



1 **Deep learning of flood forecasting by considering interpretability** 2 **and physical constraints**

3 Ting Zhang *, Ran Zhang, Jianzhu Li, Ping Feng

4 State Key Laboratory of Hydraulic Engineering Intelligent Construction and Operation, Tianjin
5 University, Tianjin 300072, China

6 Corresponding author: Ting Zhang (zhangting_hydro@tju.edu.cn)

7 **ABSTRACT**

8 Deep learning models have been proven to be effective in flood forecasting by leveraging the
9 rich time-series information in the data. However, their limited interpretability and lack of physical
10 mechanisms remain significant challenges. To address these limitations, this study introduces a
11 novel model called PHY-FTMA-LSTM, which combines the feature-time-based multi-head
12 attention mechanism with physical constraints. The PHY-FTMA-LSTM model takes four essential
13 features of runoff, rainfall, evapotranspiration, and initial soil moisture as inputs to forecast floods
14 in the Luan River Basin with a lead time of 1-6 h. It emphasizes the significance of relevant factors
15 in the input features and historical moments through the feature-time attention module. Furthermore,
16 the model enhances physical consistency by considering the monotonic relationship between the
17 input variables and the output results. The results demonstrate that the PHY-FTMA-LSTM in most
18 cases outperforms the original LSTM, the feature-time-based attention LSTM (FTA-LSTM), and
19 the feature-time-based multi-head attention LSTM (FTMA-LSTM). For a lead time of $t+1$, the
20 model achieves an NSE of 0.988, with KGE and R^2 of 0.984 and 0.988. The NSE, KGE, and R^2 also
21 reach 0.908, 0.905, and 0.911 for a lead time of $t+6$. The proposed PHY-FTMA-LSTM model
22 achieves excellent prediction accuracy, offering valuable insights for enhancing interpretability and
23 physical consistency in deep learning approaches.

24 **Keywords:** Deep learning; Flood forecasting; Physical constraints; Attention mechanism

25 **1. Introduction**

26 Floods are one of the most common and destructive natural hazards, posing a great threat to
27 human life, infrastructure, and socio-economic conditions (Kellens et al., 2013; Mourato et al.,
28 2021). Building reliable and accurate flood forecasting models is the foundation for sustainable



29 flood risk management with a focus on prevention and protection, and is one of the most challenging
30 tasks in hydrological forecasting (Birkholz et al., 2014; Zhang et al., 2016).

31 Traditional hydrological models simulate hydrological processes such as rainfall runoff with a
32 clear physical meaning, but their construction often demands rich hydro-meteorological data and
33 subsurface information. Additionally, the large number of parameters involved poses challenges in
34 determining their values, limiting their practical applicability (Chen et al., 2011). In contrast, data-
35 driven machine learning (ML) models, which do not rely on explicit consideration of the physical
36 mechanisms governing hydrological processes and only analyze the statistical relationships between
37 inputs and outputs, have been widely used in hydrology in recent years (Lima et al., 2016; Yang et
38 al., 2020; Yu et al., 2006; Zhu et al., 2005). Among them, deep learning (DL) models with multiple
39 hidden layers have demonstrated significant advantages, including convolutional neural networks
40 (CNNs), recurrent neural networks (RNNs), and their variants such as long short-term memory
41 neural networks (LSTMs), and gated recurrent units (GRUs). LSTM, a type of RNN, is specifically
42 designed for learning long-term dependencies, and its architectural enhancements effectively
43 address issues such as gradient disappearance and explosion that are inherent to traditional RNNs.
44 Consequently, LSTM has emerged as a highly favored model in flood forecasting (Cui et al., 2021a;
45 Kao et al., 2020; Luppichini et al., 2022; Lv et al., 2020).

46 The DL models, with their powerful characterization capabilities, excel in fitting observations
47 and have high prediction accuracy for hydrological problems such as flood forecasting, but they still
48 have limitations. First, the interpretability of DL models is poor (Nearing et al., 2021). The inherent
49 black-box nature of DL models makes it difficult to understand the significance of model parameters
50 and the decision-making process. The attention mechanism is an approach to enhance the
51 interpretability of DL models (Vaswani et al., 2017). Attention allows for the interpretation of
52 feature importance by selectively emphasizing critical information from a multitude of input
53 variables through attention weights. Moreover, attention weights can be visualized to gain insights
54 into the underlying reasoning behind the model's predictions. The attention mechanism has been
55 successfully applied in various domains. Song et al. (2017) proposed an end-to-end spatio-temporal
56 attention model for recognizing human actions from skeleton data, selectively attending to
57 distinguishable joints within each frame of the input, and assigning different levels of attention to



58 the output of different frames. Zhang et al. (2021) constructed an anomaly structure by incorporating
59 spatial attention and channel attention modules, which facilitated the creation of feature spaces
60 characterized by high compactness within the same class and separation between different classes,
61 resulting in the accurate classification of floral images. As for hydrological forecasting, Wang et al.
62 (2023) introduced an improved spatio-temporal attention mechanism model (STA-LSTM) for
63 predicting river water levels. By visualizing attention weights, they discovered that the hydrological
64 station closer to the outlet had greater influence, while the temporal weights decreased with
65 increasing historical moments. However, it should be noted that the discussed model (STA-LSTM)
66 considers only a single historical water level as input, neglecting the potential influence of other
67 relevant input features on the final prediction. This limitation underscores the need for further
68 research and development to explore the incorporation of multiple input features in attention
69 mechanisms for more comprehensive and accurate models.

70 Second, the DL models lack physical mechanisms. DL models primarily focus on establishing
71 a mapping relationship between inputs and outputs, overlooking the underlying physical
72 connections between them (Jiang et al., 2020). Consequently, the prediction results obtained from
73 DL models may be physically inconsistent or unreliable due to extrapolation or observation bias
74 (Reichstein et al., 2019). To address this limitation, researchers have proposed incorporating
75 physical constraints into the loss function, which serves as the optimization objective of DL models.
76 By adding physical theory as a priori knowledge, the models can be constrained to generate outputs
77 that are consistent with the underlying physical principles, thereby enhancing their physical
78 consistency. Several studies have explored this approach in different contexts. Read et al. (2019)
79 chose the law of energy conservation as a physical constraint in temperature simulation to build a
80 lake water temperature prediction model that conforms to physical theory. Wang et al. (2020)
81 proposed a theory-guided neural network (TgNN) framework for groundwater flow that
82 incorporates control equations, boundary conditions, initial conditions, and expert knowledge as
83 additional terms in the loss function to guide the training process. Xie et al. (2021) considered
84 extreme storm events, long-duration rainless events, and rainfall-runoff monotonic relationships in
85 the rainfall-runoff process at a daily scale and constrained LSTM with these three physical
86 mechanisms to improve the physical interpretability.



87 Moreover, the current inputs for the DL models in flood forecasting are mainly historical runoff,
88 rainfall, and evapotranspiration (Leedal et al., 2013; Rahimzad et al., 2021; Wan et al., 2019), but
89 the initial soil moisture is also a crucial parameter, particularly for arid watersheds (Grillakis et al.,
90 2016). The initial soil moisture directly affects the soil infiltration capacity, water input and output
91 from the soil, and ultimately, the flooding process. Therefore, the paper also explores the effect of
92 initial soil moisture on flood forecasting through the attention weight visualization matrix.

93 Based on the above research, this paper proposes a combined feature-time multi-head attention
94 mechanism and physical constraints model for flood forecasting, named PHY-FTMA-LSTM. The
95 main contributions of this work are outlined as follows: (1) The initial soil moisture in the watershed
96 is introduced as an input, alongside historical runoff, rainfall, and evapotranspiration, these four
97 input features are considered to investigate their influence on the flooding process. (2) The dual
98 attention module of features and time and multiple attention heads are used. The resulting attention
99 weight matrix is visualized to enhance the interpretability of the model, providing insights into the
100 importance of different features and time dynamics. (3) The physical constraints of flood forecasting
101 are combined with the DL models at hourly scales to enhance the physical consistency of the model.
102 By optimizing the loss function, the model incorporates the monotonic relationship between rainfall,
103 evapotranspiration, initial soil moisture, and runoff during the flooding process. This integration
104 ensures that the output aligns with physical laws.

105 The novelty of this study is that, for the first time, the attention mechanism and physical
106 constraints are simultaneously incorporated into the DL model based on the hourly scale, and the
107 important parameter of soil moisture content is added as input to forecast flood with a lead time of
108 1~6h in Luan River Basin in China as an example, which improves the prediction performance of
109 flood forecasting models while enhancing interpretability and physical law consistency. The
110 proposed PHY-FTMA-LSTM can effectively leverage key input information and produce prediction
111 results that conform to the monotonicity constraints on the water balance.

112 **2. Methods**

113 To increase the interpretability and physical consistency of DL models in flood forecasting,
114 this paper establishes a PHY-FTMA-LSTM model that combines the feature-time-based multi-head
115 attention mechanism with physical constraints (Fig. 1(a)). The attention mechanism consists of a



116 dual module: feature-based attention and time-based attention. In the feature-based attention module,
117 the model generates a feature-based attention matrix that assigns different weights to the input
118 features based on their importance. Similarly, the time-based attention module generates a time-
119 based attention matrix that assigns different weights to historical moments. By taking the dot product
120 of these two matrices, the model generates the feature-time-based attention matrix (Fig. 1(b)). To
121 enhance the modeling capability, the multi-head attention mechanism is utilized. Multiple attention
122 heads are computed in parallel, and their outputs are averaged to balance the influence of each
123 subhead. The attention weight matrix is then multiplied with the input matrix, resulting in the output
124 of the feature-time-based multi-head attention layer (Fig. 1 (c)). In addition, the physical constraints
125 of the hydrological cycle process are added to the loss function to make the output conform to the
126 physical laws. And the model is compared with the original LSTM, the feature-time-based attention
127 LSTM (FTA-LSTM), and the feature-time-based multi-head attention LSTM (FTMA-LSTM).

128 2.1. Long short-term memory neural network (LSTM)

129 The LSTM model aims to alleviate the weaknesses of ordinary RNNs in handling long-time
130 dynamics (Zhao et al., 2017). Different from the circular structure of the RNN hidden layer, the
131 hidden layer of the LSTM introduces the memory cell, which consists of an input gate, forget gate,
132 and output gate to selectively remember and forget the input data, and its structure is shown in Fig.
133 1(d). The inputs at time t include the input information x_t at t , the hidden layer state h_{t-1} , and the cell
134 state c_{t-1} at $t-1$. First, the forget gate determines the extent to which cell state c_{t-1} is discarded. Next,
135 the input gate decides how much of the current external information x_t to retain and generates the
136 candidate cell state \bar{c}_t . Then, c_t is updated based on the results of the forget and input gate. Finally,
137 the output gate decides which state features of c_t are output and generates the hidden layer state
138 variable h_t (Duan et al., 1992). The above process can be expressed as follows:

$$139 \quad f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (1)$$

$$140 \quad i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (2)$$

$$141 \quad \bar{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \quad (3)$$

$$142 \quad c_t = c_{t-1} \odot f_t + \bar{c}_t \odot i_t \quad (4)$$



143
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (5)$$

144
$$h_t = \sigma \odot \tanh(c_t) \quad (6)$$

145 where W_f , W_i , W_c , W_o are the weight vectors of the three gates and the gating unit, respectively.
146 Similarly, b_f , b_i , b_c , b_o are the bias vectors. σ is the Sigmoid activation function. \tanh is the hyperbolic
147 tangent activation function. \odot denotes the vector element product.

148 2.2. Attention mechanism

149 The attention mechanism is inspired by the concept of human visual selective attention, which
150 helps neural networks focus on important information while disregarding irrelevant details, thereby
151 establishing connections between inputs and outputs (Brauwert & Frasincar, 2023; Niu et al., 2021).
152 By incorporating the attention mechanism, the model can allocate varying degrees of attention to
153 different historical moments or feature vectors within the input sequence. This enables the model to
154 automatically identify and prioritize the most relevant input information, leading to more accurate
155 modeling of flood causes and trends. Ultimately, this improves the accuracy of flood prediction
156 results and enhances the interpretability of the model.

157 In this study, a soft attention module is introduced before the original LSTM's input. This
158 module calculates attention weight matrices separately for input features and historical moments
159 and then combines them to produce a feature-time attention weight matrix.

160 The feature-based attention module can focus on the effects of different features on predicted
161 floods and improve the model's attention to important features. In this paper, the input features are
162 runoff, rainfall, evapotranspiration, and initial soil moisture. Let the input be a two-dimensional
163 matrix $X \in R^{k \times n}$, where k and n denote the number of input features and the number of historical
164 moments, respectively, then the input matrix at time t can be regarded as n k -dimensional vectors
165 $X_t = [x'_1, x'_2, \dots, x'_k]_{1 \times k}^T$. The input features at each time step are normalized using the softmax function
166 (Eq. (7) and Eq. (8)). The attention weight matrix based on the input features is obtained by
167 synthesizing the feature weights of all historical moments.

168
$$\alpha_i^t = \text{softmax}(x_i^t) = \frac{e^{-x_i^t}}{\sum_{i=1}^k e^{-x_i^t}} \quad (7)$$



169
$$\alpha_i = [\alpha'_1, \alpha'_2, \dots, \alpha'_k]_{1 \times k}^T \quad (8)$$

170 where α'_i is the weight of the i th feature, and $\sum_{i=1}^k \alpha'_i = 1$.

171 The time-based attention module allows simulating the relationship between different time
 172 steps, focusing on the more important historical moments. The input matrix of features can be
 173 viewed as $X_k = [x_k^{t-n-1}, x_k^{t-n-2}, \dots, x_k^t]_{1 \times n}$, and the same softmax function (Eq. (9)) is used to generate
 174 the time-based attention weights (Eq. (10)), and the time weights of all features are synthesized to
 175 be the attention weight matrix based on historical moments.

176
$$\beta_k^i = \text{softmax}(x_k^i) = \frac{e^{-x_k^i}}{\sum_{i=1}^n e^{-x_k^i}} \quad (9)$$

177
$$\beta_k = [\beta_1, \beta_2, \dots, \beta_n]_{1 \times n} \quad (10)$$

178 where β_k^i is the weight of the i th time step, and $\sum_{i=1}^n \beta_k^i = 1$. Finally, the above two weight matrices
 179 are multiplied element by element to obtain the attention weight matrix that focuses on both the
 180 input features and historical moments (Eq. (11)).

181
$$FTA = FA \odot TA^T = \begin{bmatrix} \alpha_1^{t-n-1} \beta_1^{t-n-1} & \dots & \alpha_1^t \beta_1^t \\ \vdots & & \vdots \\ \alpha_k^{t-n-1} \beta_k^{t-n-1} & \dots & \alpha_k^t \beta_k^t \end{bmatrix}_{k \times n} \quad (11)$$

182 To enhance model expressiveness and interpretability, this study also employs a multi-head
 183 attention mechanism. This mechanism involves passing input sequences through m independent
 184 attention heads in parallel. Each head can be seen as a distinct representation space, enabling the
 185 model to concurrently focus on different parts of the input. As a result, the model becomes more
 186 capable of capturing the intricate relationships between inputs and gaining a deeper understanding
 187 of the input data.

188 The multi-head attention mechanism computes m sets of attention coefficients based on the
 189 number of heads, adds the output tensor of the attention heads using the Add function, and then
 190 balances the effects of different sub-heads by averaging operations. Finally, the average output
 191 tensor is multiplied by the input to get the final output, which makes the attention head weights more
 192 discriminative and better captures the relationship between sequences. The feature-time-based



193 multi-head attention weight matrix is as follows:

$$194 \quad FTMA = \frac{1}{M} \begin{bmatrix} \sum_{m=1}^M \alpha_1^{t-n-1} \beta_1^{t-n-1} & \dots & \sum_{m=1}^M \alpha_1^t \beta_1^t \\ \vdots & & \vdots \\ \sum_{m=1}^M \alpha_k^{t-n-1} \beta_k^{t-n-1} & \dots & \sum_{m=1}^M \alpha_k^t \beta_k^t \end{bmatrix}_{k \times n} \quad (12)$$

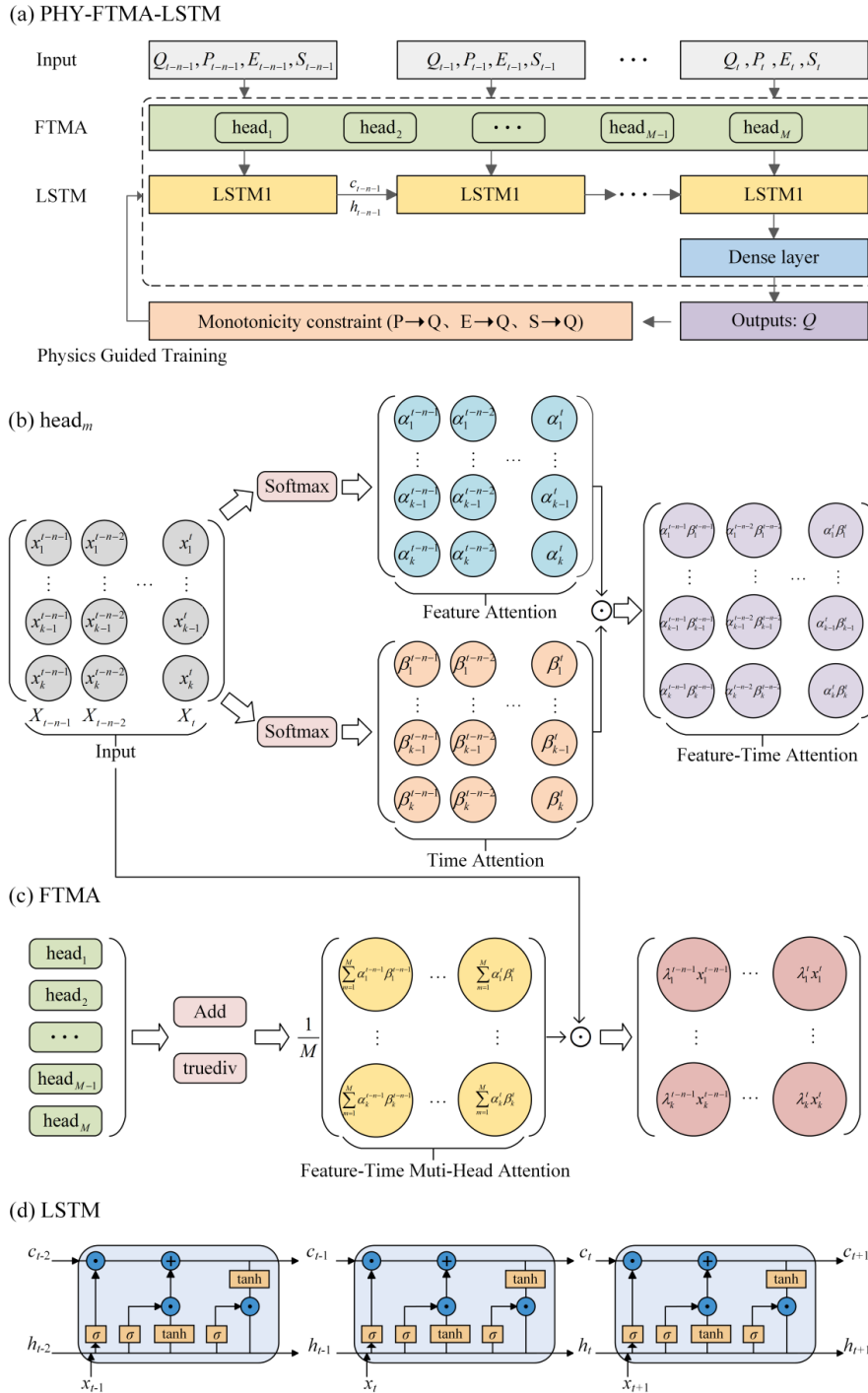
195 where M represents the number of attention heads.

196 2.3. Physical constraints

197 The LSTM is a black-box model that ignores complex physical processes, making it difficult
198 to maintain consistency with the basic principles of flood forecasting (Yokoo et al., 2022). To
199 overcome this limitation, the physical constraints can be combined with the DL models to enhance
200 the physical consistency by modifying the model loss function and transforming the prior
201 knowledge of flood forecasting into the penalty term of the loss function. A soft penalty is often
202 utilized to enforce constraints on the model's behavior (Karniadakis et al., 2021), ensuring
203 adherence to physical principles such as conservation and monotonicity.

204 In the DL models for flood forecasting, the occurrence of flooding due to heavy rainfall is
205 influenced by various factors, including rainfall intensity, evapotranspiration, infiltration, and
206 storage dynamics. When considering the input features of rainfall, evapotranspiration, and initial
207 soil moisture, it is important to maintain a monotonic relationship between each feature and the
208 resulting runoff. However, the traditional DL models disregard the physical relationships between
209 inputs and outputs. This lack of consistency with the physical principles of water balance equations
210 undermines the overall reliability of the model. Therefore, this study incorporates inequality
211 constraints to enforce the desired monotonic relationships between rainfall, evapotranspiration,
212 initial soil moisture, and runoff. Under the assumption that all other input variables remain
213 unchanged, a new time series of rainfall, evapotranspiration, and initial soil moisture is generated
214 respectively by applying a small random increase using the random.uniform function. These new
215 time series are then combined with the unchanged time series to form new input data. The difference
216 between the predicted values corresponding to the new data and the predicted values corresponding
217 to the original input data is calculated. This difference is then converted into a specific loss value
218 using the ReLU function and added to the loss function.

219



220

221 **Fig. 1.** (a) The PHY-FTMA-LSTM model architecture. (b) Feature-time-based attention matrix



222 generation process for each attention head. (c) Feature-time-based multi-head attention workflow.

223 (d) The internals of LSTM cells.

224 For rainfall, the runoff should increase if there is a slight increase in rainfall at the current time
225 step, provided that other variables are constant, and the monotonic relationship and losses for
226 rainfall-runoff are expressed as follows:

$$227 \quad f[p(t) + \Delta p, t] - f[p(t), t] \geq 0 \quad (13)$$

$$228 \quad Loss_p = \frac{1}{N_p} \sum_{i=1}^{N_p} \{\text{ReLU}\{f[p(t), t] - f[p(t) + \Delta p, t]\} \geq 0\}^2 \quad (14)$$

229 where Δp is the small increase in rainfall, $Loss_p$ is the error in the monotonic relationship of rainfall
230 runoff, N_p is the sample length of the perturbed rainfall, and ReLU is the response function.

231 For evapotranspiration, the runoff should decrease if there is a slight increase in
232 evapotranspiration at the current time step, provided that other variables are constant, and the
233 monotonic relationship and losses for evapotranspiration runoff are expressed as follows:

$$234 \quad f[e(t) + \Delta e, t] - f[e(t), t] \leq 0 \quad (15)$$

$$235 \quad Loss_e = \frac{1}{N_e} \sum_{i=1}^{N_e} \{\text{ReLU}\{f[e(t), t] - f[e(t) + \Delta e, t]\} \leq 0\}^2 \quad (16)$$

236 where Δe is the small increase in evapotranspiration, $Loss_e$ is the error in the monotonic relationship
237 of evapotranspiration runoff, N_e is the sample length of the perturbed evapotranspiration.

238 For soil moisture, the runoff should increase if the initial soil moisture of the watershed
239 increases slightly for each flood event, provided that other variables are constant, and the monotonic
240 relationship and losses between initial soil moisture and runoff are expressed as follows:

$$241 \quad f[s(t) + \Delta s, t] - f[s(t), t] \geq 0 \quad (17)$$

$$242 \quad Loss_s = \frac{1}{N_s} \sum_{i=1}^{N_s} \{\text{ReLU}\{f[s(t), t] - f[s(t) + \Delta s, t]\} \geq 0\}^2 \quad (18)$$

243 where Δs is the small increase in initial soil moisture, $Loss_s$ is the error in the monotonic relationship
244 of initial soil moisture runoff, N_s is the sample length of the perturbed initial soil moisture.

245 Based on the above physical constraints of flood forecasting, the loss function of the traditional
246 LSTM model is improved with the following equation:

$$247 \quad Loss = \lambda_{data} Loss_{data} + \lambda_p Loss_p + \lambda_e Loss_e + \lambda_s Loss_s \quad (19)$$



248 where $Loss$ is the loss function of the LSTM guided by the physical constraints of flood forecasting;
 249 $Loss_{data}$ is the mean square error of the observed and predicted values of the LSTM; λ_{data} , λ_p , λ_e ,
 250 λ_s are the weighting coefficients of different losses, respectively. To treat the three physical
 251 constraints equally, the weighting coefficients of the four losses are set to $\{0.7, 0.1, 0.1, 0.1\}$.

252 2.4. Evaluation metrics

253 To evaluate the accuracy of different models for flood forecasting, the Nash-Sutcliffe efficiency
 254 (NSE), Kling–Gupta efficiency (KGE), the coefficient of determination (R^2), root mean square error
 255 (RMSE), and mean absolute error (MAE) are selected for evaluation. The specific equations are as
 256 follows:

$$257 \quad NSE = 1 - \frac{\sum_{i=1}^n (Q_i - Q'_i)^2}{\sum_{i=1}^n (Q_i - \bar{Q}_i)^2} \quad (20)$$

$$258 \quad KGE = 1 - \sqrt{(R-1)^2 + (\alpha-1)^2 + (\beta-1)^2} \quad (21)$$

$$259 \quad R^2 = \frac{\left(\sum_{i=1}^n (Q_i - \bar{Q}_i)(Q'_i - \bar{Q}'_i) \right)^2}{\sum_{i=1}^n (Q_i - \bar{Q}_i)^2 \sum_{i=1}^n (Q'_i - \bar{Q}'_i)^2} \quad (22)$$

$$260 \quad RMSE = \sqrt{\frac{\sum_{i=1}^n (Q_i - Q'_i)^2}{n}} \quad (23)$$

$$261 \quad MAE = \frac{1}{n} \sum_{i=1}^n |Q_i - Q'_i| \quad (24)$$

262 where Q_i is the observed value; Q'_i is the predicted value; \bar{Q}_i is the observed mean value; \bar{Q}'_i is
 263 the mean value of the predicted series; α between the standard deviation of the predicted value and
 264 that of the observed value; β is the ratio between the mean of the predicted value and that of the
 265 observed value; n is the total number of samples. The NSE is commonly used to evaluate
 266 hydrological prediction models, KGE considers the contribution of mean, variance and correlation
 267 on model performance, R^2 is often used to evaluate the linear correlation between the forecast
 268 process and the observed process. The values of NSE, KGE and R^2 range from 0 to 1. The closer the
 269 result is to 1, the more accurate the forecast result is and the higher the model credibility is. RMSE
 270 and MAE are used to reflect the degree of deviation between the predicted and observed values, the

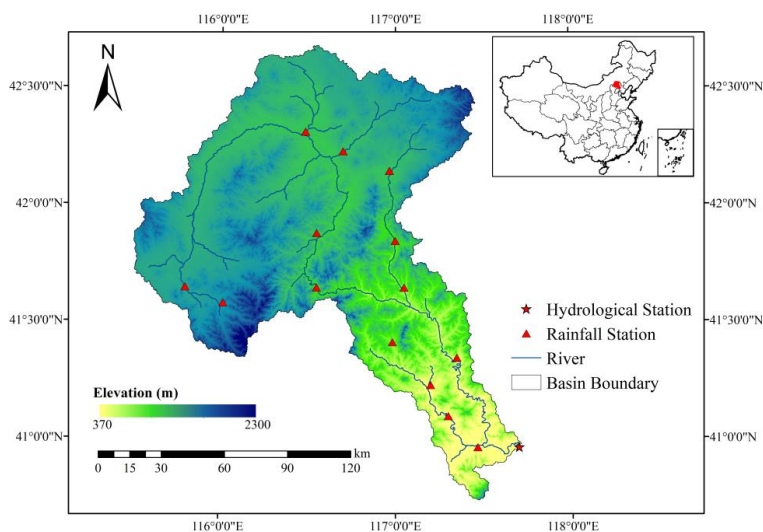


271 smaller the value the smaller the deviation.

272 **3. Study area and data**

273 3.1. Study area

274 In this study, the watershed controlled by the Sandaohezi station in the Luan River Basin was
275 selected as the study area. The Luan River originates from the northern foot of Bayangurtu Mountain
276 in Hebei Province, with a total length of 888 km, and flows through Inner Mongolia, Hebei, and
277 Liaoning provinces before injecting into the Bohai Sea at Laoting County, Hebei Province. The
278 station is in the middle reaches of the mainstream of the Luan River, controlling a watershed area
279 of 17100 km², accounting for about 40% of the total area of the Luan River basin. Geographically,
280 it is located between 115.5°E to 117.7°E longitude and 40.7°N to 42.7°N latitude. The elevation of
281 the study area ranges from 370 to 2300 m, with a high northwest to low southeast topography. Except
282 for the upstream origin of the dam plateau, the rest of the area is dominated by mountainous terrain.
283 The northwest of the basin is located in the temperate continental climate zone, precipitation is
284 scarce and concentrated in summer; the southeast is located in the temperate monsoon climate zone,
285 with cold, dry winters and hot, rainy summers. The average annual temperature of the basin ranges
286 from 5 to 12°C, and the average annual runoff is about 480 million m³. The average annual rainfall
287 is about 500mm, and the spatial and temporal distribution of rainfall within the year is uneven,
288 mainly concentrated from May to September, and the precipitation decreases from south to north.
289 Floods in the basin are mostly formed by heavy rainfall, which is short-lived and strong, making the
290 flooding process steep up and steep down, often causing disasters in the downstream areas.
291 Consequently, accurate flood forecasting is of utmost importance for effective flood control and
292 water resources management in the Luann River basin. The location of the study area and the
293 stations are shown in Fig. 2.



294
295 **Fig.2.** Geographical location of the study area and hydrological and rainfall stations.

296 3.2. Data

297 The rainfall and runoff data were obtained from the Hydrological Yearbook of the Haihe River
298 Basin, including rainfall data from 15 rainfall stations, such as Sandaohezi, Zhangbaiwan, and
299 Baorono, and runoff data from Sandaohezi hydrological station. The period covers 39 years from
300 1964 to 1989, 1991, and 2006 to 2017. There is a gap in the data for 1990 and 1992 to 2005 due to
301 incomplete data collection.

302 The evapotranspiration and soil moisture data were obtained from the Global Land Surface
303 Data Assimilation System (GLDAS) using the GLDAS-Noah model product 0.25°×0.25° spatial
304 resolution, 3h temporal resolution dataset, and the evapotranspiration data were averaged backward
305 3h, and the soil moisture data were instantaneous values. Among them, GLDAS-2.0 provides data
306 from 1964 to 2014, and GLDAS-2.1 provides data from 2015.

307 In this study, 30 flood events during the 39 years were selected (Table 1), and the collected
308 observed runoff data were linearly interpolated to 1h step data, the observed rainfall data were
309 averaged to 1h step data, and the Tyson polygon method was used to derive the areal rainfall. For
310 evapotranspiration and soil moisture, the average values were calculated for each grid in the
311 watershed at each period, where the soil moisture was taken as the initial soil moisture before the
312 onset of rainfall for each flood event. Twenty flood events were used for model training, ten flood



313 events were used for model validation.

314 Since different input features have different magnitudes, maximum-minimum normalization

315 was used to process the input data into the range [0,1], see Eq. (25).

316
$$x_{norm} = \frac{x_i - x_{min}}{x_{max} - x_{min}} \quad (25)$$

317 where x_{norm} is the normalized data, x_i is the original data, and x_{min} and x_{max} are respectively the
 318 minimum and maximum values of the original data.

319 **Table 1** Flood events used in the study.

Dataset	Flood number	Peak discharge (m ³ /s)	Year	Duration (month/day/hour)
Training	1	314.2	1964	08/01/04-08/09/12
	2	218	1964	08/13/02-08/16/00
	3	313	1965	07/17/20-07/21/12
	4	204	1966	07/27/16-07/31/20
	5	260	1968	07/27/12-07/30/22
	6	154	1969	08/20/12-08/27/12
	7	296	1971	07/17/15-07/29/08
	8	153	1972	07/19/08-07/24/08
	9	742	1973	08/12/04-08/26/08
	10	213	1975	08/11/00-08/16/08
	11	218	1978	08/25/12-09/03/08
	12	246	1982	07/22/12-07/29/16
	13	313	1983	08/04/00-08/11/20
	14	400	1985	08/24/05-08/31/04
	15	210	1986	08/08/04-08/13/08
	16	87.5	1987	08/19/12-08/23/04
	17	465	1991	06/10/04-06/18/00
	18	70.1	2008	08/10/00-08/16/00
	19	149	2010	07/30/17-08/04/20
	20	80.4	2015	07/27/16-07/31/16
Validation	21	241	1965	08/26/21-08/30/20
	22	260	1967	06/27/12-06/29/22
	23	164	1970	07/14/12-07/16/04
	24	506.7	1974	07/23/12-08/06/08
	25	313	1979	08/13/04-08/21/08
	26	132	1985	08/11/16-08/14/04
	27	212	1989	06/03/22-06/07/04
	28	205	2011	08/14/10-08/20/04



29	95.9	2013	07/21/08-07/25/16
30	84.2	2013	08/13/09-08/21/00

320 3.3. Model construction

321 This study is based on Python 3.9, and the Numpy, Pandas, and Scikit-Learn packages in
322 Python are used for data processing, and the LSTM, FTA-LSTM, FTMA-LSTM, and PHY-FTMA-
323 LSTM models are constructed using the Keras library in TensorFlow.

324 The model inputs are runoff, rainfall, evapotranspiration, and initial soil moisture for a
325 specified time step, and the outputs are the discharge from 1 to 6h of the lead time. All four models
326 use the ReLU activation function, which avoids gradient vanishing and is more effective compared
327 to the tanh and sigmoid functions. The Adam optimizer is used and the LSTM layer is a single layer,
328 with the number of attention heads set to 3 for the FTMA-LSTM and PHY-FTMA-LSTM. The mean
329 square error is the loss function of the four models, and for PHY-FTMA-LSTM it incorporates
330 physical constraints, as shown in Eq. (19). To avoid overfitting, all models use the early stopping
331 and set the maximum number of epochs to 200.

332 To construct the base models, the common values of the DL model parameters are used as the
333 initial values. The base models have an observed input time step of 12 hours, a learning rate of 0.001,
334 batch size of 64, and hidden units set to 128. After evaluating the performance of the base models,
335 parameter optimization is performed separately for each of the four models, considering that the
336 optimal parameter combinations may differ among the models. The goal is to study the effects of
337 the input time step and three hyperparameters (learning rate, batch size, and hidden units) on the
338 model performance. The ranges used for parameter optimization are as follows: input time step of
339 3 to 24 hours, learning rate of 0.00001 to 0.01, batch size of 16 to 256, and hidden units of 32 to
340 512. A single parameter is varied while the other parameters are taken as their initial values.
341 Considering the stochastic nature of the DL model running process, each of the four models is
342 repeated five times for each lead time, and the results with the best prediction performance are
343 selected for analysis.

344 4. Results

345 4.1. Model optimization

346 The LSTM, FTA-LSTM, FTMA-LSTM, and PHY-FTMA-LSTM base models are established



347 individually, and their average NSE values during the 1-6 hour lead time, measure to evaluate flood
348 prediction accuracy, are found to be 0.925, 0.930, 0.936, and 0.950, respectively. These results
349 indicate that all four base models can effectively predict flooding events. In order to determine the
350 optimal parameter combination for each model and how individual parameter variations affect the
351 model performance, the following parameters are investigated while keeping the other three
352 parameters constant: input time step, learning rate, batch size, and hidden units.

353 Regarding the input time step of observations, experiments are conducted by varying the time
354 step within a certain range. The result depicted in Figure 3(a) shows that the average NSE value for
355 all four models is highest at a time step of 12 hours and decreases with increasing time step. The
356 worst performance is observed at a time step of 24 hours. This observation suggests that longer input
357 sequences introduce more noise, and the inclusion of extraneous information adversely affects the
358 final prediction. Therefore, a 12-hour input time step is identified as the optimal choice for flood
359 forecasting in all four models and is adopted for subsequent experiments.

360 For the learning rate, tests are performed using a learning rate ranging from 0.00001 to 0.01.
361 The finding, presented in Figure 3(b), indicates that the performance of the four models is
362 comparable at learning rates of 0.01 and 0.001. However, when the learning rate is set to 0.0001 and
363 0.00001, the models exhibit slow convergence and degrade performance rapidly. Considering the
364 possibility of failure to converge at a very high learning rate, a combined analysis suggests a learning
365 rate of 0.001 as the optimal choice for all four models in the subsequent studies.

366 The batch size optimization ranges from 16 to 256. The result depicted in Figure 3(c)
367 demonstrates varying performances of the four models with different batch sizes. The LSTM model
368 achieves the highest average NSE of 0.932 at a batch size of 128. Similarly, the FTA-LSTM model
369 attained its highest average NSE of 0.932 at a batch size of 32. On the other hand, the FTMA-LSTM
370 and PHY-FTMA-LSTM models reach their highest average NSE values at a batch size of 64, with
371 0.936 and 0.950, respectively. Consequently, the optimal batch size for flood forecasting is
372 determined as 128, 32, 64, and 64 for the LSTM, FTA-LSTM, FTMA-LSTM, and PHY-FTMA-
373 LSTM models, respectively. These batch sizes are employed for subsequent studies.

374 Regarding the hidden units, tests are conducted with the count varying from 32 to 512. Figure
375 3(d) illustrates the distinct performances of the four models concerning different hidden units. The



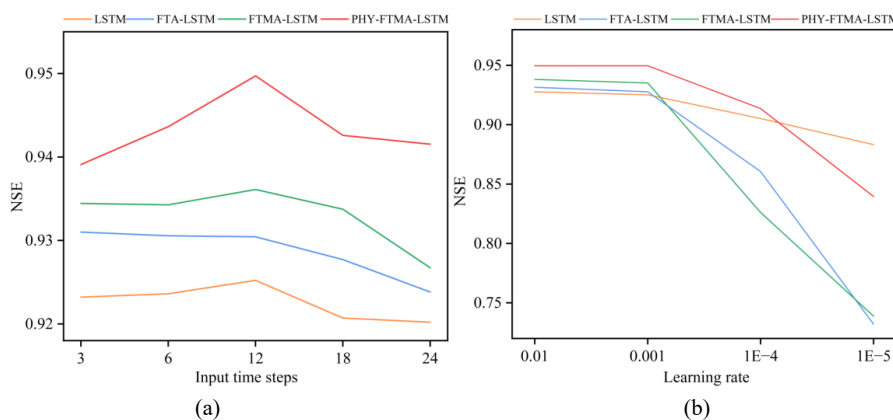
376 LSTM model achieves the highest average NSE of 0.925 with 64 hidden units. The FTA-LSTM and
 377 FTMA-LSTM models attain their highest average NSE values of 0.935 and 0.939 with 256 hidden
 378 units, respectively. In contrast, the PHY-FTMA-LSTM model reaches the highest average NSE of
 379 0.950 at 128. Accordingly, the optimal hidden units for flood prediction are identified as 64, 256,
 380 256, and 128 for the LSTM, FTA-LSTM, FTMA-LSTM, and PHY-FTMA-LSTM models,
 381 respectively.

382 Considering the above parameter optimization process, the model parameters used in the
 383 subsequent study are as follows (Table 2). Notably, the PHY-FTMA-LSTM model consistently
 384 outperforms the other three models across various parameter values, exhibiting the smallest
 385 variation in NSE. These findings indicate that the PHY-FTMA-LSTM model proposed in this paper
 386 offers the best and most stable performance.

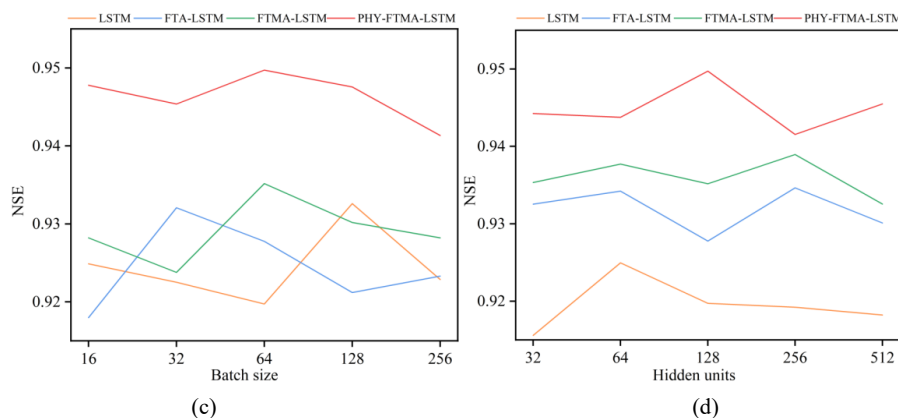
387 **Table 2** Parameters of models.

Models	Input time step	Learning rate	Batch size	Hidden units
LSTM	12	0.001	128	64
FTA-LSTM	12	0.001	32	256
FTMA-LSTM	12	0.001	64	256
PHY-FTMA-LSTM	12	0.001	64	128

388



389
 390



391
392

393 **Fig.3.** The NSE values for 6 lead times with different (a) input time steps of observations, (b)
394 learning rate, (c) batch size, and (d) hidden units.

395 4.2. Model performance evaluation

396 The LSTM, FTA-LSTM, FTMA-LSTM, and PHY-FTMA-LSTM models are constructed using
397 the optimal parameters mentioned above, the evaluation metrics of the forecasting performance of
398 the four models in the training and validation stages are shown in Table 3 and Table 4. All the metrics
399 of the four models almost outperform the validation period in the training period. And with the
400 increase of the lead time, the gap between the performance of the models in the training period and
401 the testing period gradually increases. It can be seen that the three models based on the attention
402 mechanism outperform the original LSTM model in all lead times. It indicates that the dual attention
403 module of time and feature proposed in this paper effectively focuses on the more significant
404 historical moments and feature variables, improving the performance of the LSTM model. Among
405 the attention-based models, the FTMA-LSTM model, which utilizes a multi-headed attention
406 mechanism, achieves better performance than the FTA-LSTM model with a single attention head in
407 most cases. This demonstrates that the parallel computation of the multi-head attention mechanism
408 enables the model to emphasize more important information in the input compared to the single-
409 head attention mechanism. Furthermore, the PHY-FTMA-LSTM model, which incorporates
410 physical constraints, outperforms the other three models across almost all metrics. Specifically, at
411 the lead time $t+1$, compared to the original LSTM model, the PHY-FTMA-LSTM model shows an
412 improvement in NSE, KGE, and R^2 , increasing from 0.977 to 0.988, from 0.953 to 0.984 and from
413 0.979 to 0.988, respectively. Additionally, the RMSE and MAE decrease by 27.4% and 49.6%,



414 respectively. At the lead time $t+6$, NSE increases from 0.865 to 0.908, KGE from 0.851 to 0.905,
 415 R^2 from 0.886 to 0.911, and RMSE and MAE decrease by 21.1% and 15.1%, respectively. These
 416 results mean that incorporating physical constraints enables the DL model to understand the
 417 monotonic relationship presented in the flooding process, improving forecast accuracy by enhancing
 418 the model's physical consistency.

419 As the lead time increases, the performance of all four models declines, suggesting that their
 420 robustness and generalization gradually deteriorate. However, the extent of the decline in the four
 421 model metrics varies. In terms of NSE, when transitioning from a 1-hour to a 6-hour lead time, the
 422 PHY-FTMA-LSTM model exhibits the smallest decline of 0.065 during the training period, while
 423 the LSTM, FTA-LSTM, and FTMA-LSTM models experience decreases of 0.072, 0.079, and 0.073
 424 respectively. During the validation period, the NSE value decreases by 0.080 for the PHY-FTMA-
 425 LSTM model and by 0.112, 0.109, and 0.104 for the LSTM, FTA-LSTM and FTMA-LSTM models,
 426 respectively. Maintaining high accuracy in longer lead times is crucial in practical applications.
 427 Extended lead times necessitate more comprehensive information for accurate predictions,
 428 presenting challenges for the models. Nonetheless, the PHY-FTMA-LSTM model exhibits minimal
 429 degradation, indicating its superior ability to adapt to longer lead times and maintain high precision.
 430 This superiority may be attributed to the unique characteristics and structure of the PHY-FTMA-
 431 LSTM model. It likely encompasses considerations of physical factors and key input features,
 432 enabling a better capture of flood complexity and variability. This advantage positions the model
 433 favorably in scenarios requiring predictions further into the future.

434

435 **Table 3** Performance of the four models for flood forecasting at different lead times for training.

Lead times/h	Models	NSE	KGE	R^2	RMSE	MAE
t+1	LSTM	0.977	0.964	0.980	16.14	7.14
	FTA-LSTM	0.986	0.972	0.987	12.32	5.19
	FTMA-LSTM	0.990	0.977	0.990	10.62	4.76
	PHY-FTMA-LSTM	0.992	0.984	0.992	9.65	4.03
t+2	LSTM	0.959	0.944	0.963	21.52	11.29
	FTA-LSTM	0.966	0.983	0.967	20.93	7.85
	FTMA-LSTM	0.969	0.960	0.972	18.54	8.80
	PHY-FTMA-LSTM	0.976	0.949	0.977	16.56	9.10



Lead times/h	Models	NSE	KEGE	R ²	RMSE	MAE
t+3	LSTM	0.943	0.945	0.948	25.09	13.91
	FTA-LSTM	0.949	0.943	0.952	22.05	11.02
	FTMA-LSTM	0.954	0.963	0.955	21.14	10.79
	PHY-FTMA-LSTM	0.958	0.955	0.963	20.01	11.45
t+4	LSTM	0.933	0.915	0.942	27.59	15.83
	FTA-LSTM	0.945	0.956	0.948	23.06	14.57
	FTMA-LSTM	0.948	0.953	0.949	22.12	13.75
	PHY-FTMA-LSTM	0.950	0.948	0.955	23.63	14.27
t+5	LSTM	0.929	0.917	0.929	29.16	18.91
	FTA-LSTM	0.930	0.942	0.931	27.99	16.37
	FTMA-LSTM	0.934	0.925	0.937	26.08	16.18
	PHY-FTMA-LSTM	0.937	0.931	0.937	25.58	15.19
t+6	LSTM	0.905	0.900	0.917	33.29	19.78
	FTA-LSTM	0.907	0.913	0.913	33.63	17.86
	FTMA-LSTM	0.917	0.926	0.919	30.59	15.83
	PHY-FTMA-LSTM	0.927	0.949	0.929	28.05	16.04

436 Figure 4 displays the scatter plots for the LSTM, FTA-LSTM, FTMA-LSTM, and PHY-
 437 FTMA-LSTM models during the training and validation periods. When the foresight period is 1
 438 hour, all models demonstrate predictions that closely track the ideal 1:1 line. The PHY-FTMA-
 439 LSTM model outperforms the others, exhibiting the narrowest scatter distribution. However, as the
 440 lead time increases, the scatter plots of the four models show varying degrees of deterioration,
 441 becoming more uneven and scattered. The high discharge prediction error increases in the training
 442 period, and the validation period reveals numerous underestimated discharges. Among them, the
 443 PHY-FTMA-LSTM model performs the best (with the narrowest scatter distribution), followed by
 444 the FTA-LSTM and FTMA-LSTM models. The LSTM model performs the worst. Notably, during
 445 the validation period, for longer foresight periods, the high flow scatter of all models deviates further
 446 from the ideal 1:1 line. One possible explanation is the scarcity of high flow instances in the training
 447 data. As the lead time increases, the models struggle to capture the necessary information, leading
 448 to underestimation and poorer predictions. For a foresight period of 6 hours, the scatter plots of the
 449 LSTM, FTA-LSTM, and FTMA-LSTM models both in the training and validation periods exhibit
 450 discrete distributions. In contrast, the PHY-FTMA-LSTM model's scatter plot shows the narrowest
 451 band and is closest to the ideal 1:1 line. Consequently, the PHY-FTMA-LSTM model achieves the

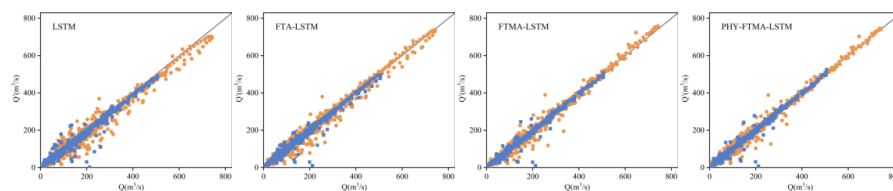


452 highest prediction accuracy, effectively reducing prediction errors for longer lead times. The FTA-
 453 LSTM and FTMA-LSTM models follow while the LSTM model performs the worst in terms of
 454 prediction accuracy.

455 **Table 4** Performance of the four models for flood forecasting at different lead times for validation.

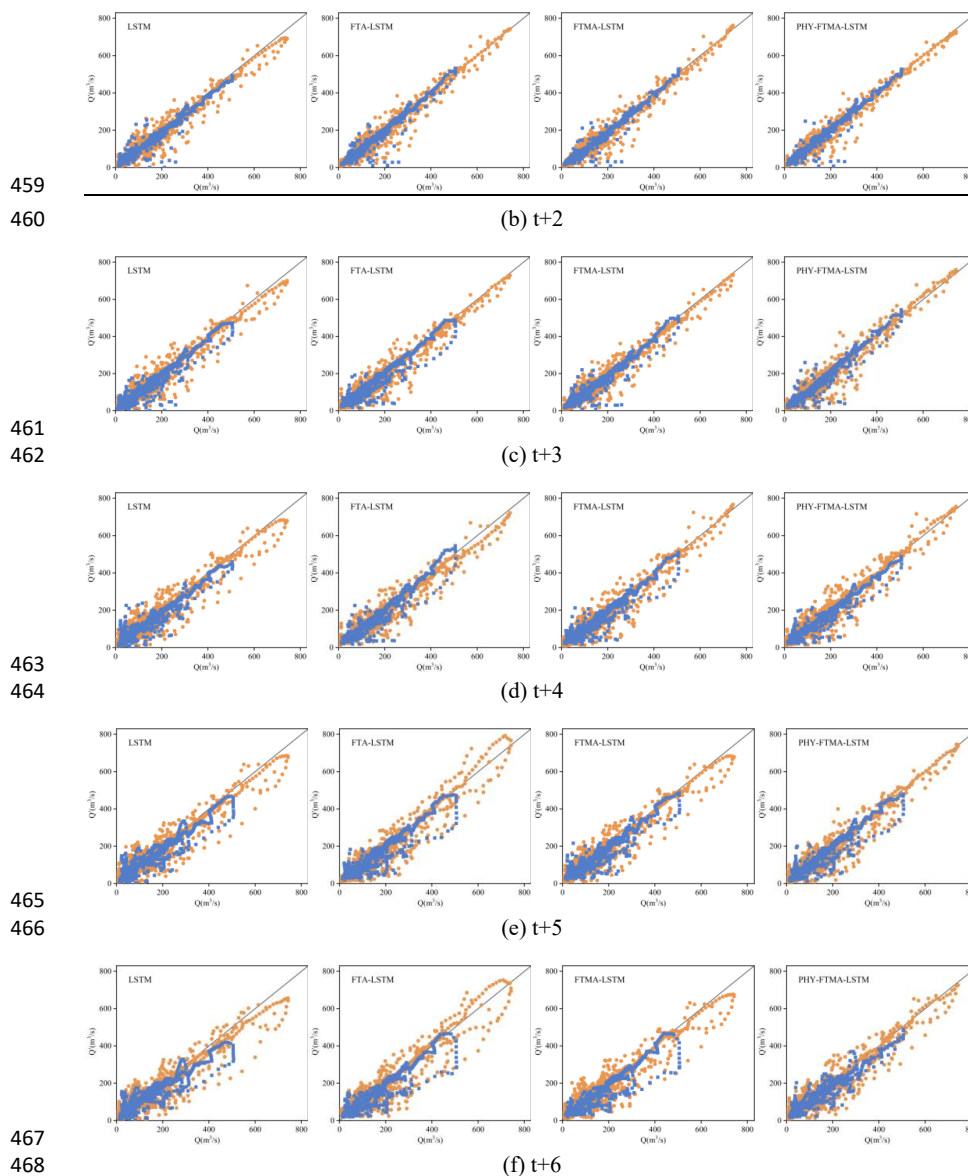
Lead times/h	Models	NSE	KEGE	R ²	RMSE	MAE
t+1	LSTM	0.977	0.953	0.979	15.84	8.45
	FTA-LSTM	0.985	0.969	0.985	12.65	6.28
	FTMA-LSTM	0.987	0.975	0.988	11.83	5.04
	PHY-FTMA-LSTM	0.988	0.984	0.988	11.50	4.26
t+2	LSTM	0.956	0.939	0.961	21.83	11.94
	FTA-LSTM	0.961	0.974	0.961	19.07	10.22
	FTMA-LSTM	0.967	0.950	0.970	18.83	9.52
	PHY-FTMA-LSTM	0.968	0.954	0.970	18.56	9.45
t+3	LSTM	0.934	0.928	0.938	27.09	14.93
	FTA-LSTM	0.942	0.927	0.943	25.07	13.49
	FTMA-LSTM	0.948	0.947	0.951	23.66	12.56
	PHY-FTMA-LSTM	0.952	0.945	0.955	21.57	12.74
t+4	LSTM	0.918	0.914	0.929	28.15	16.43
	FTA-LSTM	0.928	0.938	0.933	28.17	14.20
	FTMA-LSTM	0.931	0.946	0.933	28.44	16.24
	PHY-FTMA-LSTM	0.939	0.938	0.944	26.13	14.59
t+5	LSTM	0.898	0.890	0.900	36.43	22.83
	FTA-LSTM	0.905	0.911	0.910	32.54	19.36
	FTMA-LSTM	0.915	0.915	0.920	30.43	20.52
	PHY-FTMA-LSTM	0.918	0.930	0.919	30.33	16.65
t+6	LSTM	0.865	0.851	0.886	40.61	23.77
	FTA-LSTM	0.876	0.894	0.886	37.38	20.57
	FTMA-LSTM	0.883	0.889	0.896	36.52	20.65
	PHY-FTMA-LSTM	0.908	0.905	0.911	32.02	20.18

456



457
 458

(a) t+1



459
460

461
462

463
464

465
466

467
468

469 **Fig.4.** Scatter plots of observed and predicted discharges in the training and validation stages, in
470 which yellow represents the training stage and blue represents the validation stage.

471 4.3. Typical flood event forecast results

472 Floods in the basin are mainly two types, single-peak and double-peak, so two typical flood
473 events were selected to analyze the specific flood process: a double-peak flood event (19740723)
474 with a peak discharge of 507 m³/s and 290 m³/s, and a single-peak flood event (19790813) with a



475 peak discharge of 313 m³/s. Fig. 5 and Fig. 6 illustrate the flood processes of the two events predicted
476 by the four models. It can be observed that as the lead time increases, the prediction hydrographs
477 from all four models gradually deviate from the observed values and the three evaluation metrics
478 decrease. Notably, the LSTM model exhibits the greatest decline in prediction performance,
479 followed by the FTA-LSTM and FTMA-LSTM models. In contrast, the PHY-FTMA-LSTM model
480 demonstrates relatively better performance across the evaluated flood events.

481 Based on the analysis of prediction hydrographs, the four models exhibit better performance in
482 predicting the double-peak flood event compared to the single-peak flood event. Additionally, the
483 models demonstrate higher accuracy in predicting the rising stage of floods in contrast to the falling
484 stage. Specifically, the prediction errors increase as the duration of the flood increases, and there is
485 a time lag in predicting the occurrence of the second flood peak. When it comes to the single-peak
486 flood event, the predictions by the four models display greater fluctuations, and the time lag problem
487 is more pronounced, along with an overestimation of the peak discharge.

488 Regarding the 19740723 flood event, the LSTM model generally underestimates the discharge
489 values, and the discrepancy with the observed hydrograph gradually increases as the lead time
490 increases. Although the FTA-LSTM and FTMA-LSTM models also underestimate the discharge,
491 their errors are reduced, indicating improved performance compared to the LSTM model. In contrast,
492 the PHY-FTMA-LSTM model predicts the flood hydrograph more accurately. However, when the
493 foresight period is 6 h, the PHY-FTMA-LSTM model experiences significant prediction errors due
494 to anomalous fluctuations.

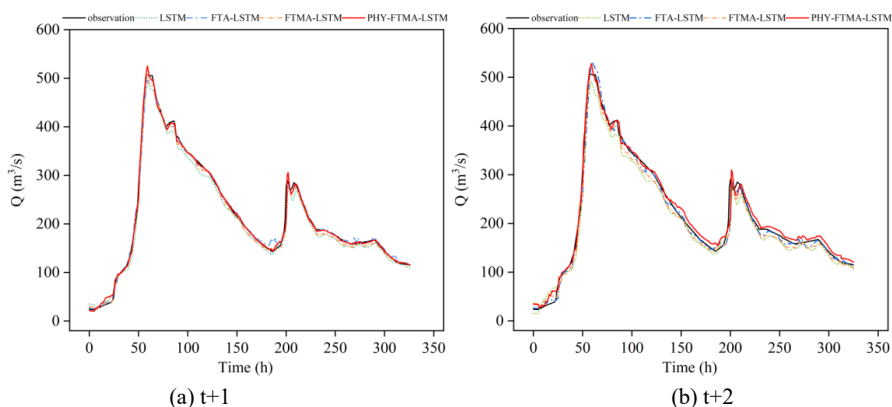
495 For the 19790813 flood event, the LSTM model demonstrates a noticeable deviation from the
496 predicted hydrograph with increasing lead times. The FTA-LSTM and FTMA-LSTM models
497 exhibit better performance, as their predicted hydrographs are closer to the observed ones. However,
498 there is some overestimation of the peak discharge in these models. Additionally, all three models
499 suffer from a more severe time lag issue in longer foresight periods. In contrast, the PHY-FTMA-
500 LSTM model shows smaller volume errors and is closer to the observed hydrograph. Nevertheless,
501 this model exhibits a more pronounced overestimation of the peak discharge.

502 In conclusion, the LSTM model exhibits poor prediction performance for longer lead times.
503 On the other hand, the FTA-LSTM, FTMA-LSTM, and PHY-FTMA-LSTM models show improved

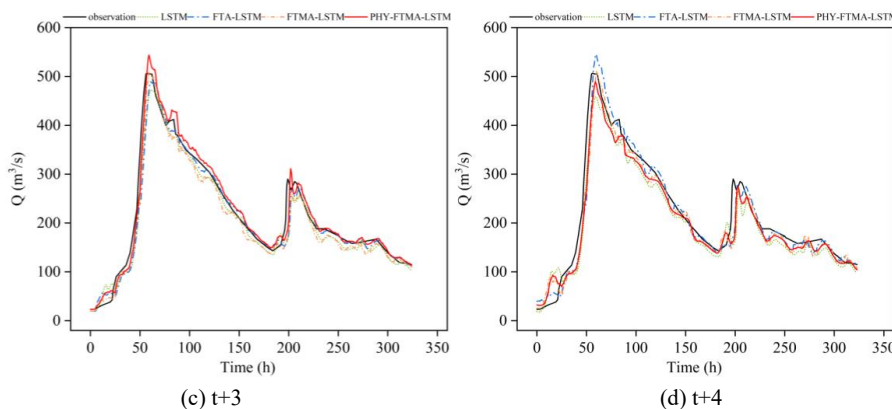


504 performance with longer lead times and higher forecasting accuracy. Among these models, the PHY-
505 FTMA-LSTM model stands out by producing better predictions for both single-peak and multi-peak
506 flood events, but it may encounter challenges with predicting anomalous fluctuations at longer lead
507 times. Additionally, the PHY-FTMA-LSTM model mitigates the issue of time lag to some extent by
508 considering the physical monotonicity relationship.

509
510



511
512



513



514

(e) t+5

(f) t+6

515

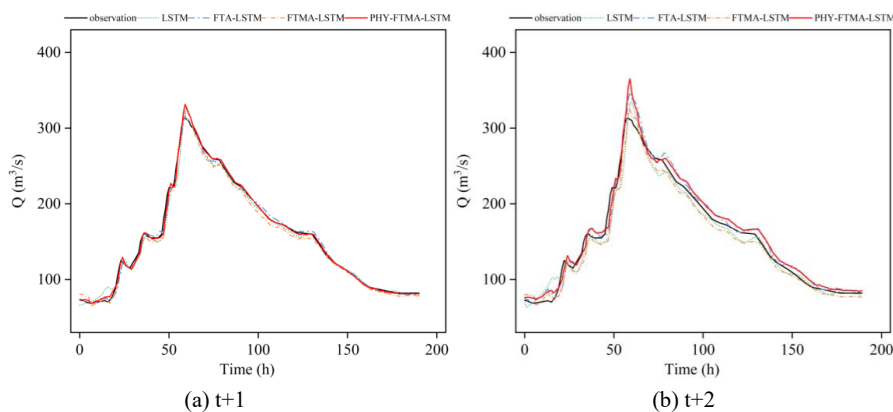
Fig.5. Comparison of observed and predicted values of the 19740723 flood event by the four

516

models.

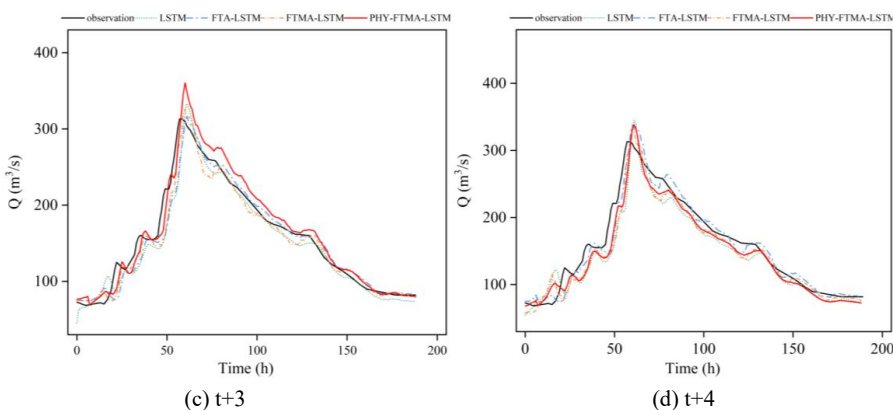
517

518



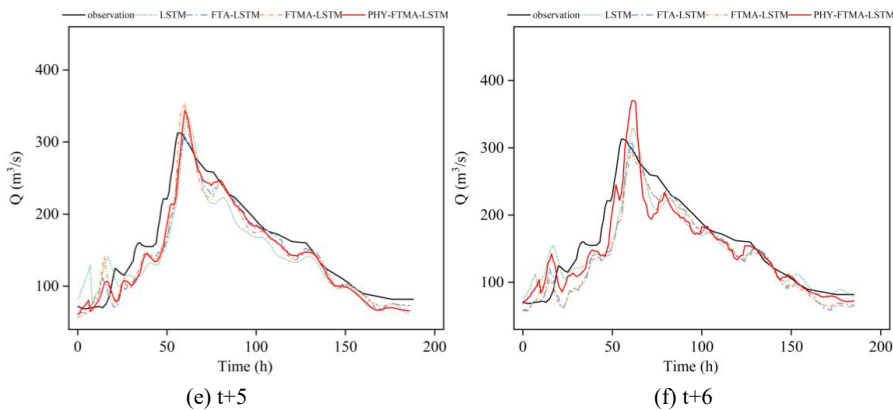
519

520



521

522



523

Fig.6. Comparison of observed and predicted values of the 19790813 flood event by the four



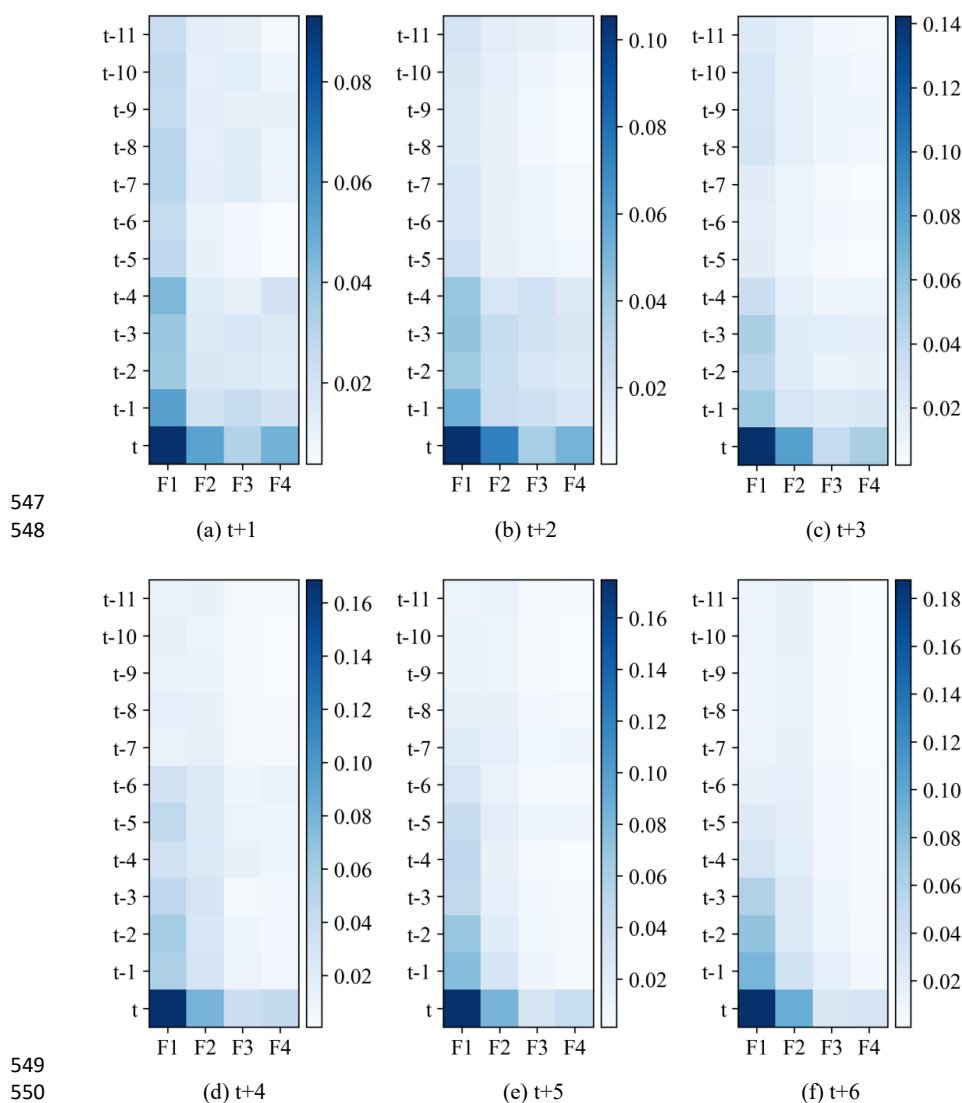
524 models.

525 4.4. Visual attention analysis

526 To investigate the changes in features and time attention of PHY-FTMA-LSTM with different
527 lead times, the attention weights of PHY-FTMA-LSTM are visualized in Fig. 7. The figure consists
528 of six subplots representing lead times ranging from $t+1$ to $t+6$.

529 From Fig. 7, it can be observed that the distribution pattern of the weights remains relatively
530 similar across different forecasting periods. The temporal attention weights decrease as the historical
531 moment increases. Among the feature-based weights, runoff has the highest proportion, followed
532 by rainfall, and finally the initial soil moisture and evapotranspiration. These results align with
533 hydrological principles, where runoff is considered the most direct manifestation of the flooding
534 process and holds the highest importance. Rainfall, as the main driver of flood formation,
535 significantly influences flooding. In contrast, the effects of initial soil moisture and
536 evapotranspiration in the basin are more indirect and therefore receive lower weights. In the case of
537 the Luan River basin, which is relatively arid, the initial soil moisture of the basin is typically not
538 saturated. During a rainfall-induced flood, there is a possibility of transitioning from infiltration-
539 excess runoff to saturation-excess runoff. Hence, special attention should be given to the role of the
540 initial soil moisture, which carries slightly greater relative importance than evapotranspiration.

541 As the forecasting horizon extends, the feature-time-based weights of the model become more
542 concentrated, with the time-based weights gradually moving forward. Consequently, the model
543 places more emphasis on the values that are closer to the current moment. Additionally, the feature-
544 based attention module exhibits a gradual increase in attention to rainfall while decreasing attention
545 to evapotranspiration and the initial soil moisture. Notably, runoff retains its status as the most
546 influential factor.



547
548

549
550

551 **Fig.7.** The visualization of feature-time-based attention weights of the PHY-FTMA-LSTM. The
552 X-coordinate variables F1 to F4 represent the input features of runoff, rainfall, evapotranspiration,
553 and initial soil moisture of the watershed, respectively. The Y-coordinate variables represent the
554 input history moments.

555 5. Discussion

556 The input time step of observations, learning rate, batch size, and hidden units are significant
557 parameters that influence the performance of the model, and the optimal parameters may vary for



558 different structural models (Xiang et al., 2020; Cao et al., 2022). In this study, four models, namely
559 LSTM, FTA-LSTM, FTMA-LSTM, and PHY-FTMA-LSTM, have been constructed. To ensure that
560 each model achieves its optimal prediction performance and to investigate the impact of different
561 parameter variations on model performance, the same parameter values are utilized to build the four
562 base models individually. After confirming that the base models meet the accuracy requirements for
563 flood forecasting, the optimal parameter combination for each model is determined. This is done by
564 selecting the parameter value associated with the highest NSE obtained through single parameter
565 tuning. The single parameters are changed while keeping the initial values of the other three
566 parameters constant. This approach ensures that the subsequent analysis reflects the best
567 performance achievable by each model's specific structure. Moreover, it enables a more explicit
568 evaluation of the performance changes resulting from the addition of attention mechanisms and
569 physical constraints to the model.

570 In terms of model performance evaluation metrics, the PHY-FTMA-LSTM model
571 demonstrates the best overall performance. However, a closer examination reveals that its KGE
572 score may not necessarily be optimal. This could be attributed to the comprehensiveness of the KGE
573 metric, which considers factors such as correlation, mean consistency, and variance consistency of
574 the flow. Fluctuations in the KGE score may arise from various uncertainties related to data quality,
575 model structure, and flood forecasting.

576 With an increase in the forecast period, the performance of the model, particularly the LSTM
577 model, shows a significant decrease, consistent with the findings reported by Xu et al. (2021). They
578 provided NSE, RMSE, and Bias indices for the LSTM model in forecast periods of 1~12 hours,
579 demonstrating that the LSTM model meets prediction requirements for short forecast periods.
580 However, as the forecast period extends, the accuracy diminishes, leading to underestimation of
581 flood peaks and significant fluctuations. Similar conclusions were drawn in the studies conducted
582 (Cui et al., 2021; Ding et al., 2020). The longer the foresight period, the lower the correlation
583 between input and output variables. The models face increased difficulty due to the lack of future
584 information and the challenges associated with flood forecasting.

585 The addition of an attention mechanism effectively enhances the accuracy of flood forecasting
586 in the original LSTM model. As the lead time increases, the temporal weights gradually shift



587 forward, causing the model to pay greater attention to values closer to the current moment. This
588 finding aligns with the conclusions of studies on temporal attention conducted by Ding et al. (2020)
589 and Wang et al. (2023). However, there is a difference between their studies and the current one, as
590 they incorporated a spatial attention module to focus on the relevance of spatial locations, while this
591 study introduces a feature attention module to highlight the importance of different input features in
592 flood forecasting.

593 Incorporating physical constraints into the model enhances the understanding of the monotonic
594 relationships between variables in the flooding process and improves the physical consistency of
595 the model. This study considers the monotonic relationships between precipitation, evaporation,
596 initial soil moisture content, and runoff in the watershed. In a study by Xie et al. (2021), three
597 physical conditions related to the rainfall-runoff forecasting process were encoded into the loss
598 function at the daily scale. Experimental results on 531 watersheds in the CAMELS dataset showed
599 that the model achieved an improvement from 0.52 to 0.61 in the NSE mean compared to the LSTM
600 model. In this study, flood forecasting is performed at a finer time scale, specifically at the hourly
601 scale, and additional monotonic relationship constraints between evapotranspiration, initial soil
602 water content, and runoff are incorporated.

603 Flood forecasting is challenged by various complex factors such as meteorological conditions
604 and rainfall patterns, and the uncertainty of these factors increases over time (Cheng et al., 2021;
605 Hu et al., 2019). Consequently, the model is prone to significant prediction errors. When the forecast
606 period extends to 6 hours, each model exhibits a significant deviation from the observed hydrograph
607 and more anomalous fluctuations. In this study, the maximum prediction period of the model is set
608 at 6 hours, and the effects of longer prediction periods need further investigation. In future research,
609 we propose exploring additional methods to address these limitations and enhance the performance
610 of our model. One potential avenue is the incorporation of error correction methods such as K
611 nearest neighbor (KNN) and backpropagation (BP) algorithms. Additionally, data assimilation
612 techniques, such as ensemble Kalman filter and particle filter methods, can be used to assimilate the
613 latest observed data and improve real-time forecasting accuracy. These approaches have the
614 potential to extend the forecasting period of flood prediction.



615 **6. Conclusions**

616 This research introduces a DL model called PHY-FTMA-LSTM, which combines feature-time-
617 based multi-head attention mechanisms with physical constraints. The primary aim is to explore
618 how incorporating interpretability and physical constraints into DL models affects flood forecasting
619 accuracy. The evaluation of the flood forecasting results from 1 to 6 h during the foresight period in
620 the Luan River basin yields the following conclusions:

621 (1) The attention mechanism that considers both features and time effectively enhances the
622 model's prediction performance, surpassing that of the original LSTM model. The FTMA-LSTM
623 model, equipped with an increased number of attention heads, further improves accuracy by
624 considering more information through parallel computation. Taking the integration of physical
625 constraints into account, the PHY-FTMA-LSTM model achieves the best performance, exhibiting
626 stable results. For a lead time of $t+1$, the NSE, KGE, R^2 , RMSE, and MAE reaches 0.988, 0.984,
627 0.988, 11.50, and 4.26, respectively. Additionally, NSE, KGE, and R^2 also could reach 0.908, 0.905,
628 and 0.911 for a lead time of $t+6$.

629 (2) The incorporation of a feature-time-based multi-head attention mechanism improves
630 interpretability by directing attention to the most valuable features and historical moments within
631 the inputs. The weight matrix visualization reveals that runoff emerges as the most influential feature
632 in flood forecasting, followed by rainfall, and finally initial soil moisture and evapotranspiration.
633 Furthermore, the weight distribution becomes more concentrated with increasing lead time.

634 (3) The model combines physical constraints by considering the monotonic relationships
635 between rainfall, evapotranspiration, initial soil moisture, and runoff at an hourly scale. This
636 augmentation significantly improves the model's predictive capacity for flood processes, including
637 flood peaks, while reducing the lag time.

638 In this study, we have successfully incorporated both the attention mechanism and physical
639 mechanism into a DL model to improve the accuracy of flood prediction while ensuring
640 interpretability and physical consistency. In future research, we recognize that there is room for
641 further enhancing the interpretability of our model. We suggest exploring alternative interpretation
642 techniques to gain deeper insights into the model's decision-making process. Furthermore, the
643 combination of physical mechanisms and DL models can be expanded by incorporating more



644 detailed basin subsurface information and exploring different integration methods that consider both
645 physical mechanisms and DL models.

646 **Code and data availability**

647 The rainfall and flood data and model codes used in this study could be available online
648 (https://github.com/zran1/PHY_FTMA_LSTM.git). The evapotranspiration and initial soil moisture
649 data are extracted from GLDAS Noah Land Surface Model (Beaudoin et al., 2019; D. Beaudoin
650 et al., 2020), which is freely available at <https://disc.gsfc.nasa.gov/datasets>.

651 **Author contributions**

652 Ting Zhang: Conceptualization, Methodology, Writing-original draft, Writing-review &
653 editing. Ran Zhang: Conceptualization, Methodology, Software, Validation, Writing-original draft.
654 Jianzhu Li: Validation, Writing-review & editing. Ping Feng: Validation, Writing-review & editing.

655 **Competing interests**

656 The contact author has declared that none of the authors has any competing interests.

657 **Acknowledgements**

658 This work was supported by the National Key Research and Development Program of China
659 (2023YFC3006501, 2023YFC3006503), National Natural Science Foundation of China (No.
660 52279022, 52079086).

661 **References:**

- 662 Beaudoin, H. and M. Rodell, NASA/GSFC/HSL (2019), GLDAS Noah Land Surface Model L4 3
663 hourly 0.25 x 0.25 degree V2.0 [Dataset]. Greenbelt, Maryland, USA, Goddard Earth Sciences
664 Data and Information Services Center (GES DISC).
665 https://disc.gsfc.nasa.gov/datasets/GLDAS_NOAH025_3H_2.0
- 666 Beaudoin, H., M. Rodell, NASA/GSFC/HSL (2020), GLDAS Noah Land Surface Model L4 3 hourly
667 0.25 x 0.25 degree V2.1 [Dataset]. Greenbelt, Maryland, USA, Goddard Earth Sciences Data and
668 Information Services Center (GES DISC).
669 https://disc.gsfc.nasa.gov/datasets/GLDAS_NOAH025_3H_2.1
- 670 Birkholz, S., Muro, M., Jeffrey, P., Smith, H. M., 2014. Rethinking the relationship between flood risk



- 671 perception and flood management. *Sci. Total Environ.* 478, 12-20.
672 <http://doi.org/10.1016/j.scitotenv.2014.01.061>
- 673 Brauwiers, G., Frasincar, F., 2023. A General Survey on Attention Mechanisms in Deep Learning. *Ieee*
674 *Trans. Knowl. Data Eng.* 35(4), 3279-3298. <http://doi.org/10.1109/TKDE.2021.3126456>
- 675 Cao, Q., Zhang, H., Zhu, F., Hao, Z., Yuan, F., 2022. Multi-step-ahead flood forecasting using an
676 improved BiLSTM-S2S model. *J. Flood Risk Manag.* 15(e128274)
677 <http://doi.org/10.1111/jfr3.12827>
- 678 Chen, Y., Ren, Q., Huang, F., Xu, H., Cluckie, I., 2011. Liuxihe Model and Its Modeling to River
679 Basin Flood. *J. Hydrol. Eng.* 16(1), 33-50. [http://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000286](http://doi.org/10.1061/(ASCE)HE.1943-5584.0000286)
- 680 Cheng, M., Fang, F., Navon, I. M., Pain, C. C., 2021. A real-time flow forecasting with deep
681 convolutional generative adversarial network: Application to flooding event in Denmark. *Phys.*
682 *Fluids.* 33(0566025) <http://doi.org/10.1063/5.0051213>
- 683 Cui, Z., Zhou, Y., Guo, S., Wang, J., Ba, H., He, S., 2021a. A novel hybrid XAJ-LSTM model for
684 multi-step-ahead flood forecasting. *Hydrol. Res.* 52(6), 1436-1454.
685 <http://doi.org/10.2166/nh.2021.016>
- 686 Ding, Y., Zhu, Y., Feng, J., Zhang, P., Cheng, Z., 2020. Interpretable spatio-temporal attention LSTM
687 model for flood forecasting. *Neurocomputing.* 403, 348-359.
688 <http://doi.org/https://doi.org/10.1016/j.neucom.2020.04.110>
- 689 Duan, Q., Sorooshian, S., Gupta, V., 1992. Effective and efficient global optimization for conceptual
690 rainfall-runoff models. *Water Resour. Res.*
- 691 Grillakis, M. G., Koutroulis, A. G., Komma, J., Tsanis, I. K., Wagner, W., Bloeschl, G., 2016. Initial
692 soil moisture effects on flash flood generation - A comparison between basins of contrasting
693 hydro-climatic conditions. *J. Hydrol.* 541(SIA), 206-217.
694 <http://doi.org/10.1016/j.jhydrol.2016.03.007>
- 695 Hu, R., Fang, F., Pain, C. C., Navon, I. M., 2019. Rapid spatio-temporal flood prediction and
696 uncertainty quantification using a deep learning method. *J. Hydrol.* 575, 911-920.
697 <http://doi.org/10.1016/j.jhydrol.2019.05.087>
- 698 Jiang, S., Zheng, Y., Solomatine, D., 2020. Improving AI System Awareness of Geoscience
699 Knowledge: Symbiotic Integration of Physical Approaches and Deep Learning. *Geophys. Res.*
700 *Lett.* 47(e2020GL08822913) <http://doi.org/10.1029/2020GL088229>



- 701 Kao, I., Zhou, Y., Chang, L., Chang, F., 2020. Exploring a Long Short-Term Memory based Encoder-
702 Decoder framework for multi-step-ahead flood forecasting. *J. Hydrol.* 583(124631)
703 <http://doi.org/10.1016/j.jhydrol.2020.124631>
- 704 Karniadakis, G. E., Kevrekidis, I. G., Lu, L., Perdikaris, P., Wang, S., Yang, L., 2021. Physics-
705 informed machine learning. *Nat. Rev. Phys.* 3(6), 422-440. [http://doi.org/10.1038/s42254-021-](http://doi.org/10.1038/s42254-021-00314-5)
706 00314-5
- 707 Kellens, W., Terpstra, T., De Maeyer, P., 2013. Perception and Communication of Flood Risks: A
708 Systematic Review of Empirical Research: Perception and Communication of Flood Risks. *Risk*
709 *Anal.* 33(1), 24-49. <http://doi.org/10.1111/j.1539-6924.2012.01844.x>
- 710 Leedal, D., Weerts, A. H., Smith, P. J., Beven, K. J., 2013. Application of data-based mechanistic
711 modelling for flood forecasting at multiple locations in the Eden catchment in the National Flood
712 Forecasting System (England and Wales). *Hydrol. Earth Syst. Sci.* 17(1), 177-185.
713 <http://doi.org/10.5194/hess-17-177-2013>
- 714 Lima, A. R., Cannon, A. J., Hsieh, W. W., 2016. Forecasting daily streamflow using online sequential
715 extreme learning machines. *J. Hydrol.* 537, 431-443. <http://doi.org/10.1016/j.jhydrol.2016.03.017>
- 716 Luppichini, M., Barsanti, M., Giannechini, R., Bini, M., 2022. Deep learning models to predict flood
717 events in fast-flowing watersheds. *Sci. Total Environ.* 813(151885)
718 <http://doi.org/10.1016/j.scitotenv.2021.151885>
- 719 Lv, N., Liang, X., Chen, C., Zhou, Y., Li, J., Wei, H., Wang, H., 2020. A long Short-Term memory
720 cyclic model with mutual information for hydrology forecasting: A Case study in the xixian basin.
721 *Adv. Water Resour.* 141(103622) <http://doi.org/10.1016/j.advwatres.2020.103622>
- 722 Mourato, S., Fernandez, P., Marques, F., Rocha, A., Pereira, L., 2021. An interactive Web-GIS fluvial
723 flood forecast and alert system in operation in Portugal. *Int. J. Disaster Risk Reduct.* 58(102201)
724 <http://doi.org/10.1016/j.ijdrr.2021.102201>
- 725 Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., Prieto, C., Gupta,
726 H. V., 2021. What Role Does Hydrological Science Play in the Age of Machine Learning?. *Water*
727 *Resour. Res.* 57(e2020WR0280913) <http://doi.org/10.1029/2020WR028091>
- 728 Niu, Z., Zhong, G., Yu, H., 2021. A review on the attention mechanism of deep learning.
729 *Neurocomputing.* 452, 48-62. <http://doi.org/10.1016/j.neucom.2021.03.091>
- 730 Rahimzad, M., Moghaddam Nia, A., Zolfonoon, H., Soltani, J., Danandeh Mehr, A., Kwon, H., 2021.



- 731 Performance Comparison of an LSTM-based Deep Learning Model versus Conventional Machine
732 Learning Algorithms for Streamflow Forecasting. *Water Resour. Manag.* 35(12), 4167-4187.
733 <http://doi.org/10.1007/s11269-021-02937-w>
- 734 Read, J. S., Jia, X., Willard, J., Appling, A. P., Zwart, J. A., Oliver, S. K., Karpatne, A., Hansen, G. J.
735 A., Hanson, P. C., Watkins, W., Steinbach, M., Kumar, V., 2019. Process-Guided Deep Learning
736 Predictions of Lake Water Temperature. *Water Resour. Res.* 55(11), 9173-9190.
737 <http://doi.org/10.1029/2019WR024922>
- 738 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvahalais, N., Prabhat., 2019.
739 Deep learning and process understanding for data-driven Earth system science. *Nature*.
740 566(7743), 195-204. <http://doi.org/10.1038/s41586-019-0912-1>
- 741 Song, S., Lan, C., Xing, J., Zeng, W., Liu, J., 2017. An End-to-End Spatio-Temporal Attention Model
742 for Human Action Recognition from Skeleton Data. THIRTY-FIRST AAAI CONFERENCE ON
743 ARTIFICIAL INTELLIGENCE, San Francisco, CA.
- 744 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., Polosukhin, I.,
745 2017. Attention Is All You Need. *ADVANCES IN NEURAL INFORMATION PROCESSING*
746 *SYSTEMS 30 (NIPS 2017)*, Long Beach, CA.
- 747 Wan, X., Yang, Q., Jiang, P., Zhong, P., 2019. A Hybrid Model for Real-Time Probabilistic Flood
748 Forecasting Using Elman Neural Network with Heterogeneity of Error Distributions. *Water*
749 *Resour. Manag.* 33(11), 4027-4050. <http://doi.org/10.1007/s11269-019-02351-3>
- 750 Wang, N., Zhang, D., Chang, H., Li, H., 2020. Deep learning of subsurface flow via theory-guided
751 neural network. *J. Hydrol.* 584, 124700.
752 <http://doi.org/https://doi.org/10.1016/j.jhydrol.2020.124700>
- 753 Wang, Y., Huang, Y., Xiao, M., Zhou, S., Xiong, B., Jin, Z., 2023. Medium-long-term prediction of
754 water level based on an improved spatio-temporal attention mechanism for long short-term
755 memory networks. *J. Hydrol.* 618(129163) <http://doi.org/10.1016/j.jhydrol.2023.129163>
- 756 Xiang, Z., Yan, J., Demir, I., 2020. A Rainfall-Runoff Model With LSTM-Based Sequence-to-
757 Sequence Learning. *Water Resour. Res.* 56(e2019WR0253261)
758 <http://doi.org/10.1029/2019WR025326>
- 759 Xie, K., Liu, P., Zhang, J., Han, D., Wang, G., Shen, C., 2021. Physics-guided deep learning for
760 rainfall-runoff modeling by considering extreme events and monotonic relationships. *J. Hydrol.*



- 761 603, 127043. <https://doi.org/https://doi.org/10.1016/j.jhydrol.2021.127043>
- 762 Xu, Y., Hu, C., Wu, Q., Li, Z., Jian, S., Chen, Y., 2021. Application of temporal convolutional network
763 for flood forecasting. *Hydrol. Res.* 52(6), 1455-1468. <http://doi.org/10.2166/nh.2021.021>
- 764 Yang, S., Yang, D., Chen, J., Santisirisomboon, J., Lu, W., Zhao, B., 2020. A physical process and
765 machine learning combined hydrological model for daily streamflow simulations of large
766 watersheds with limited observation data. *J. Hydrol.* 590(125206)
767 <http://doi.org/10.1016/j.jhydrol.2020.125206>
- 768 Yokoo, K., Ishida, K., Ercan, A., Tu, T., Nagasato, T., Kiyama, M., Amagasaki, M., 2022. Capabilities
769 of deep learning models on learning physical relationships: Case of rainfall-runoff modeling with
770 LSTM. *Sci. Total Environ.* 802(149876) <http://doi.org/10.1016/j.scitotenv.2021.149876>
- 771 Yu, P., Chen, S., Chang, I., 2006. Support vector regression for real-time flood stage forecasting. *J.*
772 *Hydrol.* 328(3-4SI), 704-716. <http://doi.org/10.1016/j.jhydrol.2006.01.021>
- 773 Zhang, H., Singh, V. P., Wang, B., Yu, Y., 2016. CEREF: A hybrid data-driven model for forecasting
774 annual streamflow from a socio-hydrological system. *J. Hydrol.* 540, 246-256.
775 <http://doi.org/10.1016/j.jhydrol.2016.06.029>
- 776 Zhang, M., Su, H., Wen, J., 2021. Classification of flower image based on attention mechanism and
777 multi-loss attention network. *Comput. Commun.* 179, 307-317.
778 <http://doi.org/10.1016/j.comcom.2021.09.001>
- 779 Zhao, Z., Chen, W., Wu, X., Chen, P. C. Y., Liu, J., 2017. LSTM network: a deep learning approach
780 for short-term traffic forecast. *Iet Intell. Transp. Syst.* 11(2), 68-75. [http://doi.org/10.1049/iet-](http://doi.org/10.1049/iet-its.2016.0208)
781 [its.2016.0208](http://doi.org/10.1049/iet-its.2016.0208)
- 782 Zhu, X., Lu, C., Wang, R., Bai, J., 2005. Artificial neural network model for flood water level
783 forecasting. *J. Hydraul. Eng.-Asce.* 36(0559-9350(2005)36:7<806:JYRGSJ>2.0.TX;2-O7), 806-
784 811.