



1 Probabilistic Hierarchical Interpolation and Interpretable Configuration for Flood Prediction

2 Mostafa Saberian¹, Vidya Samadi^{2*}, Ioana Popescu³

3 1. The Glenn Department of Civil Engineering, Clemson University, Clemson, SC

4 2. Department of Agricultural Sciences, Clemson University, Clemson, SC.

5 3. Department of Hydroinformatics and Socio-Technical Innovation, IHE Delft Institute for Water
6 Education, Delft, the Netherlands

7 *Corresponding author: samadi@clemson.edu

8 Abstract

9 The last few years have witnessed the rise of Neural Networks (NNs) applications for hydrological time
10 series modeling. By virtue of their capabilities, NN models can achieve unprecedented levels of
11 performance when learn how to solve increasingly complex rainfall-runoff processes via data, making them
12 pivotal for the development of computational hydrologic tasks such as flood predictions. The NN models
13 should, in order to be considered practical, provide a probabilistic understanding of the model mechanisms
14 and predictions and hints on what could perturb the model. In this paper, we developed two probabilistic
15 NN models, i.e., Neural Hierarchical Interpolation for Time Series Forecasting (N-HiTS) and Network-
16 Based Expansion Analysis for Interpretable Time Series Forecasting (N-BEATS) and benchmarked them
17 with long short-term memory (LSTM) for flood prediction across two headwater streams in Georgia and
18 North Carolina, USA. To generate a probabilistic prediction, a Multi-Quantile Loss was used to assess the
19 95th percentile prediction uncertainty (95PPU) of multiple flooding events. We conducted extensive flood
20 prediction experiments demonstrating the advantages of hierarchical interpolation and interpretable
21 architecture, where both N-HiTS and N-BEATS provided an average accuracy improvement of almost 5%
22 (NSE) over the LSTM benchmarking model. On a variety of flooding events with different timing and
23 magnitudes, both N-HiTS and N-BEATS demonstrated significant performance improvements over the
24 LSTM benchmark and showcased their probabilistic predictions by specifying a likelihood parameter.

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

27 Short Summary

28 Recent progress in neural network accelerated improvements in the performance of catchment modeling
29 systems. Yet flood modeling remains a very difficult task. Focusing on two headwater streams, this paper
30 developed Neural Hierarchical Interpolation for Time Series Forecasting (N-HiTS) and Network-Based
31 Expansion Analysis for Interpretable Time Series Forecasting (N-BEATS) and benchmarked them with
32 long short-term memory (LSTM) to predict flooding events. Analysis suggested that both N-HiTS and N-
33 BEATS outperformed LSTM for short-term (1 hour) flood predictions. We demonstrated how the proposed



34 N-HiTS and N-BEATS architectures can be augmented with uncertainty and sensitivity approaches to
35 provide skilled flood predictions that are interpretable without considerable loss in accuracy.

36

37 **1. Introduction**

38 The last few years have been characterized by an upsurge in the Neural Networks (NN) models. As opaque
39 NN models are increasingly being employed to make important predictions in hydrological systems, the
40 demand for creating legitimate NN models is increasing in the hydrology community. However,
41 maintaining coherence while producing accurate predictions can be a challenging problem (Olivares et al.,
42 2024). There is a general agreement on the importance of providing probabilistic NN prediction (Samadi et
43 al., 2020), especially in the case of flood prediction (Martinaitis et al., 2023).

44 Flood occurrences have witnessed an alarming surge in frequency and severity globally. Jonkman (2005)
45 studied a natural disaster database (EM-DAT, 2023) and reported that over 27 years, more than 175000
46 people died, and close to 2.2 billion were affected directly by floods worldwide. These numbers are likely
47 an underestimation due to unreported events (Nevo et al., 2022). In addition, the United Nations Office for
48 Disaster Risk Reduction reported that flooding has been the most frequent, widespread weather-related
49 natural disaster since 1995, claiming over 600,000 lives, affecting around 4 billion people globally, and
50 causing annual economic damage of more than 100 billion USD (UNISDR, 2015). This escalating trend
51 has necessitated the need for better flood prediction and management strategies. Scholars have successfully
52 implemented different flood models such as deterministic (Roelvink et al., 2009, Thompson and Frazier,
53 2014; Barnard et al., 2014; Erikson et al., 2018) and physically based flood models (Basso et al., 2016;
54 Chen et al., 2016; Pourreza-Bilondi et al., 2017; Saksena et al., 2019; Refsgaard et al., 2021) in various
55 environmental systems over the past several decades. These studies have heightened the need for precise
56 flood prediction, they have also unveiled limitations inherent in existing deterministic and physics-based
57 models. While evidence suggests that both deterministic and physics-based approaches are meaningful and
58 useful (Sukovich et al., 2014), their forecasts rest heavily on imprecise and subjective expert opinion; there
59 is a challenge for setting robust evidence-based thresholds to issue flood warnings and alerts (Palmer, 2012).
60 Moreover, many of these traditional flood models particularly physically explicit models rely heavily on a
61 particular choice of numerical approximation and describe multiple process parameterizations only within
62 a fixed spatial architecture (Clark et al., 2015). Recent NN models have shown promising results across a
63 large variety of flood modeling applications (e.g., Nevo et al., 2022; Pally and Samadi, 2022; Dasgupta et
64 al., 2023; Zhang et al., 2023) and encourage the use of such methodologies as core drivers for neural flood
65 prediction (Windheuser et al., 2023).

66 Earlier adaptations of these intelligent techniques showed promising results for flood prediction (e.g., Hsu
67 et al., 1995; Tiwari and Chatterjee, 2010). However, recent efforts have taken NN application to the next



68 level, providing uncertainty assessment (Sadeghi Tabas and Samadi, 2022) and improvements over various
69 spatio-temporal scales, regions, and processes (e.g., Kratzert et al., 2018; Park and Lee, 2023; Zhang et al.,
70 2023). Nevo et al., 2022 were the first scholars who employed long short-term memory (LSTM) for flood
71 stage prediction and inundation mapping, achieving notable success during the 2021 monsoon season. Soon
72 after, Russo et al. (2023) evaluated various NN models for predicting flood depth in urban systems,
73 highlighting the potential of data-driven models for urban flood prediction. Similarly, Defontaine et al.
74 (2023) emphasized the role of NN algorithms in enhancing the reliability of flood predictions, particularly
75 in the context of limited data availability. Windheuser et al., (2023) studied flood gauge height forecasting
76 using images and time series data for two gauging stations in Georgia, USA. They used multiple NN models
77 such as Convolutional Neural Network (ConvNet/CNN) and LSTM to forecast floods in near real-time (up
78 to 72 hours). In a sequence, Wee et al., 2023 used Impact-Based Forecasting (IBF) to propose a Flood
79 Impact-Based Forecasting system (FIBF) using flexible fuzzy inference techniques, aiding decision-makers
80 in a timely response. Zou et al. (2023) proposed a Residual LSTM (ResLSTM) model to enhance and
81 address flood prediction gradient issues. They integrated Deep Autoregressive Recurrent (DeepAR) with
82 four recurrent neural networks (RNNs), including ResLSTM, LSTM, Gated Recurrent Unit (GRU), and
83 Time Feedforward Connections Single Gate Recurrent Unit (TFC-SGRU), and showed that ResLSTM
84 achieved superior accuracy. While these studies reported the superiority of NN models for flood modeling,
85 they highlighted a number of challenges, notably (i) the limited capability of proposed NN models to
86 capture the spatial variability and magnitudes of extreme data over time, (ii) the lack of a sophisticated
87 mechanism to capture different flood magnitudes and synthesize the prediction, and (iii) inability of the NN
88 models to process data in parallel and capture the relationships between all elements in a sequential manner.
89 Recent advances in neural time series forecasting showed promising results that can be used to address the
90 above challenges for flood prediction. Recent techniques include the adoption of the attention mechanism
91 and Transformer-inspired approaches (Fan et al. 2019; Alaa and van der Schaar 2019; Lim et al. 2021)
92 along with attention-free architectures composed of deep stacks of fully connected layers (Oreshkin et al.
93 2020). All of these approaches are relatively easy to scale up in terms of flood magnitudes (small to major
94 flood predictions), compared to LSTM and have proven to be capable of capturing spatiotemporal
95 dependencies (Challu et al., 2022). In addition, these architectures can capture input-output relationships
96 implicitly while they tend to be more computationally efficient. Many state-of-the-art NN approaches for
97 flood forecasting have been established based on LSTM. There are cell states in the LSTM networks that
98 can be interpreted as storage capacity often used in flood generation schemes. In LSTM, the updating of
99 internal cell states (or storages) is regulated through a number of gates: the first gate regulates the storage
100 depletion, the second one regulates storage fluctuations, and the third gate regulates the storages outflow
101 (Tabas and Samadi, 2022). The elaborate gated design of the LSTM partly solves the long-term dependency



102 problem in flood time series prediction (Fang et al., 2020), although, the structure of LSTMs is designed in
103 a sequential manner that cannot directly connect two nonadjacent portions (positions) of a time series. This
104 indicates the fact that data dependencies can flow from left to right, rather than in both directions as in the
105 case of the attention-based and Transformer approaches.

106 In this paper, we take a step in this direction by developing attention-free architecture, i.e. Neural
107 Hierarchical Interpolation for Time Series Forecasting (N-HiTS; Challu et al., 2022) and Network-Based
108 Expansion Analysis for Interpretable Time Series Forecasting (N-BEATS; Oreshkin et al., 2020) and
109 benchmarked these models with LSTM for flood prediction. We developed fully connected N-BEATS and
110 N-HiTS architectures using multi-rate data sampling, synthesizing the flood prediction outputs via multi-
111 scale interpolation.

112 We implemented all algorithms for flood prediction on two headwater streams i.e., the Lower Dog River,
113 Georgia, and the Upper Dutchmans Creek, North Carolina, USA. We selected two study areas to ensure
114 that the results are reliable. The results of N-BEATS and N-HiTS techniques were compared with the
115 benchmarking LSTM to understand how these techniques can improve the representations of rainfall and
116 runoff dispensing over a recurrence process. Notably, this study represents a pioneering effort, as to the
117 best of our knowledge, it is the first instance in which the application of N-BEATS and N-HiTS algorithms
118 in the field of flood prediction has been explored. The scope of this research will focus on:

119 1. *Flood prediction in a hierarchical fashion with interpretable outputs:* We built N-BEATS and N-HiTS
120 for flood prediction with a very deep stack of fully connected layers to implicitly capture input-output
121 relationships with hierarchical interpolation capabilities. The predictions also involve programming the
122 algorithms with decreasing complexity and aligning their time scale with the final output through multi-
123 scale hierarchical interpolation and interpretable architecture. Predictions were aggregated in a hierarchical
124 fashion that enabled the building of a very deep neural network with interpretable configurations.

125 2. *Uncertainty quantification of the models by employing probabilistic approaches:* a Multi-Quantile
126 Loss (MQL) was used to assess the 95th percentile prediction uncertainty (95PPU) of multiple flooding
127 events. MQL was integrated as the loss function to account for probabilistic prediction. MQL trains the
128 model to produce probabilistic forecasts by predicting multiple quantiles of the distribution of future values.

129 3. *Exploring headwater stream response to flooding:* Understanding the dynamic response of headwater
130 streams to flooding is essential for managing downstream flood risks. Headwater streams constitute the
131 uppermost sections of stream networks, usually comprising 60% to 80% of a catchment area. Given this
132 substantial coverage and the tendency for precipitation to increase with elevation, headwater streams are
133 responsible for generating and controlling the majority of runoff in downstream portions (MacDonald and
134 Coe, 2007).



135 The remainder of this paper is structured as follows. Section 2 presents the case study and data, NN models,
136 performance metrics, and sensitivity and uncertainty approaches. Section 3 focuses on the results of flood
137 predictions including sensitivity and uncertainty assessment and computation efficiency. Finally, Section 4
138 concludes the paper.

139

140 **2. Methodology**

141 **2.1. Case Study and Data**

142 This research used two headwater gauging stations located at the Lower Dog River watershed, Georgia
143 (GA; USGS02337410, Dog River gauging station), and the Upper Dutchmans Creek watershed, North
144 Carolina (NC; USGS0214269560, Killian Creek gauging station). As depicted in Figures 1 and 2, the Lower
145 Dog River and the Upper Dutchmans Creek watersheds are located in the west and north parts of two
146 metropolitan cities, Atlanta and Charlotte. As shown in Figure 1, the Lower Dog River stream gauge is
147 established southeast of Villa Rica in Carroll County, where the USGS has regularly monitored discharge
148 data since 2007 in 15-minute increments. The Lower Dog River is a stream with a length of 15.7 miles
149 (25.3 km; obtained from the U.S. Geological Survey [USGS] National Hydrography Dataset high-
150 resolution flowline data), an average elevation of 851.94 meters, and the watershed area above this gauging
151 station is 66.5 square miles (172 km²; obtained from the Georgia Department of Natural Resources). This
152 watershed is covered by 15.2% residential area, 14.6% agricultural land, and ~70% forest (Munn et al.,
153 2020). Killian Creek gauging station at the Upper Dutchmans Creek watershed is established
154 in Montgomery County, North Carolina, where the USGS has regularly monitored discharge data since
155 1995 in 15-minute increments. The Upper Dutchmans Creek is a stream with a length of 4.9 miles (7.9 km),
156 an average elevation of 642.2 meters (see Table 1), and the watershed area above this gauging station is 4
157 square miles (10.3 km²) with less than 3% residential area and about 93% forested land use (the United
158 States Environmental Protection Agency).

159

160 The Lower Dog River has experienced significant flooding in the last decades. For example, in September
161 2009, the creek, along with most of northern GA, experienced heavy rainfall (5 inches, equal to 94 mm).
162 The Dog River, overwhelmed by large amounts of overland flow from saturated ground in the watershed,
163 experienced massive flooding in September 2009 (Gotvald, 2010). The river crested at 33.8 feet (10.3 m)
164 with a peak discharge of 59,900 cfs (1,700 m³/s), nearly six times the 100-year flood level (McCallum and
165 Gotvald, 2010). In addition, Dutchmans Creek has experienced significant flooding in February 2020.
166 According to local news (WCCB Charlotte, 2020), the flood in Gaston County caused significant
167 infrastructure damage and community disruption. Key impacts included the threatened collapse of the
168 Dutchman's Creek bridge in Mt. Holly and the closure of Highway 7 in McAdenville.



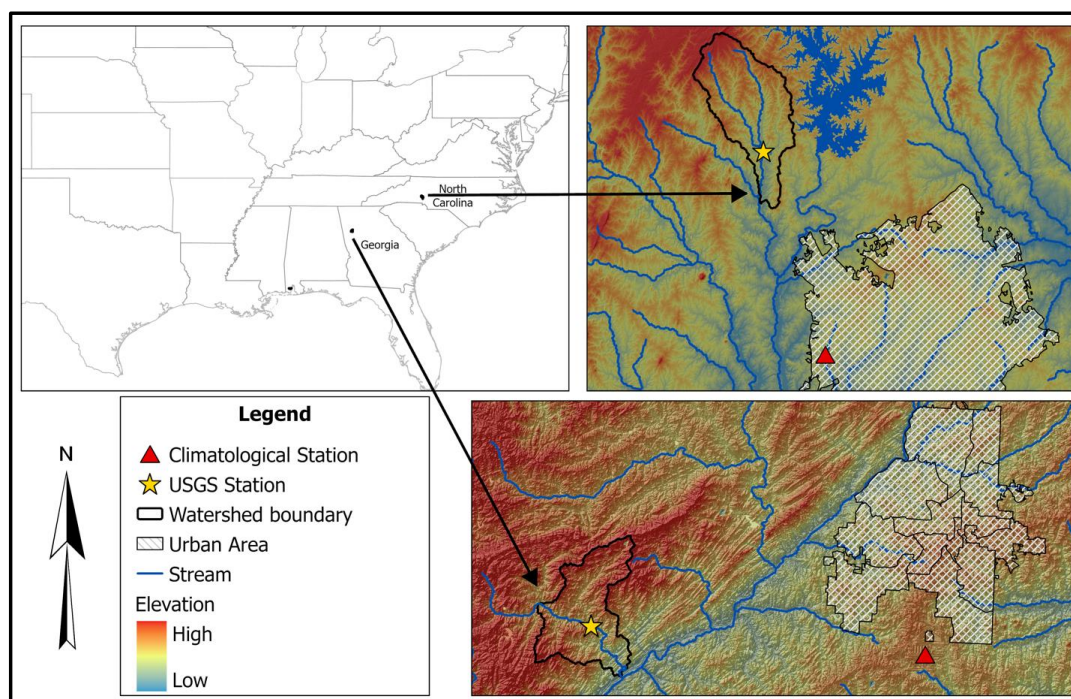
169

170

Table 1. Lower Dog River and Upper Dutchmans Creek’s physical characteristics.

Watershed	USGS Station ID Number	Average Elevation (m)	Stream Length (km)	Watershed area (km ²)
Lower Dog River watershed, GA	USGS02337410	851.9	25.3	172
Upper Dutchmans Creek watershed, NC	USGS0214269560	642.2	7.9	10.3

171



172

173 Figure 1. The Lower Dog River and The Upper Dutchmans Creek watersheds are located in GA and NC.

174 The proximity of the watersheds to Atlanta and Charlotte (urban area) are also displayed on the map.

175

176 To provide the meteorological forcing data, i.e., precipitation, temperature, and humidity, were extracted
 177 from the National Oceanic and Atmospheric Administration’s (NOAA) Local Climatological Data
 178 (LCD). We used the NOAA precipitation, temperature, and humidity data of Atlanta Hartsfield Jackson
 179 International Airport and Charlotte Douglas Airport stations as an input variable for neural network
 180 algorithms. The data has been monitored since January 1, 1948, and July 22, 1941, with an hourly interval
 181 which was used as an input variable for constructing neural networks.



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

189

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

205 **2.2. Neural Network Algorithms**

206 **2.2.1. LSTM**

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



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

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

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

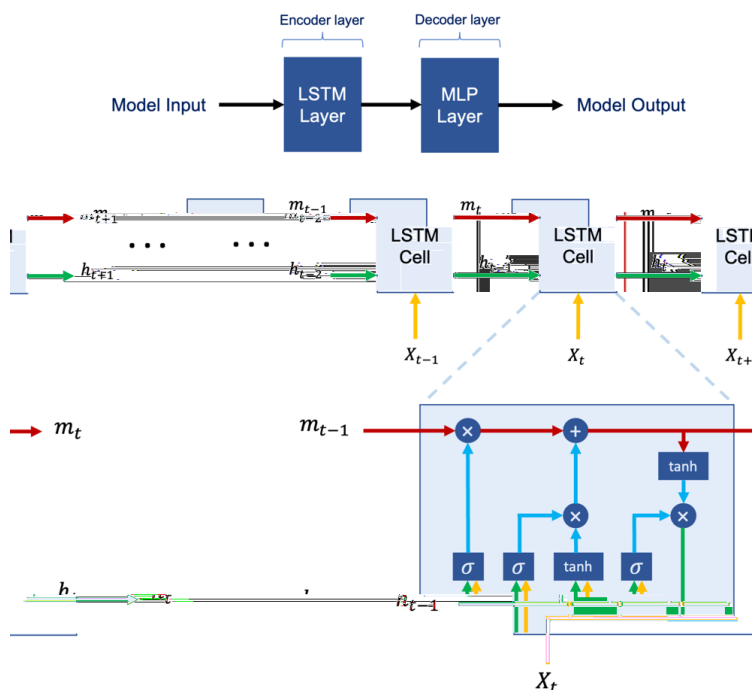
$$o_t = \sigma(A_o x_t + B_o h_{t-1} + c_o) \quad (\text{Equation 3})$$

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

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

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

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



231

232

233

234

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.

235 2.2.2. N-BEATS

236

237

238

239

240

241

242

243

244

245

246

247

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).



248

249

250

251

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

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

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

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

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

271 The second element performs the mapping of expansion coefficients θ_l^f and θ_l^b to produce outputs using
272 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:



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

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

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

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

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

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

295 **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
 296 across the forecast window, to replicate monotonic or slowly varying nature of trends:



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

297 The time vector $t = [0, 1, 2, \dots, H-2, H-1]^T / H$ is specified on a discrete grid ranging from 0 to
 298 $(H-1)/H$, projecting H steps into the future. Consequently, the trend forecast represented in matrix form is:

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

299

300 Where the polynomial coefficients, $\theta_{s,l}^f$, predicted by an FC network at layer l of stack s , are described by
 301 Equations (6) and (7). The matrix T , consisting of powers of t , is represented as $[1, t, \dots, t^p]$. When p is
 302 small, such as 2 or 3, it compels $\hat{y}_{s,l}^{tr}$ to emulate a trend (Oreshkin et al., 2020).

303 Seasonality model: In this stack $g_{s,l}^b$ and $g_{s,l}^f$ are periodic functions, to capture the cyclical and recurring
 304 characteristics of seasonality, such that $y_t = y_{t-\Delta}$, where Δ is the seasonality period. The Fourier series
 305 serves as a natural foundation for modeling periodic functions:

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

306

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

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

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

308

309 Where the Fourier coefficients $\theta_{s,l}^f$, that predicted by an FC network at layer l of stack s , are described by
 310 Equations (6) and (7). The matrix S represents sinusoidal waveforms. As a result, the forecast $\hat{y}_{s,l}^{seas}$
 311 becomes a periodic function that imitates typical seasonal patterns (Oreshkin et al., 2020).

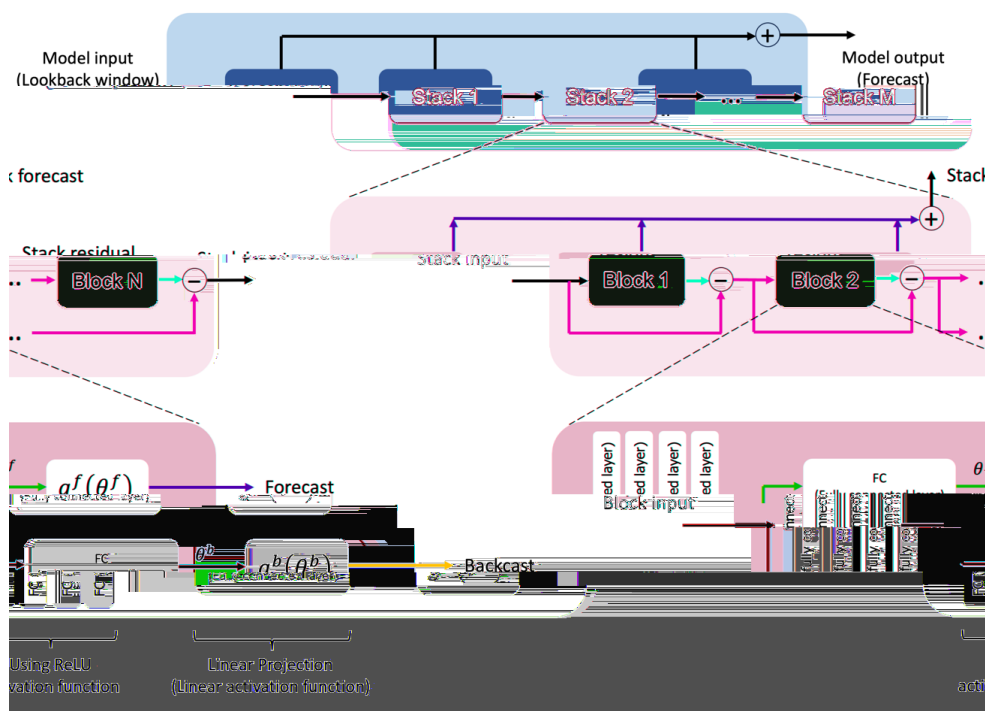


Figure 3. The N-BEATS modeling structure used in this research.

312

313

314

315 2.2.3. N-HiTS

316 N-HiTS builds upon the N-BEATS architecture but with improved accuracy and computational efficiency
 317 for long-horizon forecasting. N-HiTS utilizes multi-rate sampling and multi-scale synthesis of forecasts,
 318 leading to a hierarchical forecast structure that lowers computational demands and improves prediction
 319 accuracy (Challu et al., 2022).

320 Like N-BEATS, N-HiTS employs local nonlinear mappings onto foundational functions within numerous
 321 blocks. Each block includes an MLP that generates backcast and forecast output coefficients. The backcast
 322 output refines the input data for the following blocks, and the forecast outputs are combined to generate the
 323 final prediction. Blocks are organized into stacks, with each stack dedicated to grasping specific data
 324 attributes using its own distinct set of functions. The network's input is a sequence of L lags (look-back
 325 period), with S stacks, each containing B blocks (Challu et al., 2022).



326 In each block, a *MaxPool* layer with varying kernel sizes (k_l) is employed at the input, enabling the block
 327 to focus on specific input components of different scales. Larger kernel sizes emphasize the analysis of
 328 larger-scale, low-frequency data, aiding in improving long-term forecasting accuracy. This approach,
 329 known as multi-rate signal sampling, alters the effective input signal sampling rate for each block's MLP
 330 (Challu et al., 2022).

331 Additionally, multi-rate processing has several advantages. It reduces memory usage, computational
 332 demands, the number of learnable parameters, and helps prevent overfitting, while preserving the original
 333 receptive field. The following operation is applicable to the input $y_{t-L:t,l}$ of each block, with the first block
 334 ($l = 1$) using the network-wide input, where $y_{t-L:t,1} \equiv y_{t-L:t}$.

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

335 In many multi-horizon forecasting models, the number of neural network predictions matches the horizon's
 336 dimensionality, denoted as H . For instance, in N-BEATS, the number of predictions $|\theta_l^f| = H$. This results
 337 in a significant increase in computational demands and an unnecessary surge in model complexity as the
 338 horizon H becomes larger (Challu et al., 2022).

339 To address these challenges, N-HiTS proposes the use of temporal interpolation. This model manages the
 340 parameter counts per unit of output time ($|\theta_l^f| = \lceil r_l H \rceil$) by defining the dimensionality of the interpolation
 341 coefficients with respect to the expressiveness ratio r_l . To revert to the original sampling rate and predict
 342 all horizon points, this model employs temporal interpolation through the function g :

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

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

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

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

343 The hierarchical interpolation approach involves distributing expressiveness ratios over blocks, integrated
 344 with multi-rate sampling. Blocks closer to the input employ more aggressive interpolation, generating lower
 345 granularity signals. These blocks specialize in analyzing more aggressively subsampled signals. The final
 346 hierarchical prediction, $\hat{y}_{t+1:t+H}$, is constructed by combining outputs from all blocks, creating

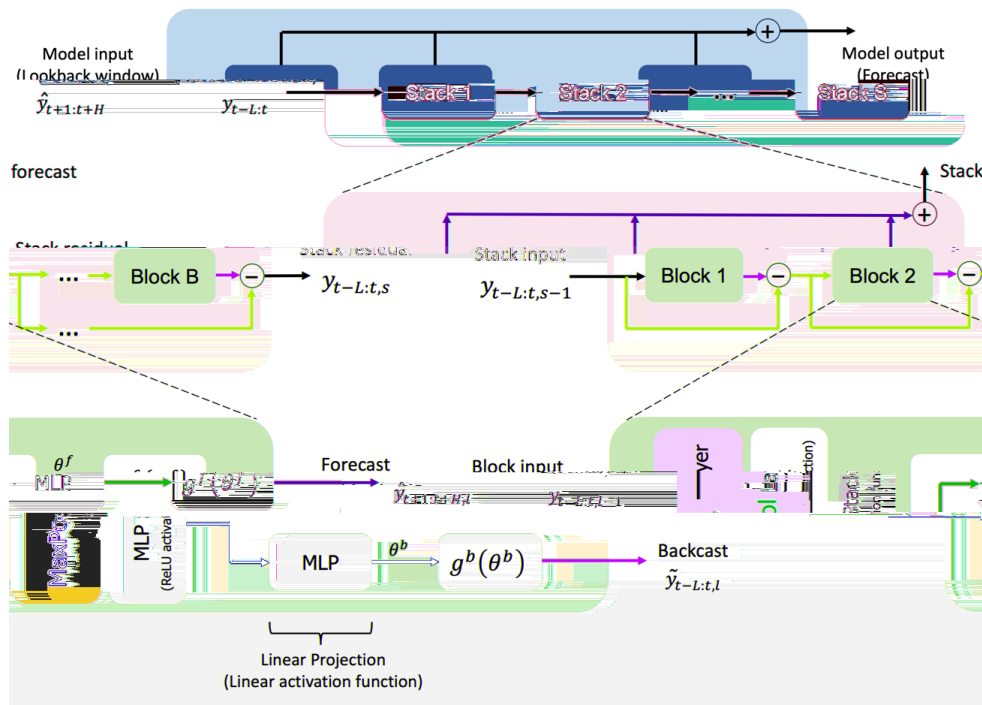


347 interpolations at various time-scale hierarchy levels. This approach maintains a structured hierarchy of
 348 interpolation granularity, with each block focusing on its own input and output scales (Challu et al., 2022).

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

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

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



355
 356 Figure 4. The structure of N-HiTS model programmed in this study. The architecture includes several
 357 Stacks, each Stack includes several Block, where each block consists of a MaxPool layer and a multi-
 358 layer which learn to produce coefficients for the backcast and forecast outputs of its basis.



359 2.3. Performance Metrics

360 To comprehensively evaluate the accuracy of flood predictions, we utilized a suite of metrics, including
361 Nash-Sutcliffe Efficiency (NSE), persistent Nash-Sutcliffe Efficiency (persistent-NSE), Root Mean Square
362 Error (RMSE), Mean Absolute Error (MAE), Peak Flow Error (PFE), and Time to Peak Error (TPE; Evin
363 et al., 2023; Lobligeois et al., 2014). These metrics collectively facilitate a rigorous assessment of the
364 model's performance in reproducing the magnitude of observed peak flows and the shape of the hydrograph.

365 The Nash–Sutcliffe model efficiency coefficient (NSE; Nash and Sutcliffe, 1970) measures the model's
366 ability to explain the variance in observed data and assesses the goodness-of-fit by comparing the observed
367 and simulated hydrographs. In hydrological studies, the NSE index is a widely accepted measure for
368 evaluating the fitting quality of models (McCuen et al., 2006). It is calculated as:

$$NSE = 1 - \frac{\sum_{i=1}^n (Q_{s_i} - Q_{o_i})^2}{\sum_{i=1}^n (Q_{o_i} - \bar{Q}_o)^2} \quad (\text{Equation 22})$$

369 Where Q_{o_i} represents observed value at time i , Q_{s_i} represents simulated value at time i , \bar{Q}_o is the mean
370 observed values and n is the number of data points. An NSE value of 1 indicates a perfect match between
371 the observed and modeled data, while lower values represent the degree of departure from a perfect fit.

372 As the models are designed to predict one hour ahead, the persistent-NSE is essential for evaluating their
373 performance. The standard NSE measures the model's sum of squared errors relative to the sum of squared
374 errors when the mean observation is used as the forecast value. In contrast, persistent-NSE uses the most
375 recent observed data as the forecast value for comparison (Nevo et al., 2022). The persistent-NSE is
376 calculated as:

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

377 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
378 observed value at the last time step ($i - 1$) and n is the number of data points. RMSE quantifies the average
379 magnitude of errors between observed and modeled values, offering insights into the absolute goodness-of-
380 fit, while MAE is a measure of the average absolute difference between the modeled values and the
381 observed values and provides a measure of the average magnitude of errors. RMSE is calculated as:



$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Q_{o_i} - Q_{s_i})^2} \quad (\text{Equation 24})$$

382 and MAE is calculated as:

$$MAE = \frac{1}{n} \sum_{i=1}^n |Q_{o_i} - Q_{s_i}| \quad (\text{Equation 25})$$

383 Where Q_{o_i} represents observed value at time i , Q_{s_i} represents simulated value at time i , and n is the number
384 of data points. RMSE and MAE provide information about the magnitude of modeling errors, with smaller
385 values indicating a better model fit.

386 PFE quantifies the magnitude disparity between observed and modeled peak flow values. The PFE metric
387 is defined as:

$$PFE = \frac{|Q_{o_{max}} - Q_{s_{max}}|}{Q_{o_{max}}} \quad (\text{Equation 26})$$

388 Where $Q_{o_{max}}$ represents the observed peak flow value, and $Q_{s_{max}}$ signifies the simulated peak flow value.
389 The PFE metric, expressed as a dimensionless value, provides a quantitative measure of the relative error
390 in predicting peak flow magnitudes concerning the observed values. A smaller PFE denotes more accurate
391 modeling of peak flow magnitudes, with a value of zero indicating a perfect match.

392 TPE assesses the temporal alignment of peak flows in the observed and modeled hydrographs. The TPE
393 metric is computed as:

$$TPE = |T_{o_{max}} - T_{s_{max}}| \quad (\text{Equation 27})$$

394 Where $T_{o_{max}}$ signifies the time at which the peak flow occurs in the observed hydrograph, and $T_{s_{max}}$
395 represents the time at which the peak flow occurs in the simulated hydrograph. TPE that is measured in
396 units of time (hours), provides insight into the precision of peak flow timing. Smaller TPE values indicate
397 a superior alignment between the observed and modeled peak flow timing, while larger TPE values indicate
398 discrepancies in the temporal occurrence of peak flows.

399 The utilization of these five metrics, PFE, persistent-NSE, TPE, NSE, and RMSE, collectively provides a
400 robust and multifaceted assessment of flood prediction performance. This approach ensures that both the



401 magnitude and timing of peak flows, as well as the overall hydrograph shape, are accurately calibrated and
402 validated.

403 **2.4. Sensitivity and Uncertainty Analysis**

404 When implementing NN models, it's crucial to understand how each parameter affects the model's
405 performance or outputs. To achieve this, we systematically excluded each parameter from the model one
406 by one (the Leave-One-Out method). For each exclusion, we retrained the model without that specific
407 parameter and then tested its performance against a test dataset. This method helps in understanding which
408 parameters are most critical to the model's performance and which ones have a lesser impact. It also allows
409 us to identify any parameters that may be redundant or have little effect on the overall outcome, thus
410 potentially simplifying the model without sacrificing accuracy.

411 In this study, we utilized probabilistic approaches to quantify the uncertainty in flood prediction. This
412 method is rooted in statistical techniques employed for the estimation of unknown probability distributions,
413 with a foundation in observed data. More specifically, we leveraged the Maximum Likelihood Estimation
414 (MLE) approach, which entails the determination of parameter values that optimize the likelihood function.
415 The likelihood function quantifies the probability of parameters taking particular values, given the observed
416 realizations.

417 Within our models, we incorporated the MQL as a probabilistic error metric. MQL performs an evaluation
418 by computing the average loss for a predefined set of quantiles. This computation is grounded in the
419 absolute disparities between predicted quantiles and their corresponding observed values. The limited
420 behavior of MQL serves as an apt metric for assessing the accuracy of predictive distribution \hat{F}_t , facilitated
421 through the Continuous Ranked Probability Score (CRPS). The computation of CRPS involves a numerical
422 integration technique that discretizes quantiles and applies a left Riemann approximation for CRPS integral
423 computation. This process culminates in the averaging of these computations over uniformly spaced
424 quantiles.

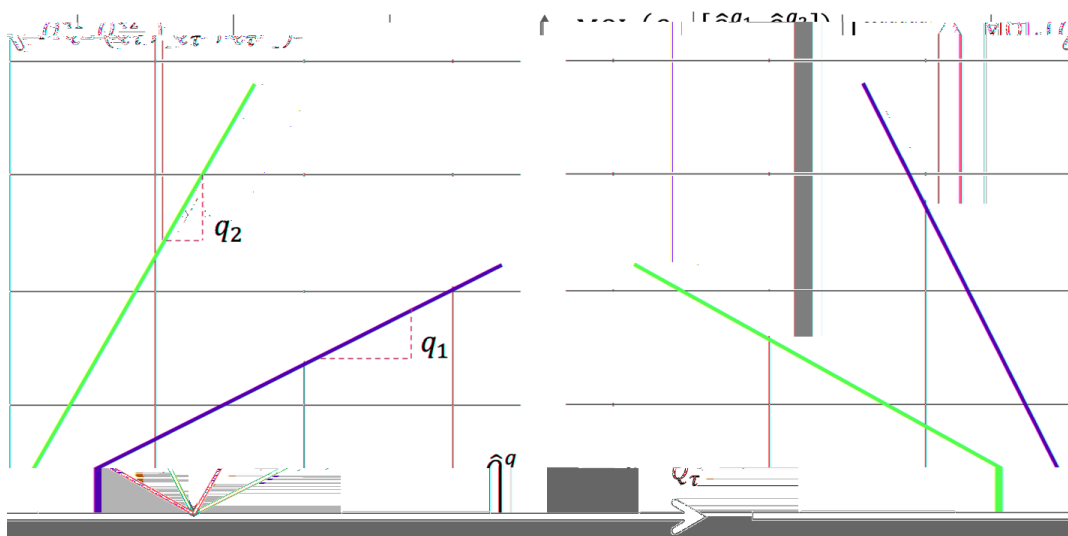
$$\text{MQL} (Q_\tau, [\hat{Q}_\tau^{q_1}, \dots, \hat{Q}_\tau^{q_i}]) = \frac{1}{n} \sum_{q_i} \text{QL} (Q_\tau, \hat{Q}_\tau^{q_i}) \quad (\text{Equation 28})$$

$$\text{CRPS} (Q_\tau, \hat{F}_t) = \int_0^1 \text{QL} (Q_\tau, \hat{Q}_\tau^{q_i}) dq \quad (\text{Equation 29})$$



$$QL(Q_\tau, \hat{Q}_\tau^q) = \frac{1}{H} \sum_{\tau=t+1}^{t+H} ((1-q)(\hat{Q}_\tau^q - Q_\tau) + q(Q_\tau - \hat{Q}_\tau^q)) \quad (\text{Equation 30})$$

425 Where Q_τ represents observed value at time τ , \hat{Q}_τ^q represents simulated value at time τ , q is the slope of the
 426 quantile loss, and H is the horizon of forecasting.



427
 428 Figure 5. The MQL function which shows loss values for different parameters of q when the true value is
 429 Q_τ .

430 Furthermore, we employed two key indices, the R-factor and the P-factor, to rigorously assess the quality
 431 of uncertainty performance in our hydrological modeling. These metrics are instrumental in quantifying the
 432 extent to which the model's predictions encompass the observed data, thereby providing valuable insights
 433 into the model's predictive accuracy and reliability.

434 The P-factor, or percentage of data within a 95PPU, is the first index used in this assessment. The P-factor
 435 quantifies the percentage of observed data that falls within the 95PPU, providing a measure of the model's
 436 predictive accuracy. The P-factor can theoretically vary from 0% to a maximum of 100%. A P-factor of
 437 100% signifies a perfect alignment between the model's predictions and the observed data within the
 438 uncertainty band. In contrast, a lower P-factor indicates a reduced ability of the model to predict data within
 439 the specified uncertainty range.



$$P - Factor = \frac{\text{Observations bracketed by 95PPU}}{\text{Number of observations}} \times 100 \quad (\text{Equation 31})$$

440 The R-factor can be computed by dividing the average width of the uncertainty band by the standard
441 deviation of the measured variable. The R-factor, with a minimum possible value of zero, provides a
442 measure of the spread of the uncertainty relative to the variability of the observed data. Theoretically, the
443 R-factor spans from 0 to infinity, and a value of zero implies that the model's predictions precisely match
444 the measured data, with the uncertainty band being very narrow in relation to the variability of the observed
445 data.

$$R - Factor = \frac{\text{Average width of 95PPU band}}{\text{Standard deviation of measured variables}} \times 100 \quad (\text{Equation 32})$$

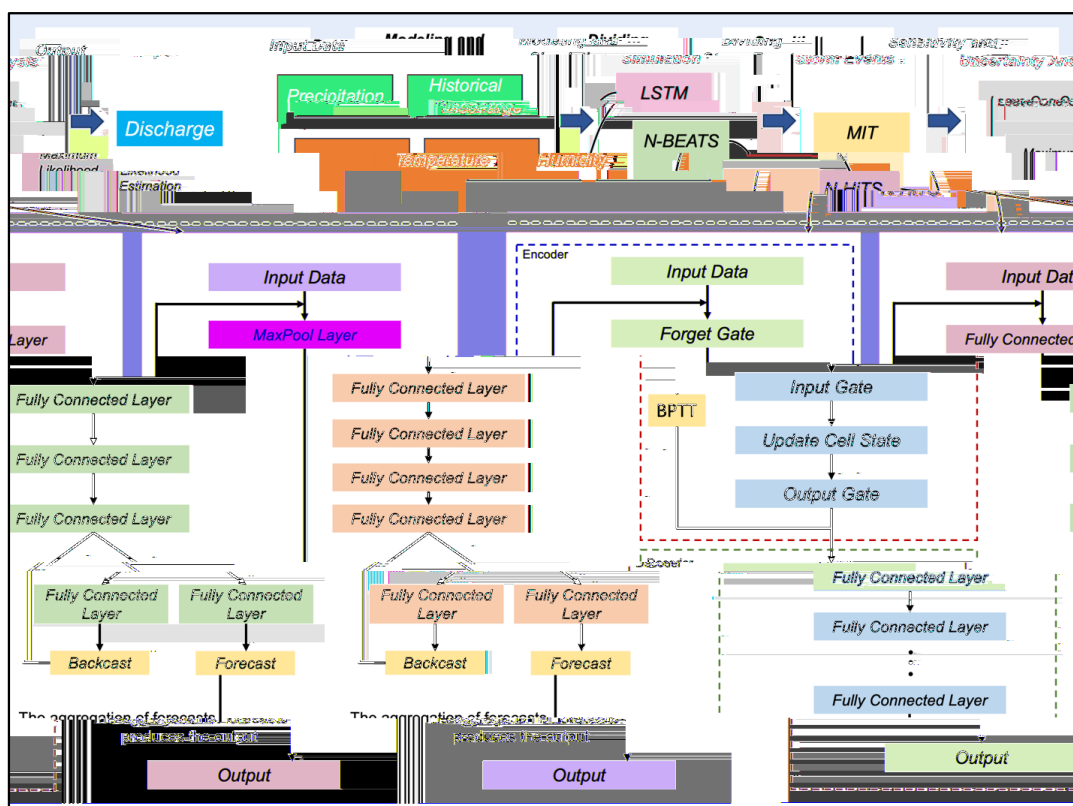
446 In practice, the quality of the model is assessed by considering the 95% prediction band with the highest P-
447 factor and the lowest R-factor. This specific band encompasses the majority of observed records, signifying
448 the model's ability to provide accurate and reliable predictions while effectively quantifying uncertainty. A
449 simulation with a P-factor of 1 and an R-factor of 0 signifies an ideal scenario where the model precisely
450 matches the measured data within the uncertainty band (Abbaspour et al., 2007).

451 Figure 6 shows the workflow of programming N-BEATS, N-HiTS, and LSTM for flood prediction. As
452 illustrated, the initial step involved cleaning and preparing the input data, which was then used to feed the
453 models. The workflow for each model and their output generation processes are depicted in Figure 6. We
454 segmented the storm events using the MIT approach, as previously described. Following this, we conducted
455 a sensitivity analysis using the Leave-One-Out method and performed uncertainty analysis using the MLE
456 approach to construct the 95PPU band. This rigorous methodology ensures a robust evaluation of model
457 performance under varying conditions and highlights the models' predictive reliability and resilience. We
458 employed the “NeuralForecast” Python package to develop the N-BEATS, N-HiTS, and LSTM models.
459 This package provides a diverse array of NN models with an emphasis on usability and robustness.

460



461



462

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

469

470 3. Results and Discussion

471 3.1. Independent Storms Delineation

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



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

482 3.2. Hyperparameter Optimization

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

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

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

502 Table 2. Optimized values for models 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	*	*

503

*Not applicable

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

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



522 average TC in the Lower Dog River watershed and Upper Dutchmans Creek watershed was found to be 19
523 and 22 hours, respectively. Through hyperparameter optimization, we determined that an input size of 24
524 hours was optimal for all the models, ensuring sufficient coverage of relevant meteorological data preceding
525 flood events.

526 **3.3. Flood Prediction and Performance Assessment**

527 In this study, we conducted a comprehensive performance evaluation of N-HiTS, N-BEATS, and
528 benchmarking LSTM models, utilizing two case studies: the Lower Dog River and the Upper Dutchmans
529 Creek watersheds. Within these case studies, we trained the models across a diverse set of storm events
530 from 01/10/2007 to 01/10/2022 (15 years) in the Lower Dog River and from 21/12/1994 to 01/10/2022 (27
531 years) in the Upper Dutchmans Creek. All algorithms were validated using flooding events that occurred
532 between 14/12/2022 and 28/03/2023. In the Dog River gauging station, two winter storms i.e., January 3rd
533 to January 5th, 2023 (Event 1) and February 17th to February 18th, 2023 (Event 2), as well as a spring flood
534 event that occurred during March 26th to March 28th, 2023 (Event 3) were selected for testing.
535 Additionally, three winter flooding events, i.e., December 14th to December 16th, 2022 (Event 4), January
536 25th and January 26th, 2023 (Event 5), and February 11th to February 13th, 2023 (Event 6), were chosen
537 to test the algorithms across the Killian Creek gauging station in the Upper Dutchmans Creek. The rainfall
538 events corresponding to these flooding events were delineated using the MIT technique discussed in Section
539 3.1.

540 Our results for the Lower Dog River case study, explicitly demonstrated the accuracy of both N-HiTS and
541 N-BEATS in generating the winter and spring flood hydrographs compared to the LSTM model across all
542 selected storm events. Although, N-HiTS prediction slightly outperformed N-BEATS during winter
543 prediction (January 3rd to January 5th, 2023). In this event, N-HiTS outperformed N-BEATS with a
544 difference of 11.6% in MAE and 20% in RMSE. The N-HiTS slight outperformance (see Tables 3 and 4)
545 is attributed to its unique structure that allows the model to discern and capture intricate patterns within the
546 data. Specifically, N-HiTS predicted flooding events hierarchically using blocks specialized in different
547 rainfall frequencies based on controlled signal projections, through expressiveness ratios, and interpolation
548 of each block. The coefficients are then used to synthesize backcast through
549 $\tilde{y}_t - L: t, l$ and forecast ($\tilde{y}_{t+1}: t + H, l$) outputs of the block as a flood value. The coefficients were locally
550 determined along the horizon, allowing N-HiTS to reconstruct nonstationary signals over time.

551 While the N-HiTS emerged as the most accurate in predicting flood hydrograph among the three models,
552 its performance was somehow comparable with N-BEATS. The N-BEATS model exhibited good
553 performance in two case studies. It consistently provided competitive results, demonstrating its capacity to



554 effectively handle diverse storm events and deliver reliable predictions. N-BEATS has a generic and
555 interpretable architecture depending on the blocks it uses. Interpretable configuration sequentially projects
556 the signal into polynomials and harmonic basis to learn trend and seasonality components while generic
557 configuration substitutes the polynomial and harmonic basis for identity basis and larger network's depth.
558 In this study, we used interpretable architecture, as it regularizes its predictions through projections into
559 harmonic and trend basis that is well-suited for flood prediction tasks. Using interpretable architecture,
560 flood prediction was aggregated in a hierarchical fashion. This enabled the building of a very deep neural
561 network with interpretable flood prediction outputs.

562 It is essential to underscore that, despite its strong performance, the N-BEATS model did not surpass the
563 N-HiTS model in terms of NSE, MAE, and RMSE for the Lower Dog River case study. Notably, the N-
564 BEATS model showcased superior results based on the PFE metric, signifying its exceptional capability in
565 accurately predicting flood peaks. However, both N-HiTS and N-BEATS models overestimated the flood
566 peak rate of Event 2 for the Lower Dog River watershed. This event, which occurred from February 17th to
567 February 18th, 2023, was flashy, short, and intense proceeded by a prior small rainfall event (from February
568 12th until February 13th) that minimized the rate of infiltration. This flash flood event caused by excessive
569 rainfall in a short period of time (<8 hours) was challenging to predict for both N-BEATS and N-HiTS
570 models. In addition, predicting the magnitude of changes in the recession curve of the third event seems to
571 be a challenge for both models. The specific part of the flood hydrograph after the precipitation event,
572 where flood diminishes during a rainless is dominated by the release of runoff from shallow aquifer systems
573 or natural storages. It seems both models showed a slight deficiency in capturing this portion of the
574 hydrograph when the rainfall amount decreases over time in the Dog River gauging station.

575 Conversely, in the Killian Creek gauging station, the N-BEATS model almost emerged as the top performer
576 in predicting the flood hydrograph based on NSE, RMSE, and PFE performance metrics (see Tables 3 and
577 4). Although, both N-BEATS and N-HiTS slightly overpredicted time to peak values for Event 5. This
578 reflects the fact that when rainfall value varies randomly around zero, it provides less to no information for
579 the algorithms to learn the fluctuations and patterns in time series data. Both N-HiTS and N-BEATS
580 provided comparable results for all events predicted in this study. N-HiTS builds upon N-BEATS by adding
581 a MaxPool layer at each block. Each block consists of an MLP layer that learns to produce coefficients for
582 the backcast and forecast outputs. This subsamples the time series and allows each stack to focus on either
583 short-term or long-term effects, depending on the pooling kernel size. Then, the partial predictions of each
584 stack are combined using hierarchical interpolation. This ability enhances N-HiTS capabilities to produce
585 drastically improved, interpretable, and computationally efficient long-horizon flood predictions.



586 In contrast, the performance of LSTM as a benchmark model lagged behind both N-HiTS and N-BEATS
587 models for all events across two case studies. Despite its extensive application in various hydrology
588 domains, the LSTM model exhibited comparatively lower accuracy when tasked with predicting flood
589 responses during different storm events. Focusing on NSE, MAE, RMSE, and PFE metrics, it is noteworthy
590 that all three models, across both case studies, consistently succeeded in capturing peak flow rates at the
591 appropriate timing. All models demonstrated commendable results with respect to the TPE metric. In most
592 scenarios, TPE revealed a value of 0, signifying that the models accurately pinpointed the peak flow rate
593 precisely at the expected time. In some instances, TPE reached a value of 1, showing a deviation of one
594 hour in predicting the peak flow time. This deviation is deemed acceptable, particularly considering the
595 utilization of short, intense rainfall for our analysis.

596 Our investigation into the performance of the three distinct forecasting models yielded compelling results
597 pertaining to their ability to generate 95PPU, as quantified by the P-factor and R-factor. These factors serve
598 as critical indicators for assessing the reliability and precision of the uncertainty bands produced by the
599 MLE. Our findings demonstrated that the N-HiTS and N-BEATS models outperformed the LSTM model
600 in mathematically defining uncertainty bands, in terms of R-factor metric. The R-factor, a crucial metric
601 for evaluating the average width of the uncertainty band, consistently favored the N-HiTS and N-BEATS
602 models over their counterparts. This finding was consistent across a diverse range of storm events. Coupling
603 MLE with the N-HiTS and N-BEATS models demonstrated superior performance in generating 95PPU
604 when assessed through the P-factor metric. The P-factor represents another vital aspect of uncertainty
605 quantification, focusing on the precision of the uncertainty bands.

606
607 Figures 8 and 9 present graphical depictions of the predicted flood with uncertainty assessment for each
608 model as well as Flow Duration Curve (FDC) across two gauging stations. As illustrated, the uncertainty
609 bands skillfully bracketed most of the observational data, reflecting the fact that MLE was successful in
610 reducing errors in flood prediction. FDC analysis also revealed that N-HiTS and N-BEATS models
611 skillfully predicted the flood hydrograph, however, both models were particularly successful in predicting
612 moderate to high flood events (1800-6000 and >6000 cfs). In the FDC plots, the x-axis denotes the
613 exceedance probability, expressed as a percentage, while the y-axis signifies flood in cubic feet per second.
614 Notably, these plots reveal distinctive patterns in the performance of the N-HiTS, N-BEATS, and LSTM
615 models. Within the lower exceedance probability range, particularly around the peak flow, the N-HiTS and
616 N-BEATS models demonstrated a clear superiority over the LSTM model, closely aligning with the
617 observed data. This observed trend is consistent when examining the corresponding hydrographs. Across
618 all events, the flood hydrographs generated by N-HiTS and N-BEATS exhibited a closer resemblance to



619 the observed data, particularly in the vicinity of the peak timing and rate, compared to the hydrographs
 620 produced by the LSTM model. These findings underscore the enhanced predictive accuracy and reliability
 621 of the N-HITS and N-BEATS models, particularly in predicting moderate to high flood events as well as
 622 critical hydrograph features such as peak flow rate and timing. The alignment of model-generated FDCs
 623 and hydrographs with observed data in the proximity of peak flow further establishes the efficacy of N-
 624 HiTS and N-BEATS in accurately reproducing the dynamics of flood generation mechanisms across two
 625 headwater streams.

626
 627

Table 3. Accuracy and uncertainty metrics for the 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
	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
	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
	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 %



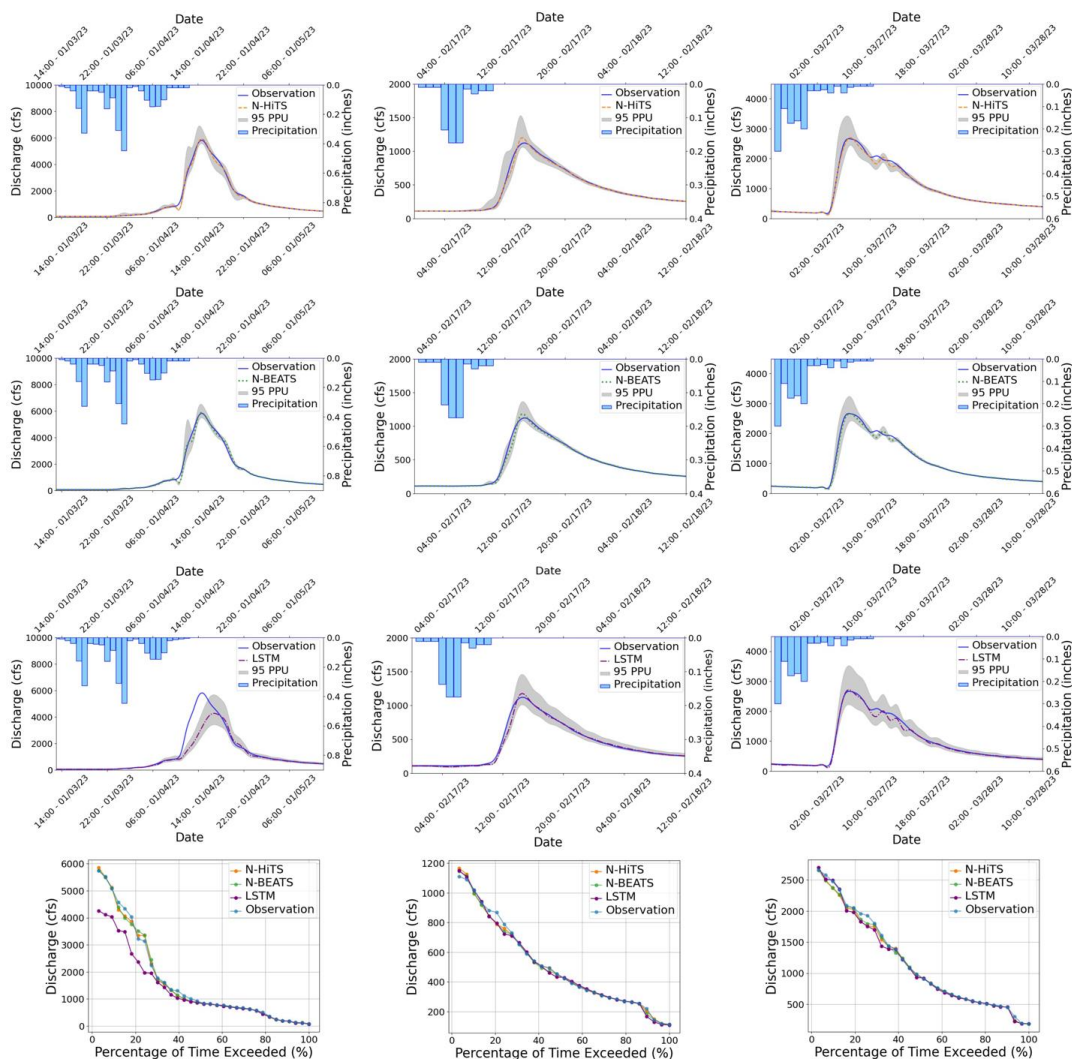
628 **R-Factor** 0.37 0.51 0.6

629 Table 4. Accuracy and uncertainty metrics for the Killian Creek flood predictions.

Model	Performance Metric	Event 4	Event 5	Event 6	
N-HiTS	NSE	99.08 %	97.13 %	99.08 %	
	Persistent-NSE				
	RMSE	28.8	46.0	19.0	
	MAE	17.9	23.8	11.5	
	PFE	0.017	0.008	0.020	
	TPE (hours)	0	0	0	
	R-Factor	0.39	0.48	0.45	
N-BEATS	NSE	99.26 %	97.36 %	98.96 %	
	Persistent-NSE				
	RMSE	25.7	44.2	20.2	
	MAE	18.3	25.9	14.0	
	PFE	0.006	0.008	0.019	
	TPE (hours)	0	0	0	
	R-Factor	0.43	0.53	0.43	
LSTM	NSE	0.952	0.892	0.935	
	Persistent-NSE				
	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	
	R-Factor	0.66	0.7	0.65	

630

631



632

Event 1

633

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

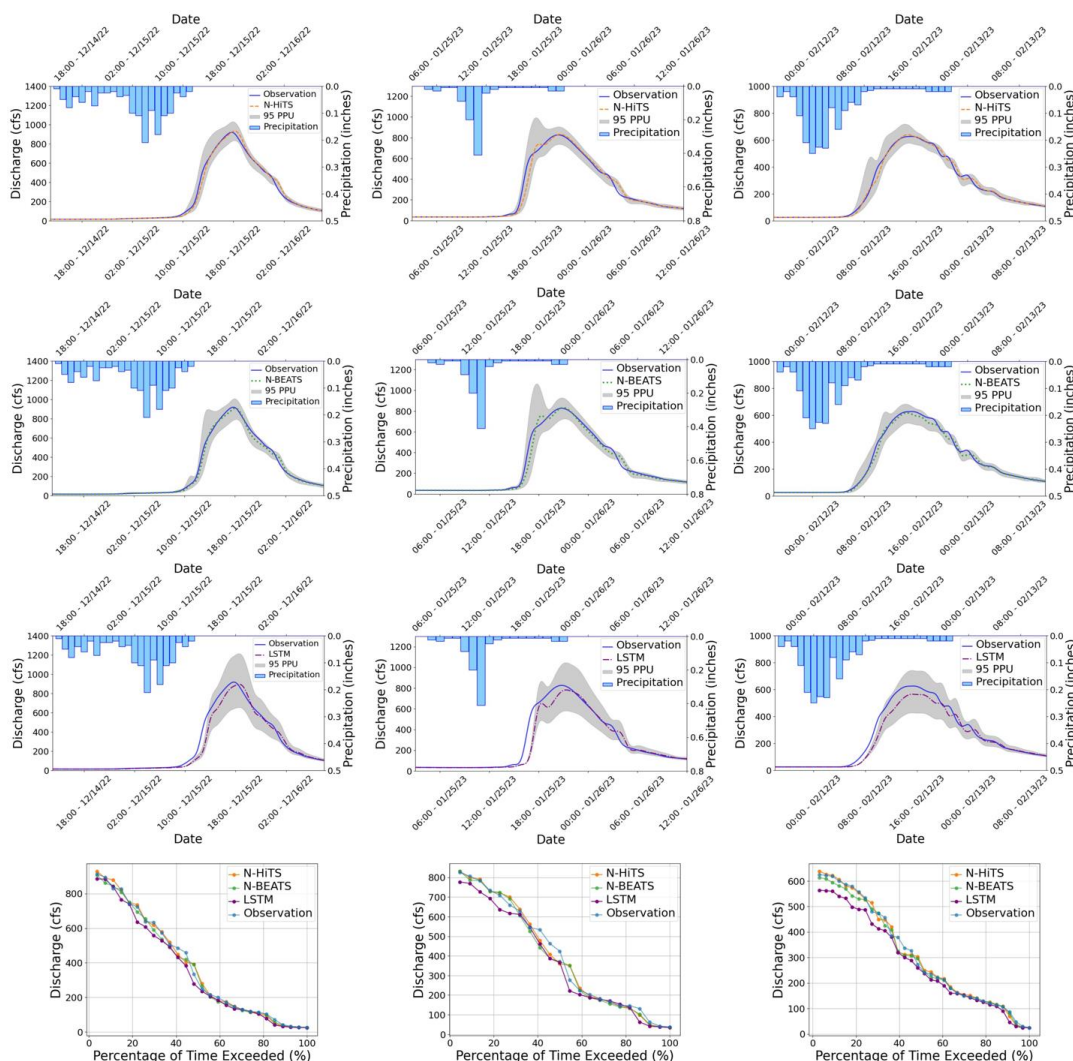
634

Event 2

635

Event 3

636



637
 638
 639

Event 4 **Event 5** **Event 6**
 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.

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



644 As detailed in Table 5, we observed a discernible trend in the R-Factor of the N-HiTS model as the input
645 size was increased. Specifically, there was a decrease in the R-Factor as the input size expanded. This trend
646 underscores the influence of input size on model performance, particularly in terms of 95PPU and accuracy.

647 Overall, uncertainty analysis revealed that the integration of MLE with N-HiTS and N-BEATS models
648 demonstrated superior performance in generating 95PPU, effectively reducing errors in flood prediction.
649 The MLE approach was more successful in reducing 95PPU bands of N-HiTS and N-BEATS models
650 compared to the LSTM, as indicated by the R-factor and P-factor. The N-BEATS model demonstrated a
651 narrower uncertainty band (lower R-factor value), while the N-HiTS model provided higher precision.
652 Furthermore, incorporating data with various sizes into the N-HiTS model led to a reduction in 95PPU and
653 an improvement in the R-Factor, highlighting the significance of input size in enhancing model accuracy
654 and reducing prediction 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

655

656 3.4. Sensitivity Analysis

657 In this study, we conducted a comprehensive sensitivity analysis of the N-HiTS, N-BEATS, and LSTM
658 models to evaluate their responsiveness to meteorological variables, specifically precipitation, humidity,
659 and temperature. The goal was to assess how the omission of input parameters impacts the overall
660 modeling performance compared to their full-variable counterparts.

661 To execute this analysis, we systematically trained each model by excluding meteorological variables one
662 or more at a time, subsequently evaluating their predictive performance using the entire testing dataset.

663 The results of our analysis indicated that N-HiTS and N-BEATS models exhibited minimal sensitivity to
664 meteorological variables, as evidenced by the negligible impact on their performance metric (NSE) upon
665 parameter exclusion.



666 Notably, as shown in Table 6, the performance of the N-HiTS model displayed a marginal deviation
 667 under variable omission, while the N-BEATS model exhibited consistent performance irrespective of the
 668 inclusion or exclusion of meteorological variables. The structure of this algorithm is based on backward
 669 and forward residual links for univariate time series point forecasting which does not take into account
 670 other parameters in the prediction task. These findings suggest that the predictive capabilities of N-HiTS
 671 and N-BEATS models predominantly rely on historical flood data, underscoring their resilience in
 672 prediction in the absence of specific meteorological inputs. This resilience to meteorological variability
 673 underscores the robustness of the N-HiTS and N-BEATS models, positioning them as viable tools and
 674 perhaps appropriate for real-time flood forecasting tasks where direct meteorological data may be limited
 675 or unavailable.

676

677 Table 6. NSE values for N-HiTS and N-BEATS models by excluding meteorological variables one or
 678 more at a time.

Model	Excluded Variables	NSE
N-HiTS	Using all variables	99.55 %
	Without Precipitation	99.34 %
	Without Humidity	99.51 %
	Without Temperature	99.49 %
	Discharge only prediction	99.3 %
N-BEATS	Using all variables	99.42 %
	Without Precipitation	99.42 %
	Without Humidity	99.42 %
	Without Temperature	99.42 %
	Discharge only prediction	99.42 %
LSTM	Using all variables	99.2 %
	Without Precipitation	97.93 %
	Without Humidity	99.13 %
	Without Temperature	98.27 %
	Discharge only prediction	97.6 %

679



680 **3.5 Computational Efficiency**

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

691 Conversely, during the testing period, the N-HiTS model emerged as the fastest and delivered the most
 692 efficient results in comparison to the other models. Notably, N-HiTS displayed a predicting time that was
 693 33% faster than LSTM and 32% faster than N-BEATS. This finding highlights the computational efficiency
 694 of the N-HiTS model in the context of predicting processes. Our experiments unveiled an interesting
 695 contrast in the computational performance of these models. While LSTM excelled in terms of training time,
 696 it lagged behind when it came to the testing period.

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

703

704 Table 7. Computational costs of N-HiTS, N-BEATS, and LSTM models in the Dog River and Killian
 705 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

706



707 **4. Conclusion**

708 This study examined multiple NN algorithms for flood prediction. We selected two headwater streams with
709 minimal human impacts to understand how NN approaches can capture flood magnitude and timing for
710 these natural systems. In conclusion, our study represents a pioneering effort in exploring and advancing
711 the application of NN algorithms, specifically the N-HiTS and N-BEATS models, in the field of flood
712 prediction. In our case studies, both N-HiTS and N-BEATS models achieved state-of-the-art results,
713 outperforming LSTM as a recurrent model. These benchmarking results are arguably a pivotal part of this
714 paper. However, the N-BEATS model slightly emerged as a powerful and interpretable tool for flood
715 prediction in most selected events.

716 In addition, the results of the experiments described above demonstrated that N-HiTS multi-rate input
717 sampling and hierarchical interpolation along with N-BEATS interpretable configuration are effective in
718 learning location-specific runoff generation behaviors. Both algorithms with an MLP-based deep neural
719 architecture with backward and forward residual links can sequentially project the data signal into
720 polynomials and harmonic basis needed to predict intense storm behaviors with varied magnitudes. The
721 innovation in this study – besides benchmarking the LSTM model for headwater streams – was to tackle
722 volatility and memory complexity challenges, by locally specializing flood sequential predictions into the
723 data signal's frequencies with interpretability, and hierarchical interpolation and pooling. Both N-HiTS and
724 N-BEATS models offered similar performance as compared with the LSTM but also offered a level of
725 interpretability about how the model learns to differentiate aspects of complex watershed-specific behaviors
726 via data. Both models also support multivariate series (and covariates) by flattening the model inputs to a
727 1-D series and reshaping the outputs to a tensor of appropriate dimensions. This approach provides
728 flexibility to handle arbitrary numbers of features. Furthermore, both N-HiTS and N-BEATS models also
729 support producing probabilistic predictions by specifying a likelihood parameter. In terms of sensitivity
730 analysis, both N-HiTS and N-BEATS models maintain consistent performance even when trained without
731 specific meteorological inputs. This resilience underscores these models' ability to generate accurate
732 predictions using historical flood data alone, making them valuable tools for flood prediction, especially in
733 data-poor watersheds or even for real-time flood prediction when near real-time meteorological inputs are
734 limited or unavailable. In terms of computational efficiency, both N-HiTS and N-BEATS are trained almost
735 at the same pace; however, N-HiTS predicted the test data much quicker than N-BEATS. Unlike N-HiTS
736 and N-BEATS, LSTM excelled in reducing training time due to its simplicity and limited number of
737 parameters.

738 Moving forward, it is worth mentioning that predicting the magnitude of the recession curve of flood
739 hydrographs was particularly challenging for all models. We argue that this is because the relation between



740 base flow and time is particularly hard to calibrate due to ground-water effluent that is controlled by
741 geological and physical conditions (vegetation, wetlands, wet meadows) in headwater streams. In addition,
742 the situations of runoff occurrence are diverse and have a high measurement variance with high frequency
743 that can make it difficult for NN algorithms to fully capture discrete representation learning on time series.
744 In future studies, it will be important to develop strategies to derive analogs to the interpretable
745 configuration as well as multi-rate input sampling, hierarchical interpolation, and backcast residual
746 connections that allow for the dynamic representation of flood times series data with different frequencies
747 and nonlinearity. A dynamic representation of flood time series is, at least in principle, possible by
748 generating additive predictions in different bands of the time-series signals, reducing memory footprint and
749 compute time, and improving architecture parsimony and accuracy. This would allow the model to “learn”
750 interpretability and hierarchical representations from raw data to reduce complexity as the information
751 flows through the network. Lastly, one could explore the idea of enhancing N-HITS and N-BEATS (or NN
752 algorithms, in general) performance with uncertainty quantification by using more robust Bayesian
753 inference such as Bayesian Model Averaging (BMA) with fixed and flexible prior distributions (see Samadi
754 et al., 2020) and/or Markov Chain Monte-Carlo optimization methods (Duane et al., 1987) addressing both
755 aleatoric and epistemic uncertainties. We leave these approaches for future discussion and exploration in
756 the context of flood neural time series prediction.

757

758 **5. Acknowledgements**

759 This research is supported by the US National Science Foundation Directorate of Engineering (Grant #
760 CMMI 2125283). The authors acknowledge and appreciate Thorsten Wagener (University of University of
761 Potsdam, Germany) discussion and feedback on this manuscript. Clemson University (USA) is
762 acknowledged for generous allotment of computing time on the Palmetto cluster.

763

764 **6. Code/Data availability**

765 The historical discharge data used in this study are from the USGS
766 (https://waterdata.usgs.gov/nwis/uv/?referred_module=sw), meteorological data from USDA
767 (<https://www.ncdc.noaa.gov/cdo-web/datatools/lcd>). We uploaded the datasets used in this
768 research to Zenodo, accessible via <https://zenodo.org/records/13342838>. For modeling, we used
769 the NeuralForecast package (Olivares et al., 2022), available at:
770 <https://github.com/Nixtla/neuralforecast>.

771



772 **7. Author contribution:** MS performed the analyses and wrote the initial draft. VS performed and
773 interpreted the results, conceptualization, writing and editing the paper, and funding acquisition. IP edited
774 the manuscript and helped with the interpretation. All the authors substantially contributed to the final draft.
775

776 **8. Competing interests:** The contact author has declared that none of the authors has any competing
777 interests.

778

779 **9. References**

780 Abbaspour, K.C., Yang, J., Maximov, I., Siber, R., Bogner, K., Mieleitner, J., Zobrist, J., Srinivasan, R.,
781 2007. Modelling hydrology and water quality in the pre-alpine/alpine Thur watershed using SWAT.
782 *Journal of Hydrology* 333, 413–430. <https://doi.org/10.1016/j.jhydrol.2006.09.014>

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

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

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

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

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

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

800 Clark, M.P., Nijssen, B., Lundquist, J.D., Kavetski, D., Rupp, D.E., Woods, R.A., Freer, J.E., Gutmann,
801 E.D., Wood, A.W., Brekke, L.D., Arnold, J.R., Gochis, D.J., Rasmussen, R.M., 2015. A unified
802 approach for process-based hydrologic modeling: 1. Modeling concept. *Water Resources Research*
803 51, 2498–2514. <https://doi.org/10.1002/2015WR017198>

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



- 806 Dasgupta, A., Arnal, L., Emerton, R., Harrigan, S., Matthews, G., Muhammad, A., O'Regan, K., Pérez-
807 Ciria, T., Valdez, E., van Osnabrugge, B., Werner, M., Buontempo, C., Cloke, H., Pappenberger,
808 F., Pechlivanidis, I.G., Prudhomme, C., Ramos, M.-H., Salamon, P., n.d. Connecting hydrological
809 modelling and forecasting from global to local scales: Perspectives from an international joint
810 virtual workshop. *Journal of Flood Risk Management* n/a, e12880.
811 <https://doi.org/10.1111/jfr3.12880>
- 812 Defontaine, T., Ricci, S., Lapeyre, C., Marchandise, A., Pape, E.L., 2023. Flood forecasting with Machine
813 Learning in a scarce data layout. *IOP Conf. Ser.: Earth Environ. Sci.* 1136, 012020.
814 <https://doi.org/10.1088/1755-1315/1136/1/012020>
- 815 Duane, S., Kennedy, A.D., Pendleton, B.J., Roweth, D., 1987. Hybrid Monte Carlo. *Physics Letters B*
816 195, 216–222. [https://doi.org/10.1016/0370-2693\(87\)91197-X](https://doi.org/10.1016/0370-2693(87)91197-X)
- 817 Erikson, L.H., Espejo, A., Barnard, P.L., Serafin, K.A., Hegermiller, C.A., O'Neill, A., Ruggiero, P.,
818 Limber, P.W., Mendez, F.J., 2018. Identification of storm events and contiguous coastal sections
819 for deterministic modeling of extreme coastal flood events in response to climate change. *Coastal*
820 *Engineering* 140, 316–330. <https://doi.org/10.1016/j.coastaleng.2018.08.003>
- 821 Evin, G., Le Lay, M., Fouchier, C., Mas, A., Colleoni, F., Penot, D., Garambois, P.-A., Laurantin, O.,
822 2023. Evaluation of hydrological models on small mountainous catchments: impact of the
823 meteorological forcings. <https://doi.org/10.5194/egusphere-2023-845>
- 824 Fan, C., Zhang, Y., Pan, Y., Li, X., Zhang, C., Yuan, R., Wu, D., Wang, W., Pei, J., Huang, H., 2019.
825 Multi-Horizon Time Series Forecasting with Temporal Attention Learning, in: *Proceedings of the*
826 *25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '19.*
827 *Association for Computing Machinery, New York, NY, USA*, pp. 2527–2535.
828 <https://doi.org/10.1145/3292500.3330662>
- 829 Fang, K., Kifer, D., Lawson, K., Shen, C., 2020. Evaluating the Potential and Challenges of an
830 Uncertainty Quantification Method for Long Short-Term Memory Models for Soil Moisture
831 Predictions. *Water Resources Research* 56, e2020WR028095.
832 <https://doi.org/10.1029/2020WR028095>
- 833 Global assessment report on disaster risk reduction 2015 | UNDRR, 2015. URL:
834 <http://www.undrr.org/publication/global-assessment-report-disaster-risk-reduction-2015> (accessed
835 6.5.24).
- 836 Gotvald, A.J., 2010, *Historic flooding in Georgia, 2009: U.S. Geological Survey Open-File Report 2010–*
837 *1230*, 19 p.
- 838 Hochreiter, S., Younger, A.S., Conwell, P.R., 2001. Learning to Learn Using Gradient Descent, in:
839 Dorffner, G., Bischof, H., Hornik, K. (Eds.), *Artificial Neural Networks — ICANN 2001*. Springer,
840 Berlin, Heidelberg, pp. 87–94. https://doi.org/10.1007/3-540-44668-0_13



- 841 Hsu, K., Gupta, H.V., Sorooshian, S., 1995. Artificial Neural Network Modeling of the Rainfall-Runoff
842 Process. *Water Resources Research* 31, 2517–2530. <https://doi.org/10.1029/95WR01955>
- 843 Jonkman, S.N., 2005. Global Perspectives on Loss of Human Life Caused by Floods. *Nat Hazards* 34,
844 151–175. <https://doi.org/10.1007/s11069-004-8891-3>
- 845 Kingma, D.P., Ba, J., 2017. Adam: A Method for Stochastic Optimization.
846 <https://doi.org/10.48550/arXiv.1412.6980>
- 847 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., Herrnegger, M., 2018. Rainfall–runoff modelling using
848 Long Short-Term Memory (LSTM) networks. *Hydrology and Earth System Sciences* 22, 6005–
849 6022. <https://doi.org/10.5194/hess-22-6005-2018>
- 850 Lim, B., Arık, S.Ö., Loeff, N., Pfister, T., 2021. Temporal Fusion Transformers for interpretable multi-
851 horizon time series forecasting. *International Journal of Forecasting* 37, 1748–1764.
852 <https://doi.org/10.1016/j.ijforecast.2021.03.012>
- 853 Lobligeois, F., Andréassian, V., Perrin, C., Tabary, P., Loumagne, C., 2014. When does higher spatial
854 resolution rainfall information improve streamflow simulation? An evaluation using 3620 flood
855 events. *Hydrology and Earth System Sciences* 18, 575–594. [https://doi.org/10.5194/hess-18-575-](https://doi.org/10.5194/hess-18-575-2014)
856 2014
- 857 MacDonald, L.H., Coe, D., 2007. Influence of Headwater Streams on Downstream Reaches in Forested
858 Areas. *Forest Science* 53, 148–168. <https://doi.org/10.1093/forestscience/53.2.148>
- 859 Martinaitis, S.M., Wilson, K.A., Yussouf, N., Gourley, J.J., Vergara, H., Meyer, T.C., Heinselman, P.L.,
860 Gerard, A., Berry, K.L., Vergara, A. and Monroe, J., 2023. A path toward short-term probabilistic
861 flash flood prediction. *Bulletin of the American Meteorological Society*, 104(3), pp.E585-E605.
- 862 McCallum, B.E., and Gotvald, A.J., 2010, Historic flooding in northern Georgia, September 16–22, 2009:
863 U.S. Geological Survey Fact Sheet 2010–3061, 4 p.
- 864 McCuen, R.H., Knight, Z., Cutter, A.G., 2006. Evaluation of the Nash–Sutcliffe Efficiency Index. *Journal*
865 *of Hydrologic Engineering* 11, 597–602. [https://doi.org/10.1061/\(ASCE\)1084-](https://doi.org/10.1061/(ASCE)1084-0699(2006)11:6(597))
866 0699(2006)11:6(597)
- 867 Munn, M., Sheibley, R., Waite, I., Meador, M., 2020. Understanding the relationship between stream
868 metabolism and biological assemblages. *Freshwater Science* 39, 680–692.
869 <https://doi.org/10.1086/711690>
- 870 Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual models part I — A discussion
871 of principles. *Journal of Hydrology* 10, 282–290. [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6)
- 872 Nevo, S., Morin, E., Gerzi Rosenthal, A., Metzger, A., Barshai, C., Weitzner, D., Voloshin, D., Kratzert,
873 F., Elidan, G., Dror, G., Begelman, G., Nearing, G., Shalev, G., Noga, H., Shavitt, I., Yuklea, L.,
874 Royz, M., Giladi, N., Peled Levi, N., Reich, O., Gilon, O., Maor, R., Timnat, S., Shechter, T.,



- 875 Anisimov, V., Gigi, Y., Levin, Y., Moshe, Z., Ben-Haim, Z., Hassidim, A., Matias, Y., 2022. Flood
876 forecasting with machine learning models in an operational framework. Hydrology and Earth
877 System Sciences 26, 4013–4032. <https://doi.org/10.5194/hess-26-4013-2022>
- 878 NRCS (2009). Part 630 Hydrology National Engineering Handbook, Chapter 15: Time of Concentration.
- 879 Olivares, K. G., Challú, C., Garza, F., Mergenthaler Canseco, M., & Dubrawski, A. (2022).
880 NeuralForecast: User friendly state-of-the-art neural forecasting models. PyCon Salt Lake City,
881 Utah, US 2022. Retrieved from <https://github.com/Nixtla/neuralforecast>
- 882 Olivares, K.G., Meetei, O.N., Ma, R., Reddy, R., Cao, M., Dicker, L., 2024. Probabilistic hierarchical
883 forecasting with deep Poisson mixtures. International Journal of Forecasting 40, 470–489.
884 <https://doi.org/10.1016/j.ijforecast.2023.04.007>
- 885 Oreshkin, B.N., Carpov, D., Chapados, N., Bengio, Y., 2020. N-BEATS: Neural basis expansion analysis
886 for interpretable time series forecasting. <https://doi.org/10.48550/arXiv.1905.10437>
- 887 Pally, R.J., Samadi, V., 2021. Application of image processing and convolutional neural networks for
888 flood image classification and semantic segmentation. Environmental Modelling & Software 148,
889 105285. <https://doi.org/10.1016/j.envsoft.2021.105285>
- 890 Palmer, T.N., 2012. Towards the probabilistic Earth-system simulator: a vision for the future of climate
891 and weather prediction. Quarterly Journal of the Royal Meteorological Society 138, 841–861.
892 <https://doi.org/10.1002/qj.1923>
- 893 Park, K., Lee, E.H., 2024. Urban flood vulnerability analysis and prediction based on the land use using
894 Deep Neural Network. International Journal of Disaster Risk Reduction 101, 104231.
895 <https://doi.org/10.1016/j.ijdrr.2023.104231>
- 896 Pourreza-Bilondi, M., Samadi, S.Z., Akhoond-Ali, A.-M., Ghahraman, B., 2017. Reliability of Semiarid
897 Flash Flood Modeling Using Bayesian Framework. Journal of Hydrologic Engineering 22,
898 05016039. [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0001482](https://doi.org/10.1061/(ASCE)HE.1943-5584.0001482)
- 899 Refsgaard, J.C., Stisen, S., Koch, J., 2022. Hydrological process knowledge in catchment modelling –
900 Lessons and perspectives from 60 years development. Hydrological Processes 36, e14463.
901 <https://doi.org/10.1002/hyp.14463>
- 902 Roelvink, D., Reniers, A., van Dongeren, A., van Thiel de Vries, J., McCall, R., Lescinski, J., 2009.
903 Modelling storm impacts on beaches, dunes and barrier islands. Coastal Engineering 56, 1133–
904 1152. <https://doi.org/10.1016/j.coastaleng.2009.08.006>
- 905 Russo, S., Perraudin, N., Stalder, S., Perez-Cruz, F., Leitao, J.P., Obozinski, G., Wegner, J.D., 2023. An
906 evaluation of deep learning models for predicting water depth evolution in urban floods.
907 <https://doi.org/10.48550/arXiv.2302.10062>



- 908 Safaei-Moghadam, A., Tarboton, D., Minsker, B., 2023. Estimating the likelihood of roadway pluvial
909 flood based on crowdsourced traffic data and depression-based DEM analysis. *Natural Hazards and*
910 *Earth System Sciences* 23, 1–19. <https://doi.org/10.5194/nhess-23-1-2023>
- 911 Saksena, S., Dey, S., Merwade, V., Singhofen, P.J., 2020. A Computationally Efficient and Physically
912 Based Approach for Urban Flood Modeling Using a Flexible Spatiotemporal Structure. *Water*
913 *Resources Research* 56, e2019WR025769. <https://doi.org/10.1029/2019WR025769>
- 914 Samadi, S., Pourreza-Bilondi, M., Wilson, C. a. M.E., Hitchcock, D.B., 2020. Bayesian Model Averaging
915 With Fixed and Flexible Priors: Theory, Concepts, and Calibration Experiments for Rainfall-Runoff
916 Modeling. *Journal of Advances in Modeling Earth Systems* 12, e2019MS001924.
917 <https://doi.org/10.1029/2019MS001924>
- 918 Scott, J., n.d. Widespread Flooding After Severe Storms - WCCB Charlotte's CW. Available at:
919 <https://www.wccbcharlotte.com/2020/02/08/widespread-flooding-after-severe-storms/> (accessed
920 6.11.24).
- 921 Sukovich, E.M., Ralph, F.M., Barthold, F.E., Reynolds, D.W., Novak, D.R., 2014. Extreme Quantitative
922 Precipitation Forecast Performance at the Weather Prediction Center from 2001 to 2011. *Weather*
923 *and Forecasting* 29, 894–911. <https://doi.org/10.1175/WAF-D-13-00061.1>
- 924 Tabas, S.S., Samadi, S., 2022. Variational Bayesian dropout with a Gaussian prior for recurrent neural
925 networks application in rainfall–runoff modeling. *Environ. Res. Lett.* 17, 065012.
926 <https://doi.org/10.1088/1748-9326/ac7247>
- 927 Thompson, C.M., Frazier, T.G., 2014. Deterministic and probabilistic flood modeling for contemporary
928 and future coastal and inland precipitation inundation. *Applied Geography* 50, 1–14.
929 <https://doi.org/10.1016/j.apgeog.2014.01.013>
- 930 Tiwari, M.K., Chatterjee, C., 2010. Development of an accurate and reliable hourly flood forecasting
931 model using wavelet-bootstrap-ANN (WBANN) hybrid approach. *Journal of Hydrology* 394, 458–
932 470. <https://doi.org/10.1016/j.jhydrol.2010.10.001>
- 933 Watershed Report | Office of Water | US EPA, n.d. Available at:
934 <https://watersgeo.epa.gov/watershedreport/?comid=9224629> (accessed 6.9.24).
- 935 Wee, G., Chang, L.-C., Chang, F.-J., Mat Amin, M.Z., 2023. A flood Impact-Based forecasting system by
936 fuzzy inference techniques. *Journal of Hydrology* 625, 130117.
937 <https://doi.org/10.1016/j.jhydrol.2023.130117>
- 938 Windheuser, L., Karanjit, R., Pally, R., Samadi, S., Hubig, N.C., 2023. An End-To-End Flood Stage
939 Prediction System Using Deep Neural Networks. *Earth and Space Science* 10, e2022EA002385.
940 <https://doi.org/10.1029/2022EA002385>



- 941 Zhang, L., Qin, H., Mao, J., Cao, X., Fu, G., 2023. High temporal resolution urban flood prediction using
942 attention-based LSTM models. *Journal of Hydrology* 620, 129499.
943 <https://doi.org/10.1016/j.jhydrol.2023.129499>
- 944 Zhang, Y., Pan, D., Griensven, J.V., Yang, S.X., Gharabaghi, B., 2023. Intelligent flood forecasting and
945 warning: a survey. *ir* 3, 190–212. <https://doi.org/10.20517/ir.2023.12>
- 946 Zou, Y., Wang, J., Lei, P., Li, Y., 2023. A novel multi-step ahead forecasting model for flood based on time
947 residual LSTM. *Journal of Hydrology* 620, 129521. <https://doi.org/10.1016/j.jhydrol.2023.129521>
- 948