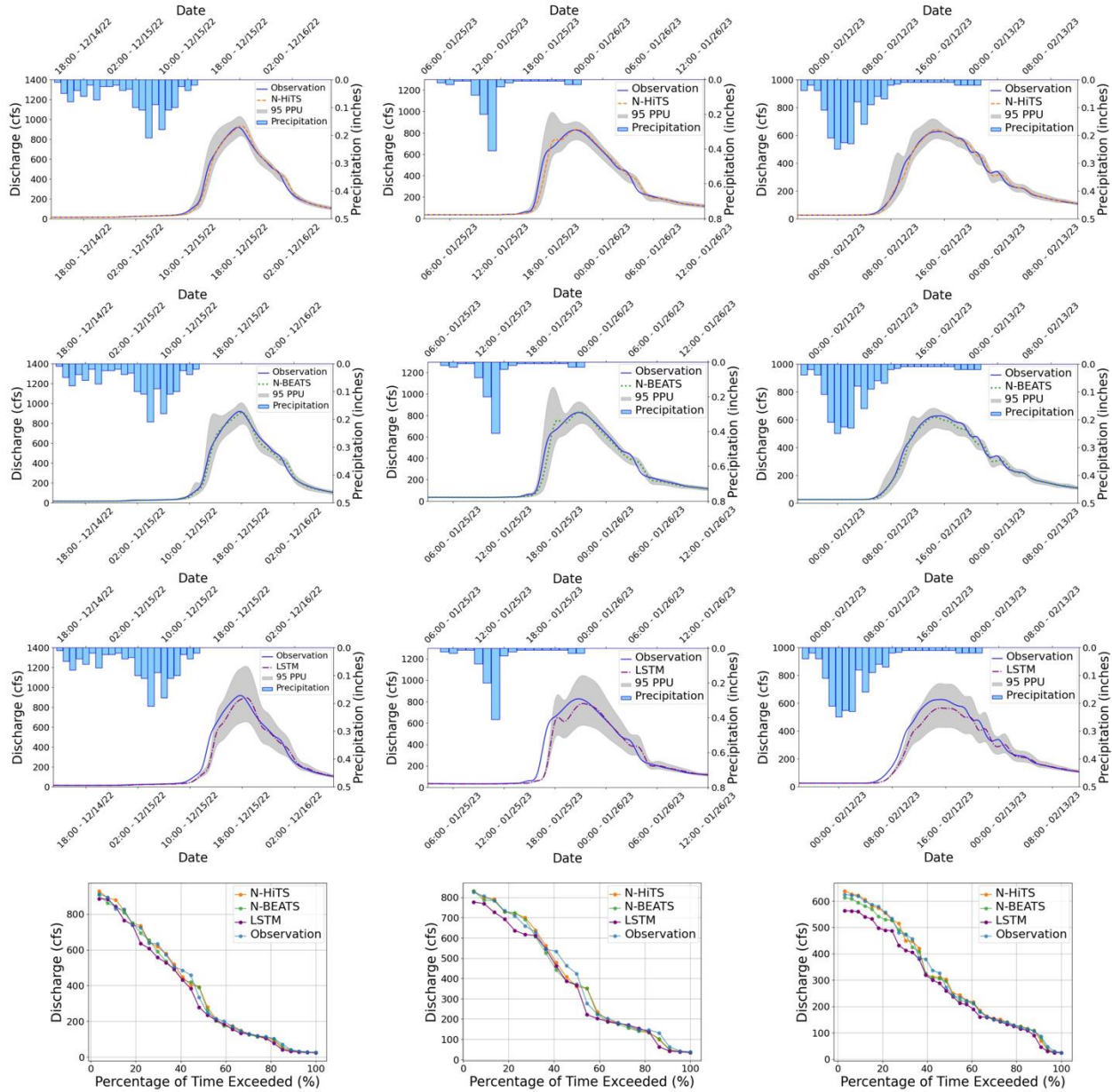652                          flooding events in the Lower Dog River gauging station.

Figure 8. 95 PPU band and FDC plots of N-HiTS, N-BEATS, and LSTM models for the three selected flooding events in the Killian Creek gauging station.

In our investigation, we conducted an analysis to assess the impact of varying input sizes on the performance of the N-HiTS, as the best model. We implemented four different durations as input sizes to observe the corresponding differences in modeling performance. Notably, one of the key metrics affected by changes in input size was 95PPU, which exhibited a general decrease with increasing input size.

As detailed in Table 5, we observed a discernible trend in the R-Factor of the N-HiTS model as the input size was increased. Specifically, there was a decline in the R-Factor as the input size expanded. This trend

662 underscores the influence of input size on model performance, particularly in terms of 95PPU band and

663 accuracy.

664 Overall, uncertainty analysis revealed that coupling MLE with N-HiTS and N-BEATS models

665 demonstrated superior performance in generating 95PPU, effectively reducing errors in flood prediction.

666 The MLE approach was more successful in reducing 95PPU bands of N-HiTS and N-BEATS models

667 compared to the LSTM, as indicated by the R-Factor and P-Factor. The N-BEATS model demonstrated a

668 narrower uncertainty band (lower R-Factor value), while the N-HiTS model provided higher precision.

669 Furthermore, incorporating data with various sizes into the N-HiTS model led to a narrower 95PPU and an

670 improvement in the R-Factor, highlighting the significance of input size in enhancing model accuracy and

671 reducing uncertainty.

Table 5. N-HiTS's R-Factor results for three storm events in each case study, using 1 hour, 2 hours, 12 hours, and 24 hours input size in training.

| Input Size | 1 hour | 6 hours | 12 hours | 24 hours |
|---|---|---|---|---|
| **Dog River, GA - Event 1** | 0.314 | 0.337 | 0.29 | 0.272 |
| **Dog River, GA - Event 2** | 0.35 | 0.413 | 0.403 | 0.402 |
| **Dog River, GA - Event 3** | 0.358 | 0.459 | 0.374 | 0.336 |
| **Killian Creek, NC - Event 4** | 0.491 | 0.422 | 0.426 | 0.388 |
| **Killian Creek, NC - Event 5** | 0.584 | 0.503 | 0.557 | 0.483 |
| **Killian Creek, NC - Event 6** | 0.482 | 0.42 | 0.446 | 0.454 |

672

673 **3.4. Sensitivity Analysis**

674 In this study, we conducted a comprehensive sensitivity analysis of the N-HiTS, N-BEATS, and LSTM

675 models to evaluate their responsiveness to meteorological variables, specifically precipitation, humidity,

676 and temperature. The goal was to assess how the omission of input parameters impacts the overall

677 modeling performance compared to their full-variable counterparts.

678 To execute this analysis, we systematically trained each model by excluding meteorological variables one

679 or more at a time, subsequently evaluating their predictive performance using the entire testing dataset.

680 The results of our analysis indicated that N-HiTS and N-BEATS models exhibited minimal sensitivity to

681 meteorological variables, as evidenced by the negligible impact on their performance metric (i.e., NSE,

682 Persistent-NSE, KGE, RMSE, and MAE) upon parameter exclusion.

683 Notably, as shown in Table 6, the performance of the N-HiTS model displayed a marginal deviation

684 under variable omission, while the N-BEATS model exhibited consistent performance irrespective of the

32

685    inclusion or exclusion of meteorological variables. The structure of this algorithm is based on backward

686    and forward residual links for univariate time series point forecasting which does not take into account

687    other parameters in the prediction task.  These findings suggest that the predictive capabilities of N-HiTS

688    and N-BEATS models predominantly rely on historical flood data. Both models demonstrated strong

689    performance even without incorporating precipitation, temperature, or humidity data, underscoring their

690    ability in flood prediction in the absence of specific meteorological inputs. This capability underscores the

691    robustness of the N-HiTS and N-BEATS models, positioning them as viable tools and perhaps

692    appropriate for real-time flood forecasting tasks where direct meteorological data may be limited or

693    unavailable.

694

695    Table 6. Performance metrics' values for N-HiTS, N-BEATS, and LSTM models by excluding

696    meteorological variables one or more at a time.

| Model | Excluded Variables | NSE | Persistent-NSE | KGE | RMSE | MAE |
|---|---|---|---|---|---|---|
| **N-HiTS** | Using all variables | 0.996 | 0.92 | 0.988 | 22.66 | 4.19 |
| | Without Precipitation | 0.993 | 0.91 | 0.97 | 23.28 | 4.31 |
| | Without Humidity | 0.995 | 0.914 | 0.976 | 22.87 | 4.22 |
| | Without Temperature | 0.995 | 0.921 | 0.985 | 22.43 | 4.14 |
| | Discharge only prediction | 0.993 | 0.911 | 0.972 | 23.21 | 4.29 |
| **N-BEATS** | Using all variables | 0.994 | 0.978 | 0.992 | 11.80 | 2.13 |
| | Without Precipitation | 0.994 | 0.978 | 0.991 | 11.86 | 2.17 |
| | Without Humidity | 0.994 | 0.978 | 0.991 | 11.81 | 2.16 |
| | Without Temperature | 0.994 | 0.978 | 0.991 | 11.82 | 2.16 |
| | Discharge only prediction | 0.994 | 0.978 | 0.991 | 11.96 | 2.17 |
| | Using all variables | 0.992 | 0.865 | 0.926 | 29.52 | 8.15 |

| | | | | | | |
|---|---|---|---|---|---|---|
| | Without Precipitation | 0.979 | 0.665 | 0.892 | 39.46 | 19.83 |
| | Without Humidity | 0.991 | 0.843 | 0.925 | 31.73 | 9.15 |
| **LSTM** | Without Temperature | 0.983 | 0.628 | 0.872 | 48.95 | 11.49 |
| | Discharge only prediction | 0.976 | 0.576 | 0.692 | 52.28 | 33.5 |

697

## 3.5 Computational Efficiency

The computational efficiency of the N-HiTS, N-BEATS, and LSTM models, as well as a comparative analysis, is presented in Table 7. The study encompassed the entire process of training and predicting over the testing period, employing the optimized hyperparameters as previously described. Regarding the training time, it is noteworthy that the LSTM model exhibited the quickest performance. Specifically, LSTM demonstrated a training time that was 71% faster than N-HiTS and 93% faster than N-BEATS in the Lower Dog River watershed, while it was respectively,126% and 118% faster than N-HiTS and N-BEATS in the Upper Dutchmans Creek, over training dataset. This is because LSTM has a simple architecture compared to the N-BEATS and N-HiTS and does not require multivariate features, hierarchical interpolation, and multi-rate data sampling. Perhaps, this outcome underscores the computational advantage of LSTM over other algorithms.

Conversely, during the testing period, the N-HiTS model emerged as the fastest and delivered the most efficient results in comparison to the other models. Notably, N-HiTS displayed a predicting time that was 33% faster than LSTM and 32% faster than N-BEATS. This finding highlights the computational efficiency of the N-HiTS model in the context of predicting processes. Our experiments unveiled an interesting contrast in the computational performance of these models. While LSTM excelled in terms of training time, it lagged behind when it came to the testing period.

In the grand scheme of computational efficiency, model accuracy, and uncertainty analysis results, it becomes evident that the superiority of the N-HiTS and N-BEATS models in terms of accuracy and uncertainty analysis holds paramount importance. This significance is accentuated by the critical nature of flood prediction, where precision and certainty are pivotal. Therefore, computational efficiency must be viewed in the context of the broader objectives, with the accuracy and reliability of flood predictions taking precedence in ensuring the safety and preparedness of the affected regions.

721

722   Table 7. Computational costs of N-HiTS, N-BEATS, and LSTM models in the Dog River and Killian
723                              Creek gauging stations.

| Model | Training Time over Train Datasets (seconds) | | Predicting Time over Test Datasets (seconds) | |
|---|---|---|---|---|
|  | Lower Dog River | Upper Dutchmans Creek | Lower Dog River | Upper Dutchmans Creek |
| N-HiTS | 256.032 | 374.569 | 1533.029 | 1205.526 |
| N-BEATS | 288.511 | 361.599 | 2028.068 | 1482.305 |
| LSTM | 149.173 | 165.827 | 2046.140 | 1792.444 |

724

725   **4. Conclusion**

726   This study examined multiple NN algorithms for flood prediction. We selected two headwater streams with
727   minimal human impacts to understand how NN approaches can capture flood magnitude and timing for
728   these natural systems. In conclusion, our study represents a pioneering effort in exploring and advancing
729   the application of NN algorithms, specifically the N-HiTS and N-BEATS models, in the field of flood
730   prediction. In our case studies, both N-HiTS and N-BEATS models achieved state-of-the-art results,
731   outperforming LSTM as a benchmark model, particularly in one-hour prediction. While a one-hour lead
732   time may seem brief, it is highly significant for accurate flash flood prediction particularly in  an area with
733   a proximity to large metropolitan cities, where rapid response is critical.  These benchmarking results are
734   arguably a pivotal part of this research. However, the N-BEATS model slightly emerged as a powerful and
735   interpretable tool for flood prediction in most selected events.
736   In addition, the results of the experiments described above demonstrated that N-HiTS multi-rate input
737   sampling and hierarchical interpolation along with N-BEATS interpretable configuration are effective in
738   learning location-specific runoff generation behaviors. Both algorithms with an MLP-based deep neural
739   architecture with backward and forward residual links can sequentially project the data signal into
740   polynomials and harmonic basis needed to predict intense storm behaviors with varied magnitudes. The
741   innovation in this study – besides benchmarking the LSTM model for headwater streams – was to tackle
742   volatility and memory complexity challenges, by locally specializing flood sequential predictions into the
743   data signal's frequencies with interpretability, and hierarchical interpolation and pooling. Both N-HiTS and
744   N-BEATS models offered similar performance as compared with the LSTM but also offered a level of
745   interpretability about how the model learns to differentiate aspects of complex watershed-specific behaviors
746   via data. The interpretability of N-HiTS and N-BEATS models stems from their designs. N-HiTS aims to
747   enhance the accuracy of long-term time-series forecasts through hierarchical interpolation and multi-scale

748    data sampling, allowing it to focus on different data patterns, which prioritizes features essential to

749    understand flood magnitudes. N-BEATS leverages interpretable configurations with trend and seasonality

750    projections, enabling it to decompose time series data into intuitive components. N-BEATS interpretable

751    architecture is recommended for scarce data settings (such as flooding event), as it regularizes its

752    predictions through projections unto harmonic and trend basis. These approaches improve model

753    transparency by allowing understanding of how each part of the model contributes to the final prediction,

754    particularly when applied to complex flood patterns. Both models also support multivariate series (and

755    covariates) by flattening the model inputs to a 1-D series and reshaping the outputs to a tensor of appropriate

756    dimensions. This approach provides flexibility to handle arbitrary numbers of features. Furthermore, both

757    N-HiTS and N-BEATS models also support producing probabilistic predictions by specifying a likelihood

758    parameter. In terms of sensitivity analysis, both N-HiTS and N-BEATS models maintain consistent

759    performance even when trained without specific meteorological inputs. Although, during some flashy

760    floods, the models encountered challenges in capturing the peak flows and the dynamics of the recession

761    curve, which is directly related to groundwater contribution to flood hydrograph, both models were

762    technically insensitive to rainfall data as an input variable. This suggests the fact that both algorithms can

763    learn patterns in discharge data without requiring meteorological input. This ability underscores these

764    models' robustness in generating accurate predictions using historical flood data alone, making them

765    valuable tools for flood prediction, especially in data-poor watersheds or even for real-time flood prediction

766    when near real-time meteorological inputs are limited or unavailable. In terms of computational efficiency,

767    both N-HiTS and N-BEATS are trained almost at the same pace; however, N-HiTS predicted the test data

768    much quicker than N-BEATS. Unlike N-HiTS and N-BEATS, LSTM excelled in reducing training time

769    due to its simplicity and limited number of parameters.

770    Moving forward, it is worth mentioning that predicting the magnitude of the recession curve of flood

771    hydrographs was particularly challenging for all models. We argue that this is because the relation between

772    base flow and time is particularly hard to calibrate due to ground-water effluent that is controlled by

773    geological and physical conditions (vegetation, wetlands, wet meadows) in headwater streams. In addition,

774    the situations of runoff occurrence are diverse and have a high measurement variance with high frequency

775    that can make it difficult for the algorithms to fully capture discrete representation learning on time series.

776    In future studies, it will be important to develop strategies to derive analogs to the interpretable

777    configuration as well as multi-rate input sampling, hierarchical interpolation, and backcast residual

778    connections that allow for the dynamic representation of flood times series data with different frequencies

779    and nonlinearity. A dynamic representation of flood time series is, at least in principle, possible by

780    generating additive predictions in different bands of the time-series signals, reducing memory footprint and

781    compute time, and improving architecture parsimony and accuracy. This would allow the model to "learn"

782 interpretability and hierarchical representations from raw data to reduce complexity as the information
783 flows through the network. Moreover, it is noteworthy that while a single station offers valuable localized
784 data, particularly for smaller watersheds such as headwater streams where runoff is closely tied to
785 immediate meteorological conditions, it may not fully capture the spatial heterogeneity of larger
786 watersheds. For our specific case, the methods applied herein captured runoff magnitude and dynamics in
787 small watersheds using a single station. However, we recognize that for broader areas, incorporating
788 spatially distributed data would likely enhance model accuracy. Lastly, one could explore the idea of
789 enhancing N-HiTS and N-BEATS (or NN algorithms, in general) performance with uncertainty
790 quantification by using more robust Bayesian inference such as Bayesian Model Averaging (BMA) with
791 fixed and flexible prior distributions (see Samadi et al., 2020) and/or Markov Chain Monte-Carlo
792 optimization methods (Duane et al., 1987) addressing both aleatoric and epistemic uncertainties. We leave
793 these approaches for future discussion and exploration in the context of flood neural time series prediction.
794

## 5. Acknowledgements

802

## 6. Open Research

804 The historical discharge data used in this study are from the USGS
805 (https://waterdata.usgs.gov/nwis/uv/?referred_module=sw), meteorological data from USDA
806 (https://www.ncdc.noaa.gov/cdo-web/datatools/lcd). We have uploaded the datasets and codes
807 used in this research to Zenodo, accessible via https://zenodo.org/records/13343364. For
808 modeling, we used the NeuralForecast package (Olivares et al., 2022), available at:
809 https://github.com/Nixtla/neuralforecast.
810

## 7. References

812 Abbaspour, K.C., Yang, J., Maximov, I., Siber, R., Bogner, K., Mieleitner, J., Zobrist, J., Srinivasan, R.,
813     2007. Modelling hydrology and water quality in the pre-alpine/alpine Thur watershed using SWAT.
814     Journal of Hydrology 333, 413–430. https://doi.org/10.1016/j.jhydrol.2006.09.014

815    Alaa, A.M., van der Schaar, M., 2019. Attentive State-Space Modeling of Disease Progression, in:
816        Advances in Neural Information Processing Systems. Curran Associates, Inc.

817    Asquith, W.H., Roussel, M.C., Thompson, D.B., Cleveland, T.G., Fang, X., 2005. Summary of
818        dimensionless Texas hyetographs and distribution of storm depth developed for Texas Department
819        of Transportation research project 0–4194 (No. 0–4194–4). Texas Department of Transportation.

820    Barnard, P.L., van Ormondt, M., Erikson, L.H., Eshleman, J., Hapke, C., Ruggiero, P., Adams, P.N.,
821        Foxgrover, A.C., 2014. Development of the Coastal Storm Modeling System (CoSMoS) for
822        predicting the impact of storms on high-energy, active-margin coasts. Nat Hazards 74, 1095–1125.
823        https://doi.org/10.1007/s11069-014-1236-y

824    Basso, S., Schirmer, M., Botter, G., 2016. A physically based analytical model of flood frequency curves.
825        Geophysical Research Letters 43, 9070–9076. https://doi.org/10.1002/2016GL069915

826    Challu, C., Olivares, K.G., Oreshkin, B.N., Garza, F., Mergenthaler-Canseco, M., Dubrawski, A., 2022.
827        N-HiTS: Neural Hierarchical Interpolation for Time Series Forecasting.
828        https://doi.org/10.48550/arXiv.2201.12886

829    Chen, Y., Li, J., Xu, H., 2016. Improving flood forecasting capability of physically based distributed
830        hydrological models by parameter optimization. Hydrology and Earth System Sciences 20, 375–
831        392. https://doi.org/10.5194/hess-20-375-2016

832    Clark, M.P., Nijssen, B., Lundquist, J.D., Kavetski, D., Rupp, D.E., Woods, R.A., Freer, J.E., Gutmann,
833        E.D., Wood, A.W., Brekke, L.D., Arnold, J.R., Gochis, D.J., Rasmussen, R.M., 2015. A unified
834        approach for process-based hydrologic modeling: 1. Modeling concept. Water Resources Research
835        51, 2498–2514. https://doi.org/10.1002/2015WR017198

836    CRED, n.d. EM-DAT - The international disaster database [WWW Document]. URL
837        https://www.emdat.be/ (accessed 6.5.24).

838    Dasgupta, A., Arnal, L., Emerton, R., Harrigan, S., Matthews, G., Muhammad, A., O'Regan, K., Pérez-
839        Ciria, T., Valdez, E., van Osnabrugge, B., Werner, M., Buontempo, C., Cloke, H., Pappenberger,
840        F., Pechlivanidis, I.G., Prudhomme, C., Ramos, M.-H., Salamon, P., n.d. Connecting hydrological
841        modelling and forecasting from global to local scales: Perspectives from an international joint
842        virtual workshop. Journal of Flood Risk Management n/a, e12880.
843        https://doi.org/10.1111/jfr3.12880

844    Defontaine, T., Ricci, S., Lapeyre, C., Marchandise, A., Pape, E.L., 2023. Flood forecasting with Machine
845        Learning in a scarce data layout. IOP Conf. Ser.: Earth Environ. Sci. 1136, 012020.
846        https://doi.org/10.1088/1755-1315/1136/1/012020

847    Duane, S., Kennedy, A.D., Pendleton, B.J., Roweth, D., 1987. Hybrid Monte Carlo. Physics Letters B
848        195, 216–222. https://doi.org/10.1016/0370-2693(87)91197-X

849 Erikson, L.H., Espejo, A., Barnard, P.L., Serafin, K.A., Hegermiller, C.A., O'Neill, A., Ruggiero, P.,
850    Limber, P.W., Mendez, F.J., 2018. Identification of storm events and contiguous coastal sections
851    for deterministic modeling of extreme coastal flood events in response to climate change. Coastal
852    Engineering 140, 316–330. https://doi.org/10.1016/j.coastaleng.2018.08.003

853 Evin, G., Le Lay, M., Fouchier, C., Mas, A., Colleoni, F., Penot, D., Garambois, P.-A., Laurantin, O.,
854    2023. Evaluation of hydrological models on small mountainous catchments: impact of the
855    meteorological forcings. https://doi.org/10.5194/egusphere-2023-845

856 Fan, C., Zhang, Y., Pan, Y., Li, X., Zhang, C., Yuan, R., Wu, D., Wang, W., Pei, J., Huang, H., 2019.
857    Multi-Horizon Time Series Forecasting with Temporal Attention Learning, in: Proceedings of the
858    25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '19.
859    Association for Computing Machinery, New York, NY, USA, pp. 2527–2535.
860    https://doi.org/10.1145/3292500.3330662

861 Fang, K., Kifer, D., Lawson, K., Shen, C., 2020. Evaluating the Potential and Challenges of an
862    Uncertainty Quantification Method for Long Short-Term Memory Models for Soil Moisture
863    Predictions. Water Resources Research 56, e2020WR028095.
864    https://doi.org/10.1029/2020WR028095

865 Global assessment report on disaster risk reduction 2015 | UNDRR, 2015. URL:
866    http://www.undrr.org/publication/global-assessment-report-disaster-risk-reduction-2015 (accessed
867    6.5.24).

868 Gotvald, A.J., 2010, Historic flooding in Georgia, 2009: U.S. Geological Survey Open-File Report 2010–
869    1230, 19 p.

870 Gupta, H.V., Kling, H., Yilmaz, K.K., Martinez, G.F., 2009. Decomposition of the mean squared error
871    and NSE performance criteria: Implications for improving hydrological modelling. Journal of
872    Hydrology 377, 80–91. https://doi.org/10.1016/j.jhydrol.2009.08.003

873 Hochreiter, S., Younger, A.S., Conwell, P.R., 2001. Learning to Learn Using Gradient Descent, in:
874    Dorffner, G., Bischof, H., Hornik, K. (Eds.), Artificial Neural Networks — ICANN 2001. Springer,
875    Berlin, Heidelberg, pp. 87–94. https://doi.org/10.1007/3-540-44668-0_13

876 Hsu, K., Gupta, H.V., Sorooshian, S., 1995. Artificial Neural Network Modeling of the Rainfall-Runoff
877    Process. Water Resources Research 31, 2517–2530. https://doi.org/10.1029/95WR01955

878 Jonkman, S.N., 2005. Global Perspectives on Loss of Human Life Caused by Floods. Nat Hazards 34,
879    151–175. https://doi.org/10.1007/s11069-004-8891-3

880 Kingma, D.P., Ba, J., 2017. Adam: A Method for Stochastic Optimization.
881    https://doi.org/10.48550/arXiv.1412.6980

882    Kratzert, F., Klotz, D., Brenner, C., Schulz, K., Herrnegger, M., 2018. Rainfall–runoff modelling using
883        Long Short-Term Memory (LSTM) networks. Hydrology and Earth System Sciences 22, 6005–
884        6022. https://doi.org/10.5194/hess-22-6005-2018

885    Lim, B., Arık, S.Ö., Loeff, N., Pfister, T., 2021. Temporal Fusion Transformers for interpretable multi-
886        horizon time series forecasting. International Journal of Forecasting 37, 1748–1764.
887        https://doi.org/10.1016/j.ijforecast.2021.03.012

888    Lobligeois, F., Andréassian, V., Perrin, C., Tabary, P., Loumagne, C., 2014. When does higher spatial
889        resolution rainfall information improve streamflow simulation? An evaluation using 3620 flood
890        events. Hydrology and Earth System Sciences 18, 575–594. https://doi.org/10.5194/hess-18-575-
891        2014

892    MacDonald, L.H., Coe, D., 2007. Influence of Headwater Streams on Downstream Reaches in Forested
893        Areas. Forest Science 53, 148–168. https://doi.org/10.1093/forestscience/53.2.148

894    Martinaitis, S.M., Wilson, K.A., Yussouf, N., Gourley, J.J., Vergara, H., Meyer, T.C., Heinselman, P.L.,
895        Gerard, A., Berry, K.L., Vergara, A. and Monroe, J., 2023. A path toward short-term probabilistic
896        flash flood prediction. Bulletin of the American Meteorological Society, 104(3), pp.E585-E605.

897    McCallum, B.E., and Gotvald, A.J., 2010, Historic flooding in northern Georgia, September 16–22, 2009:
898        U.S. Geological Survey Fact Sheet 2010–3061, 4 p.

899    McCuen, R.H., Knight, Z., Cutter, A.G., 2006. Evaluation of the Nash–Sutcliffe Efficiency Index. Journal
900        of Hydrologic Engineering 11, 597–602. https://doi.org/10.1061/(ASCE)1084-
901        0699(2006)11:6(597)

902    Munn, M., Sheibley, R., Waite, I., Meador, M., 2020. Understanding the relationship between stream
903        metabolism and biological assemblages. Freshwater Science 39, 680–692.
904        https://doi.org/10.1086/711690

905    Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual models part I — A discussion
906        of principles. Journal of Hydrology 10, 282–290. https://doi.org/10.1016/0022-1694(70)90255-6

907    Nevo, S., Morin, E., Gerzi Rosenthal, A., Metzger, A., Barshai, C., Weitzner, D., Voloshin, D., Kratzert,
908        F., Elidan, G., Dror, G., Begelman, G., Nearing, G., Shalev, G., Noga, H., Shavitt, I., Yuklea, L.,
909        Royz, M., Giladi, N., Peled Levi, N., Reich, O., Gilon, O., Maor, R., Timnat, S., Shechter, T.,
910        Anisimov, V., Gigi, Y., Levin, Y., Moshe, Z., Ben-Haim, Z., Hassidim, A., Matias, Y., 2022. Flood
911        forecasting with machine learning models in an operational framework. Hydrology and Earth
912        System Sciences 26, 4013–4032. https://doi.org/10.5194/hess-26-4013-2022

913    NRCS (2009). Part 630 Hydrology National Engineering Handbook, Chapter 15: Time of Concentration.

914    Olivares, K. G., Challú, C., Garza, F., Mergenthaler Canseco, M., & Dubrawski, A. (2022).
915        NeuralForecast: User friendly state-of-the-art neural forecasting models. PyCon Salt Lake City,
916        Utah, US 2022. Retrieved from https://github.com/Nixtla/neuralforecast

917  Olivares, K.G., Meetei, O.N., Ma, R., Reddy, R., Cao, M., Dicker, L., 2024. Probabilistic hierarchical
918      forecasting with deep Poisson mixtures. International Journal of Forecasting 40, 470–489.
919      https://doi.org/10.1016/j.ijforecast.2023.04.007

920  Oreshkin, B.N., Carpov, D., Chapados, N., Bengio, Y., 2020. N-BEATS: Neural basis expansion analysis
921      for interpretable time series forecasting. https://doi.org/10.48550/arXiv.1905.10437

922  Pally, R.J., Samadi, V., 2021. Application of image processing and convolutional neural networks for
923      flood image classification and semantic segmentation. Environmental Modelling & Software 148,
924      105285. https://doi.org/10.1016/j.envsoft.2021.105285

925  Palmer, T.N., 2012. Towards the probabilistic Earth-system simulator: a vision for the future of climate
926      and weather prediction. Quarterly Journal of the Royal Meteorological Society 138, 841–861.
927      https://doi.org/10.1002/qj.1923

928  Park, K., Lee, E.H., 2024. Urban flood vulnerability analysis and prediction based on the land use using
929      Deep Neural Network. International Journal of Disaster Risk Reduction 101, 104231.
930      https://doi.org/10.1016/j.ijdrr.2023.104231

931  Pourreza-Bilondi, M., Samadi, S.Z., Akhoond-Ali, A.-M., Ghahraman, B., 2017. Reliability of Semiarid
932      Flash Flood Modeling Using Bayesian Framework. Journal of Hydrologic Engineering 22,
933      05016039. https://doi.org/10.1061/(ASCE)HE.1943-5584.0001482

934  Refsgaard, J.C., Stisen, S., Koch, J., 2022. Hydrological process knowledge in catchment modelling –
935      Lessons and perspectives from 60 years development. Hydrological Processes 36, e14463.
936      https://doi.org/10.1002/hyp.14463

937  Roelvink, D., Reniers, A., van Dongeren, A., van Thiel de Vries, J., McCall, R., Lescinski, J., 2009.
938      Modelling storm impacts on beaches, dunes and barrier islands. Coastal Engineering 56, 1133–
939      1152. https://doi.org/10.1016/j.coastaleng.2009.08.006

940  Russo, S., Perraudin, N., Stalder, S., Perez-Cruz, F., Leitao, J.P., Obozinski, G., Wegner, J.D., 2023. An
941      evaluation of deep learning models for predicting water depth evolution in urban floods.
942      https://doi.org/10.48550/arXiv.2302.10062

943  Safaei-Moghadam, A., Tarboton, D., Minsker, B., 2023. Estimating the likelihood of roadway pluvial
944      flood based on crowdsourced traffic data and depression-based DEM analysis. Natural Hazards and
945      Earth System Sciences 23, 1–19. https://doi.org/10.5194/nhess-23-1-2023

946  Saksena, S., Dey, S., Merwade, V., Singhofen, P.J., 2020. A Computationally Efficient and Physically
947      Based Approach for Urban Flood Modeling Using a Flexible Spatiotemporal Structure. Water
948      Resources Research 56, e2019WR025769. https://doi.org/10.1029/2019WR025769

949  Samadi, S., Pourreza-Bilondi, M., Wilson, C. a. M.E., Hitchcock, D.B., 2020. Bayesian Model Averaging
950      With Fixed and Flexible Priors: Theory, Concepts, and Calibration Experiments for Rainfall-Runoff

Modeling. Journal of Advances in Modeling Earth Systems 12, e2019MS001924. https://doi.org/10.1029/2019MS001924

Scott, J., n.d. Widespread Flooding After Severe Storms - WCCB Charlotte's CW. Available at: https://www.wccbcharlotte.com/2020/02/08/widespread-flooding-after-severe-storms/ (accessed 6.11.24).

Sukovich, E.M., Ralph, F.M., Barthold, F.E., Reynolds, D.W., Novak, D.R., 2014. Extreme Quantitative Precipitation Forecast Performance at the Weather Prediction Center from 2001 to 2011. Weather and Forecasting 29, 894–911. https://doi.org/10.1175/WAF-D-13-00061.1

Tabas, S.S., Samadi, S., 2022. Variational Bayesian dropout with a Gaussian prior for recurrent neural networks application in rainfall–runoff modeling. Environ. Res. Lett. 17, 065012. https://doi.org/10.1088/1748-9326/ac7247

Thompson, C.M., Frazier, T.G., 2014. Deterministic and probabilistic flood modeling for contemporary and future coastal and inland precipitation inundation. Applied Geography 50, 1–14. https://doi.org/10.1016/j.apgeog.2014.01.013

Tiwari, M.K., Chatterjee, C., 2010. Development of an accurate and reliable hourly flood forecasting model using wavelet-bootstrap-ANN (WBANN) hybrid approach. Journal of Hydrology 394, 458–470. https://doi.org/10.1016/j.jhydrol.2010.10.001

Watershed Report | Office of Water | US EPA, n.d. Available at: https://watersgeo.epa.gov/watershedreport/?comid=9224629 (accessed 6.9.24).

Wee, G., Chang, L.-C., Chang, F.-J., Mat Amin, M.Z., 2023. A flood Impact-Based forecasting system by fuzzy inference techniques. Journal of Hydrology 625, 130117. https://doi.org/10.1016/j.jhydrol.2023.130117

Windheuser, L., Karanjit, R., Pally, R., Samadi, S., Hubig, N.C., 2023. An End-To-End Flood Stage Prediction System Using Deep Neural Networks. Earth and Space Science 10, e2022EA002385. https://doi.org/10.1029/2022EA002385

Zafarmomen, N., Alizadeh, H., Bayat, M., Ehtiat, M., Moradkhani, H., 2024. Assimilation of Sentinel-Based Leaf Area Index for Modeling Surface-Ground Water Interactions in Irrigation Districts. Water Resources Research 60, e2023WR036080. https://doi.org/10.1029/2023WR036080

Zhang, L., Qin, H., Mao, J., Cao, X., Fu, G., 2023. High temporal resolution urban flood prediction using attention-based LSTM models. Journal of Hydrology 620, 129499. https://doi.org/10.1016/j.jhydrol.2023.129499

Zhang, Y., Pan, D., Griensven, J.V., Yang, S.X., Gharabaghi, B., 2023. Intelligent flood forecasting and warning: a survey. ir 3, 190–212. https://doi.org/10.20517/ir.2023.12

984    Zou, Y., Wang, J., Lei, P., Li, Y., 2023. A novel multi-step ahead forecasting model for flood based on time
985        residual LSTM. Journal of Hydrology 620, 129521. https://doi.org/10.1016/j.jhydrol.2023.129521

986