



# 1 Extending Medium-Range Global Flood 2 Forecasts: The Google Global Flood 3 Forecasting Model Version 2

4 Deborah Cohen <sup>1</sup>, Rony Amira <sup>1</sup>, Rom Aschner <sup>1</sup>, Yuval Carny <sup>1</sup>, Ben Feinstein <sup>1</sup>,  
5 Hadas Fester <sup>1</sup>, Shmulik Fronman <sup>2</sup>, Martin Gauch <sup>3</sup>, Oren Gilon <sup>1</sup>, Rotem Green <sup>1</sup>,  
6 Avinatan Hassidim <sup>1</sup>, Daniel Klotz <sup>4</sup>, Frederik Kratzert <sup>4</sup>, Dan Korenfeld <sup>1</sup>, Gila Loike <sup>1</sup>,  
7 Amit Markel <sup>1</sup>, Yossi Matias <sup>5</sup>, Rotem Mayo <sup>1</sup>, Asher Metzger <sup>1</sup>, Benny Mosheyev <sup>1</sup>,  
8 Aviel Niego <sup>1</sup>, Stephanie Rees <sup>6</sup>, Emily Reinstein <sup>7</sup>, Amitay Sicherman <sup>1</sup>, Guy Shalev <sup>8</sup>,  
9 Omri Shefi <sup>1</sup>, Yuval Shildan <sup>1</sup>, Ido Zemach <sup>1</sup>, Oleg Zlydenko <sup>1</sup>, Grey Nearing <sup>3,\*</sup>

10

11 <sup>1</sup>Google Research, Electra Tower, Yigal Alon St 98, Tel Aviv-Yafo, 6789141, Israel

12 <sup>2</sup>Google Research, 6 Pancras Sq, London N1C 4AG, United Kingdom

13 <sup>3</sup>Google Research, Brandschenkestrasse 110, 8002 Zürich, Switzerland

14 <sup>4</sup>Google Research, Graben 19, 1010 Wien, Austria

15 <sup>5</sup>Google Research, 1600 Amphitheatre Pkwy, Mountain View, CA 94043 USA

16 <sup>6</sup>Google, 1-13 St Giles High St, London WC2H 8AG, United Kingdom

17 <sup>7</sup>Google, 111 8th Ave, New York, NY 10011 USA

18 <sup>8</sup>Google DeepMind, Electra Tower, Yigal Alon St 98, Tel Aviv-Yafo, 6789141, Israel

19

20 \*Corresponding Author: Grey Nearing <[nearing@google.com](mailto:nearing@google.com)>

21

## 22 Abstract

23 This paper evaluates an updated flood forecasting system that significantly extends reliable  
24 lead times. We evaluated this updated model (v2) against the prior system (v1) and established  
25 third-party benchmarks across 1,223 global test basins. The primary finding is that the v2  
26 system extends the reliable predictive horizon by 6 days in gauged basins and 1 day in  
27 ungauged basins relative to the v1 nowcast, as measured by the Nash Sutcliffe Efficiency. Along  
28 with this paper, we release an open-source codebase for training both the v1 and v2 forecast  
29 models with the open-source Caravan dataset.

## 30 1 Introduction

31 The hydrology research and operations communities are now several years into the adoption  
32 cycle for machine learning (ML) based simulation and prediction models. The initial phase of  
33 this adoption cycle was characterized by establishing benchmarks relative to traditional  
34 conceptual rainfall-runoff models (Kratzert et al., 2019; Mai et al., 2022; Arsenault et al., 2022),  
35 particularly in data-scarce and ungauged settings (Kratzert et al., 2019a; Prieto et al., 2019;



1 Feng et al., 2023; Nearing et al., 2024). Since then, the field has matured as research diversified  
2 into questions related to interpretability (e.g., Kratzert et al., 2019; Lees et al., 2022; Başağaoğlu  
3 et al., 2022; De la Fuente et al., 2023; Dai et al., 2023; Hosseini et al., 2025; Rosati et al., 2025),  
4 uncertainty quantification (e.g., Klotz et al., 2021; Chaudhary et al., 2022; Nourani et al., 2022; Liu  
5 et al., 2023), data assimilation (e.g., Feng et al., 2020; Nearing et al., 2021; Deng et al., 2023;  
6 Sabzipour et al., 2023), and hybrid process-based machine learning (e.g., Hoedt et al., 2021; Jia  
7 et al., 2021; Tsai et al., 2021; Feng et al., 2022; Höge et al., 2022; Frame et al., 2023; Acuña  
8 Espinoza et al., 2025).

9 We use operational machine learning hydrology models for global-scale riverine flood  
10 forecasting (Nevo et al., 2022; Nearing et al., 2024). This paper demonstrates a synthesis of  
11 approaches to help mitigate data-related issues and increase the reliability of streamflow  
12 forecasts. We specifically evaluate the technical implementation of version 2 (v2) of the Google  
13 Global Flood Forecasting system. This system is run operationally and powers the river  
14 forecasting component of Google's FloodHub ([g.co/floodhub](https://g.co/floodhub)).

15 Upgrades between v1 (Nevo et al., 2022) and v2 systems primarily focus on improvements  
16 around three data-related challenges related to training data availability, temporally limited data  
17 records, and input data distribution shifts.

18 This paper serves two primary purposes: (1) to provide transparency to the scientific  
19 community about operational flood forecasting system progress and challenges, and (2) to  
20 facilitate research around ML-based flood forecasting by providing open source resources. The  
21 GoogleHydrology open-source code repository  
22 (<https://github.com/google-research/flood-forecasting>) implements a version of the model  
23 architectures and training pipeline evaluated in this paper, and we additionally released a  
24 multi-decade historical data archive, called the Google Runoff Reanalysis and Reforecast  
25 (GRRR) dataset, which contains historical simulations and historical reforecasts at over 1 million  
26 locations globally.

## 27 **2 Model Description**

### 28 **2.1 Data**

29 The river forecast model is trained to predict daily average streamflow at a basin outlet. It uses  
30 two types of input data: static catchment attributes, which describe physical characteristics  
31 that are implicitly assumed to be time-invariant (e.g., topography, soil type, land cover), and  
32 dynamic meteorological forcings, which represent the time-varying weather conditions (e.g.,  
33 precipitation, temperature) driving the hydrological response. Sources for these three types of  
34 data are described in the following subsections.

#### 35 **2.1.1 Static Catchment Attributes**



1 We use a set of 92 static catchment attributes to characterize the time-invariant physical  
 2 properties of each basin. These attributes are spatially averaged over catchment area and  
 3 derived from the HydroATLAS dataset (Linke et al., 2019) plus hydro-climatic statistics  
 4 calculated from ERA5-Land reanalysis (Muñoz-Sabater et al., 2021). A detailed description of all  
 5 attributes is provided in **Table 1**. The Caravan GitHub repository  
 6 (<https://github.com/kratzert/Caravan>) contains the code to calculate all ERA5 and climate  
 7 attributes.

8 At a conceptual level, static attributes provide information about:

- 9 • **Topography:** Mean elevation, slope, and relief characteristics.
- 10 • **Climate Indices:** Aridity, seasonality, and precipitation frequency/duration statistics.
- 11 • **Land Cover:** Fractions of forest, urban areas, water bodies, snow cover, and various  
 12 vegetation classes.
- 13 • **Soil Properties:** Soil texture (sand, silt, clay content), organic carbon, and hydraulic  
 14 properties.
- 15 • **Anthropogenic Factors:** Human footprint index, GDP, population density, and road  
 16 density.

17 **Table 1:** Full list of static catchment attributes incorporated into the model as input features.

Attribute Name	Description
<b>HydroATLAS Attributes</b>	
aet_mm_uyr	Actual evapotranspiration; multi-year annual average (mm/yr).
ari_ix_uav	Aridity index; multi-year average.
crp_pc_use	Cropland extent; percentage of sub-basin area.
ele_mt_uav	Elevation; multi-year average in meters above sea level.
ero_kh_uav	Soil erosion; multi-year average in kg/ha/yr.
for_pc_use	Forest cover; percentage of sub-basin area.
gdp_ud_usu	Gross Domestic Product (GDP); total sum for sub-basin (USD).



gla_pc_use	Glacier cover; percentage of sub-basin area.
glc_pc_u01 - glc_pc_u22	Land cover classes (e.g., forests, grasslands, urban, water bodies); percentage of sub-basin area. See full list for specific classes.
hft_ix_u09	Human Footprint index for year 2009 (multi-year average).
hft_ix_u93	Human Footprint index for year 1993 (multi-year average).
inu_pc_ult	Inundation extent (long-term); percentage of sub-basin area.
inu_pc_umn	Inundation extent (minimum); percentage of sub-basin area.
inu_pc_umx	Inundation extent (maximum); percentage of sub-basin area.
ire_pc_use	Irrigated area; percentage of sub-basin area.
kar_pc_use	Karst area; percentage of sub-basin area.
lka_pc_use	Lake area; percentage of sub-basin area.
lkv_mc_usu	Lake volume; total sum in sub-basin (million m <sup>3</sup> ).
nli_ix_uav	Nighttime lights index; multi-year average.
pac_pc_use	Protected area cover; percentage of sub-basin area.
pet_mm_uyr	Potential evapotranspiration; multi-year annual average (mm/yr).
pnv_pc_u01 - pnv_pc_u15	Potential Natural Vegetation classes; percentage of sub-basin area.
pop_ct_usu	Population count; total sum within sub-basin.



ppd_pk_uav	Population density; multi-year average (people/km <sup>2</sup> ).
pre_mm_uyr	Precipitation; multi-year annual average (mm/yr).
prm_pc_use	Permafrost extent; percentage of sub-basin area.
pst_pc_use	Pasture extent; percentage of sub-basin area.
rdd_mk_uav	Road density; multi-year average (m/km <sup>2</sup> ).
rev_mc_usu	Reservoir volume; total sum in sub-basin (million m <sup>3</sup> ).
ria_ha_usu	Reservoir area; total sum in sub-basin (hectares).
riv_tc_usu	River volume; total sum in sub-basin (thousand m <sup>3</sup> ).
snw_pc_uyr	Snow cover; percentage multi-year annual average.
swc_pc_uyr	Soil water content; percentage multi-year annual average.
tmp_dc_uyr	Air temperature; multi-year annual average (degrees Celsius).
urb_pc_use	Urban area; percentage of sub-basin area.
wet_pc_u01 - wet_pc_u09	Wetland classes (e.g., lakes, reservoirs, rivers, marshes); percentage of sub-basin area.
wet_pc_ug1, wet_pc_ug2	Global wetland classes (Lakes/Reservoirs, Floodplains/Marshes).
<b>ERA5-Land Attributes</b>	
aridity	Aridity index derived from ERA5-Land reanalysis.
frac_snow	Fraction of precipitation falling as snow.



high_prec_dur	Average duration of high precipitation events.
high_prec_freq	Frequency of high precipitation events.
low_prec_dur	Average duration of low precipitation events.
low_prec_freq	Frequency of low precipitation events.
moisture_index	Moisture index derived from ERA5-Land.
p_mean	Long-term mean precipitation.
pet_mean	Long-term mean potential evapotranspiration from ERA5-Land.
seasonality	Seasonality index of precipitation from ERA5-Land.
<b>Derived Attributes</b>	
area	Calculated drainage area.

<sup>1</sup> Climate variables derived from FAO evapotranspiration used the guidelines outlined by Allen et al., (1998).

## 2.1.2 Dynamic Meteorological Inputs

Kratzert et al. (2021) found that using multiple meteorological data sources in a single LSTM model provides performance improvements. Accordingly, our model ingests daily time-series data from multiple global weather data sources. Historical records of our dynamic meteorological forcings are available from the open source Caravan MultiMet extension (Shalev et al., 2024) in sufficient quantity to train global streamflow models with the GoogleHydrology codebase.

The dynamic inputs are:

### 1. ECMWF HRES (European Centre for Medium-Range Weather Forecasts):

- High-Resolution operational deterministic forecasts:
  - Total precipitation
  - 2-meter temperature
  - Surface net solar radiation
  - Surface net thermal radiation



- 1           ■ Surface pressure
- 2           ■ Snow fraction
- 3 2. **CPC (NOAA Climate Prediction Center):** Global Unified Gauge-Based Analysis of Daily
- 4     Precipitation (Chen et al., 2008).
- 5       ○ Precipitation
- 6 3. **GraphCast (Google DeepMind):** AI-based medium-range global weather forecasting
- 7     mode (Lam et al., 2023)l. Variables include:
- 8       ○ Total precipitation
- 9       ○ 2-meter temperature
- 10 4. **IMERG (NASA):** Integrated Multi-satellite Retrievals for GPM Early Run (Huffman et al.,
- 11     2020).
- 12       ○ Precipitation

13 All meteorological data are aggregated from their native temporal resolutions to daily. The  
14 method of aggregation depends on the physical units of each data variable (rates vs.  
15 quantities).

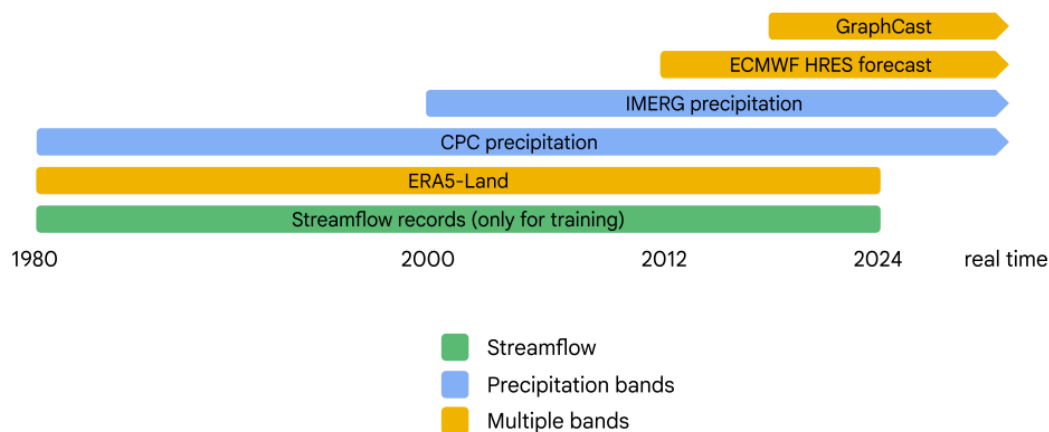
16 The model uses these inputs in two distinct modes: hindcast and forecast. The hindcast period  
17 represents observed historical weather up to the issue time of a forecast and is used to spin up  
18 the model's states. For this phase, we use ECMWF HRES and GraphCast weather forecast data  
19 products in hindcast mode by extracting the 0-day lead time forecast for each issue time. Only  
20 ECMWF HRES and GraphCast products are available as true forecasts, meaning that CPC and  
21 IMERG are not used in the forecast phase of the model. The distinction between the hindcast  
22 and forecast phases of our model is analogous to the difference between Quantitative  
23 Precipitation Estimates (QPE) and Quantitative Precipitation Forecasts (QPF) in operational  
24 hydrology (Liu et al., 2012).

25 As illustrated in **Figure 1**, meteorological input products have varying availability timelines for  
26 training. While streamflow labels and reanalysis proxies (ERA5-Land) are available for several  
27 decades (we make an arbitrary decision to limit data records to the period 1980 - 2024),  
28 operational forecast archives for HRES and GraphCast begin later (approx. 2012 and 2016,  
29 respectively).

30 To take advantage of the full historical streamflow record (1980–2024) for training, we employ  
31 a data imputation strategy called feature unioning. For time periods where operational  
32 forecasts are unavailable from either HRES or Graphcast, we substitute these with ERA5-Land  
33 reanalysis data. Because HRES shares the same underlying physical model as ERA5, and  
34 GraphCast is trained on ERA5, the assumption is that statistical distributions of the reanalysis



1 data serve as an effective proxy for the forecast inputs during the earlier training years.



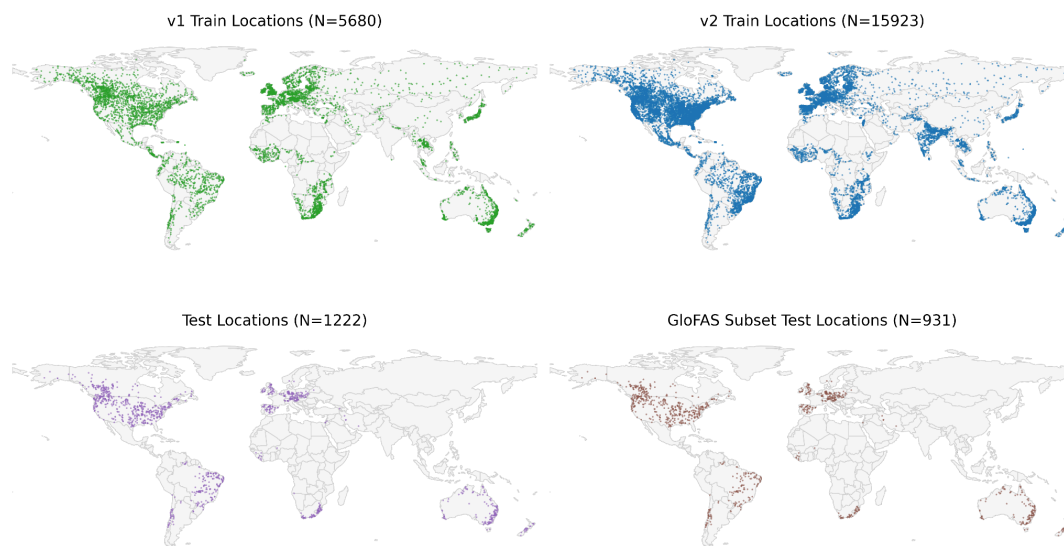
2

3 **Figure 1:** Temporal availability of the dynamic input datasets and target streamflow data.

#### 4 2.1.3 Streamflow Target Data

5 The target variable is daily streamflow, obtained for training the v2 system from the Caravan  
6 dataset (Kratzert et al., 2023), Global Runoff Data Center (GRDC; Färber et al., 2024), and  
7 BANDAS (Solís-Alvarado et al., 2015). The v1 benchmark was trained on data from the GRDC. All  
8 GRDC test gauges from the v1 evaluation by Nearing et al. (2024) are present in the v2 and  
9 third-party benchmark test sets, ensuring a direct comparison. The global spatial distribution of  
10 these expanded training and testing gauges is illustrated in **Figure 2**.

11 It is important to acknowledge that Caravan is an open-community dataset based on data from  
12 regional large-sample hydrology datasets that are often referred to as Catchment Attributes  
13 and Meteorology for Large Sample studies (CAMELS) datasets (Newman et al., 2015; Addor et  
14 al., 2017). Caravan relies fundamentally on contributions like those from Denmark (Liu et al.,  
15 2024), Chile (Alvarez-Garretón et al., 2018), Iceland (Helgason et al., 2024), Israel (Morin, 2023),  
16 Switzerland (Höge et al., 2023), Spain (Casado Rodríguez, 2023), and the GRDC (Färber et al.,  
17 2024).



1

2 **Figure 2:** Spatial distribution of training and testing gauges globally. The top row illustrates the  
3 expansion of training data locations from the v1 dataset to the expanded v2 dataset due to  
4 Caravan. The bottom row maps the specific holdout test set (left pane), and the fraction of  
5 those test locations where we have GloFAS data.

## 6 2.2 Model Architectures

7 The model used in our v1 system is an encoder-decoder (ED-LSTM) architecture described by  
8 Nevo et al. (2022) and evaluated by Nearing et al. (2024). The model used in the v2 system is a  
9 Mean Embedding (ME-LSTM) architecture developed by Gauch et al. (2025). The ME-LSTM is  
10 more robust to missing data and addresses a forecast initialization challenge described in the  
11 following subsection.

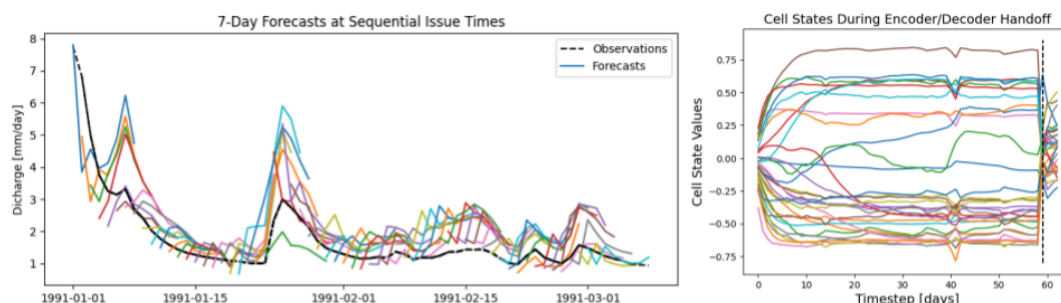
### 12 2.2.1 Encoder-Decoder LSTM (ED-LSTM)

13 The ED-LSTM runs separate LSTMs for the hindcast vs. forecast periods. The hidden and cell  
14 states of the hindcast LSTM are used to initialize the forecast LSTM through a small neural  
15 network. This allows for different weights in each of the LSTM gates for hindcast vs. forecast  
16 meteorological inputs, and addresses the challenge of distribution differences between  
17 observation-based or analysis-based features vs. numerical weather predictions. This model is  
18 described in detail in prior publications (Nevo et al., 2022; Nearing et al., 2024).

19 One motivation for the development of a new model architecture for the v2 system was a  
20 specific failure mode observed in the ED-LSTM (Nearing et al., 2023a). Empirically, the handoff  
21 between the states of the hindcast and forecast LSTMs can cause a discontinuity during  
22 forecast initialization. A toy example illustrating this forecast artifact is shown in the left panel of  
23 **Figure 3**. The cause of this artifact is that the forecast LSTM is sometimes initialized with a cell



1 state that is outside of a steady-state distribution and the model spends the first few time  
2 steps of the forecast horizon annealing its internal states to the new input climatology rather  
3 than responding strictly to hydrological drivers. This is illustrated in the right panel of **Figure 3**.



4

5 **Figure 3.** Forecast initialization artifacts and their underlying cause. The left panel shows an  
6 example of an artifact in ED-LSTM forecasts and the right panel shows the cause of this  
7 artifact. This is a toy example for illustration purposes that uses a short data record from a  
8 randomly selected basin.

### 9 2.2.2 Mean Embedding Forecast LSTM (ME-LSTM)

10 The ME-LSTM addresses limitations of the ED-LSTM by using a stacked LSTM framework to  
11 handle input distribution shifts and embedding layers to handle missing inputs. The model  
12 architecture is described in detail by Gauch et al. (2025) and illustrated in **Figure 4**.

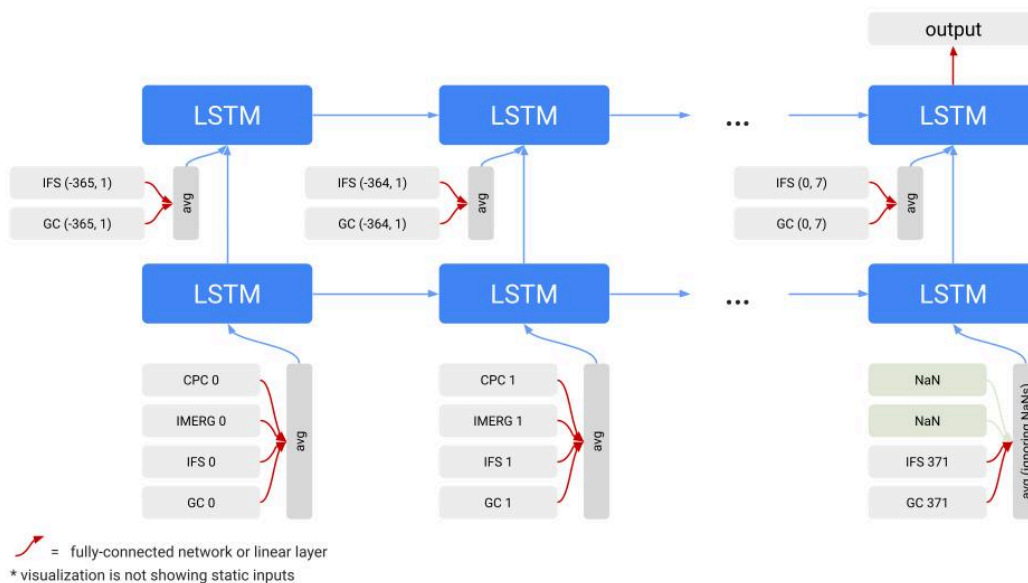
13 Instead of concatenating all meteorological inputs into a single vector at each timestep, which  
14 fails if any input is missing, the ME-LSTM treats different meteorological products (e.g., HRES,  
15 GraphCast, IMERG) as independent input groups. Each product group is processed by a  
16 dedicated embedding network that projects inputs into a shared latent space. Embedding  
17 networks are multilayer perceptrons. Static catchment attributes are concatenated to dynamic  
18 inputs of each product group before embedding. These embeddings are aggregated using a  
19 masked mean operation, which computes the average of the embeddings at each time step,  
20 ignoring any product that contains missing data.

21 The ME-LSTM then employs two LSTM layers arranged in a stack that run over the full time  
22 sequence (hindcast + forecast), eliminating the handoff discontinuity. The first layer processes  
23 a sequence of aggregated embeddings derived from hindcast data and runs over the full  
24 sequence, i.e., hindcast (spinup) + forecast sequence length. The second LSTM layer predicts  
25 streamflow. It receives two inputs at every time step: the aggregated forecast embeddings and  
26 the hidden state output from the first LSTM layer.

27 Both v1 and v2 models evaluated in this paper use a mixture density layer to make probabilistic  
28 predictions. At each location and timestep, each model predicts the parameters of a countable  
29 mixture of asymmetric Laplacians (CMAL) distribution (Klotz et al., 2022). We take the mean of



- 1 the predicted distribution to be the deterministic model prediction that we evaluate in this
- 2 paper.



3  
 4 **Figure 4:** Architecture of the Mean Embedding Forecast LSTM (ME-LSTM).

### 5 2.2.3 Hyperparameters

6 Model hyperparameters for both v1 and v2 models are provided in GoogleHydrology  
 7 configuration files released with this paper. Hyperparameters for the v2 system model were  
 8 chosen by approximate hyperparameter tuning using a procedure analogous to what was  
 9 described in Appendix A by Kratzert et al. (2024). We do not report this hypertuning  
 10 experiment in this paper, however users are free to do their own hyperparameter tuning with  
 11 the open source GoogleHydrology codebase provided. The most important hyperparameters  
 12 for the v2 model are detailed in **Table 2**.

13 **Table 2:** Main Hyperparameters of the ME-LSTM Model Architecture.

Component	Specification
Hindcast LSTM Hidden Size	512
Forecast LSTM Hidden Size	512
Sequence Length ( $T_{seq}$ )	365 days
Forecast Horizon ( $T_{lead}$ )	7 days



### Static Embedding Net

Network: Fully-Connected  
Depth: 3-Layer  
Neurons: 100, 100, 20  
Activations: tanh, tanh, linear

### Hindcast Embedding Nets

Network: Fully-Connected  
Depth: 2-Layer  
Neurons: 100, 20  
Activations: tanh, linear

### Forecast Embedding Nets

Network: Fully-Connected  
Depth: 4-Layer  
Neurons: 20, 20, 20, 20  
Activations: tanh, tanh, tanh, linear

## 1 2.3 Training

2 The v2 system model was trained using an Adam optimizer with a batch size of 512 with the  
3 following:

- 4 • **Loss Function:** We used the CMAL likelihood loss described by Klotz et al. (2022).
- 5 • **Noise Injection:** We injected Gaussian noise with a standard deviation of 0.005 into  
6 normalized target variables during training.
- 7 • **Learning Rate:** We used 0.0005.
- 8 • **Gradient Clipping:** Gradients were clipped to a norm of 1.0.
- 9 • **Input feature dropout:** We randomly dropped timeseries input features with probability  
10 0.1 to increase the robustness of our model against missing data, as described in Gauch et  
11 al. (2025).

12 We used 125 epochs with a batch limit of 2000 updates per epoch.

## 13 3 Model Evaluation

14 Our evaluation assesses forecast skill in both gauged and ungauged settings. Gauged  
15 prediction means that streamflow observations from the evaluation gauges are used in training  
16 (although out-of-sample in time), whereas ungauged prediction assumes no local streamflow  
17 data is used for training. For the ungauged setting, the v1 system used random k-fold (k=10)  
18 cross-validation, whereas the v2 system used a single holdout test set.

### 19 3.1 Confounding Factors in Operational System



## 1 Upgrades

2 A key limitation in evaluating operational system upgrades is the simultaneous evolution of  
3 multiple factors like data and architecture. Consequently, the performance deltas between the  
4 v1 and v2 systems are a compound effect from the architectural shift (ME-LSTM) and the  
5 expanded geographic training sample from Caravan. We additionally evaluate an ablated  
6 version of the v2 system without GraphCast to disaggregate performance gains from the new  
7 input meteorological forcing data.

## 8 3.2 Evaluation Metrics

9 Gauch et al. (2023) outlined a set of evaluation metrics for hydrograph prediction that can  
10 explain expert-driven model intercomparison. In the interest of tractability, here we used only a  
11 subset of those metrics that allow us to highlight particular characteristics of model  
12 performance. Since all of the model runs produced for this study are publicly available, readers  
13 are welcome to explore these results using different metrics.

14 In particular, we evaluated performance here using the Nash-Sutcliffe Efficiency (NSE; Nash &  
15 Sutcliffe, 1970) and a decomposition of the Kling-Gupta Efficiency (KGE; Gupta et al., 2009):

- 16 • **Nash-Sutcliffe Efficiency (NSE):** This is a widely used normalized statistic that  
17 determines the relative magnitude of the residual variance compared to the measured  
18 data variance. An NSE of 1.0 indicates perfect correspondence between predictions  
19 and observations. An NSE of 0 indicates that the model predictions are as accurate as  
20 the mean of the observed data, and an NSE less than zero indicates that the observed  
21 mean is a better predictor than the model.
- 22 • **Kling-Gupta Efficiency (KGE) Decomposition:** The KGE metric incorporates  
23 correlation, variability bias, and mean bias into a single metric. As with the NSE, a KGE  
24 metric of 1.0 indicates a perfect model, and a KGE of 0 indicates that the model  
25 predictions are as accurate as the observation mean. We look at these three  
26 constituent components independently:
  - 27 ○ **Correlation ( $r$ ):** The Pearson correlation coefficient evaluates the phase timing  
28 and temporal dynamics of the hydrograph.
  - 29 ○ **Bias Ratio ( $\beta$ ):** This ratio quantifies the long-term volumetric water balance error  
30 (mass conservation). This is calculated as the mean of the simulated predictions  
31 divided by the mean of the observed values.
  - 32 ○ **Variability Ratio ( $\gamma$ ):** This ratio measures the model's ability to capture the  
33 amplitude of flow extremes. It is calculated as the coefficient of variation of the  
34 simulated predictions divided by the coefficient of variation of the observed



1 values.

## 2 **3.3 Time Periods & Cross Validation**

3 Our benchmarking experiments assess model performance in both gauged and ungauged  
4 settings. When we say that a model is gauged or makes gauged predictions, we mean that the  
5 model is used to make predictions in the same locations where it received training data. When  
6 we say that a model is ungauged or makes ungauged predictions, we mean that the test  
7 locations are withheld from training. In both cases (gauged and ungauged), we use a temporal  
8 cross-validation split.

9 The test period for all models is January 1, 2016 through December 31, 2023. Both v1 and v2  
10 models were trained with cross validation over individual calendar years in that time period.  
11 Since both models use a hindcast or spin-up sequence length of 365 days, we withheld one  
12 year of training data before and after each hold-out year for each temporal cross-validation  
13 split to prevent data leakage. As an example, for testing on the year 2020, we held out from  
14 training all data during the period January 1, 2019 through December 31, 2021, and added  
15 predictions from the period January 1, 2020 through December 31, 2020 to the timeseries or  
16 test predictions that all metrics were calculated over. By repeating this strategy for each test  
17 year, we constructed out-of-sample predictions for the entire time period of 2016 - 2023. The  
18 training period for all models is 1980 - 2023, minus the appropriate 2 or 3 year holdout for each  
19 cross validation split.

20 Because we use existing model runs for all of the benchmarks, described in the next section,  
21 we constructed sets of test basins by choosing locations that exist in the v1 model runs by  
22 Nearing et al (2024), in the locations used for our v2 operational system, and in the  
23 GEOGLOWS dataset. This intersection resulted in 1,222 shared test basins, of which 931 also  
24 exist in the GloFAS test set. The bottom row of **Figure 2** shows these test locations. For the  
25 GloFAS benchmark, geolocation matching between streamflow gauges and the GloFAS river  
26 network was performed by ECMWF as part of the Nearing et al. (2024) publication.

27 The v2 models were trained using a single spatial split (with and without GraphCast) while the  
28 v1 models were trained using random k-fold (k=10) cross-validation as described by Nearing et  
29 al. (2024).

## 30 **3.4 Benchmarks**

31 To contextualize the performance of the v2 system, we compare our results against three  
32 reference systems: the v1 operational system, GloFAS, and GEOGLOWS.

33 The first benchmark is the version 1 (v1) operational system of the Google Global Flood  
34 Forecasting system. This system was described by Nevo et al. (2022) and evaluated by Nearing  
35 et al. (2024). The forecast model in this system is the ED-LSTM model described in **Section**



1 **2.2.1.** For this evaluation we used the open data release from Nearing et al. (2024).

2 The second benchmark is the Global Flood Awareness System (GloFAS), which serves as the  
3 global flood forecasting system of the Copernicus Emergency Management Service (CEMS).  
4 GloFAS couples land surface modeling with global routing networks to produce probabilistic  
5 river discharge forecasts (Alfieri et al., 2013; Harrigan et al., 2023). For our evaluations, we use  
6 the current operational version, GloFAS version 4. Readers are referred to the corresponding  
7 literature and CEMS documentation for detailed descriptions of the model's structure and  
8 calibration.

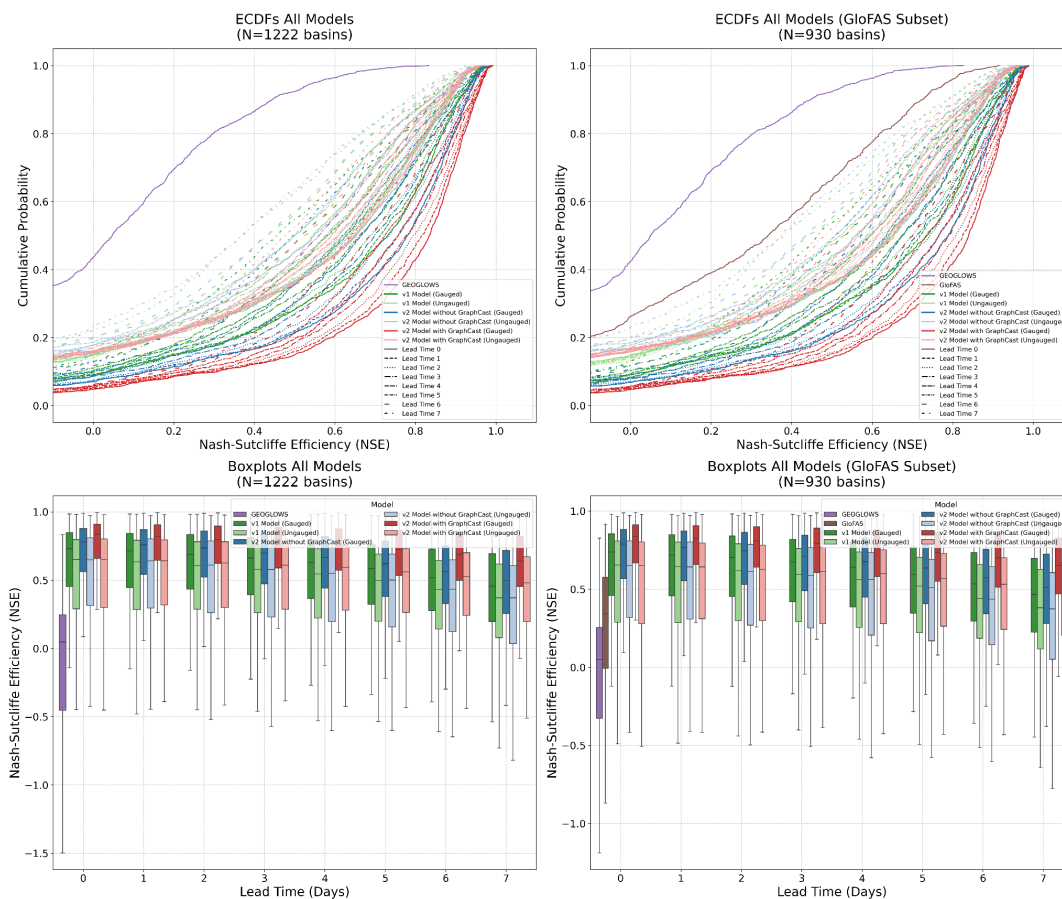
9 The third benchmark is the GEOGLOWS ECMWF Streamflow Routing model. GEOGLOWS  
10 provides global-scale hydrologic forecasting by mapping ECMWF runoff predictions to a  
11 high-resolution global river network using the Routing Application for Parallel computation of  
12 Discharge (RAPID) framework (Souffront et al., 2019). For details regarding the GEOGLOWS  
13 configuration and routing methodology, we refer readers to the associated academic literature  
14 and the model's online documentation.

15 Geolocation matching between streamflow gauges and the GloFAS river network was  
16 performed by the GloFAS team as part of the Nearing et al. (2024) publication. To align gauges  
17 with GEOGLOWS prediction points, we matched catchment polygons to the lumped upstream  
18 area polygons of TDX-Hydro subbasins, the hydrofabric used by GEOGLOWS. For each  
19 polygon pair, we calculated the Intersection over Union (IoU) of their areas and selected the  
20 match with the highest spatial agreement, enforcing a minimum IoU threshold of 0.8. The  
21 bottom row of **Figure 2** shows the locations of these test gauges used for evaluation and  
22 benchmarking.

## 23 **4 Results**

### 24 **4.1 Global Performance Distributions**

25 **Figure 5** shows the NSE distributions of all models at all lead times in both gauged and  
26 ungauged settings. These figures provide a high-level overview and are difficult to interpret, so  
27 we will dissect these results in more detail in following subsections. The main takeaway from  
28 **Figure 5** is that v2 represents an improvement in accuracy compared with v1, and both  
29 systems perform significantly better than the traditional benchmarks of GloFAS and  
30 GeoGloWS.

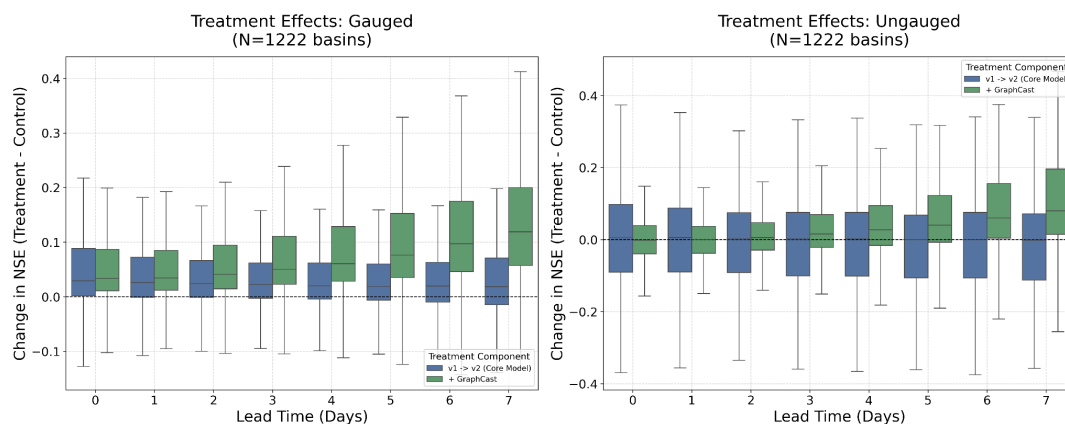


1

2 **Figure 5:** Global performance distributions of Nash-Sutcliffe Efficiency (NSE) across evaluated  
 3 models and lead times. The top row shows Empirical Cumulative Distribution Functions (ECDFs)  
 4 for NSE across all lead times. The bottom row shows the same data as boxplots. The left column  
 5 shows evaluation results for the full test set, and the right column shows evaluation results for  
 6 the portion of the test set with GloFAS data.

## 7 4.2 Disaggregating Improvements

8 **Figure 6** disaggregates the improvements provided by the ME-LSTM architecture and  
 9 expanded training data from the predictive skill injected by the GraphCast meteorological  
 10 forcings. The core v1-to-v2 upgrade (blue boxes) provides improvement in NSE across all lead  
 11 times in the gauged setting, but not in the ungauged setting. The addition of GraphCast (green  
 12 boxes) has a larger overall effect that is smaller at short lead times but grows with a longer  
 13 forecast horizon. This effect exists in both the gauged and ungauged setting, but is larger for  
 14 gauged predictions.

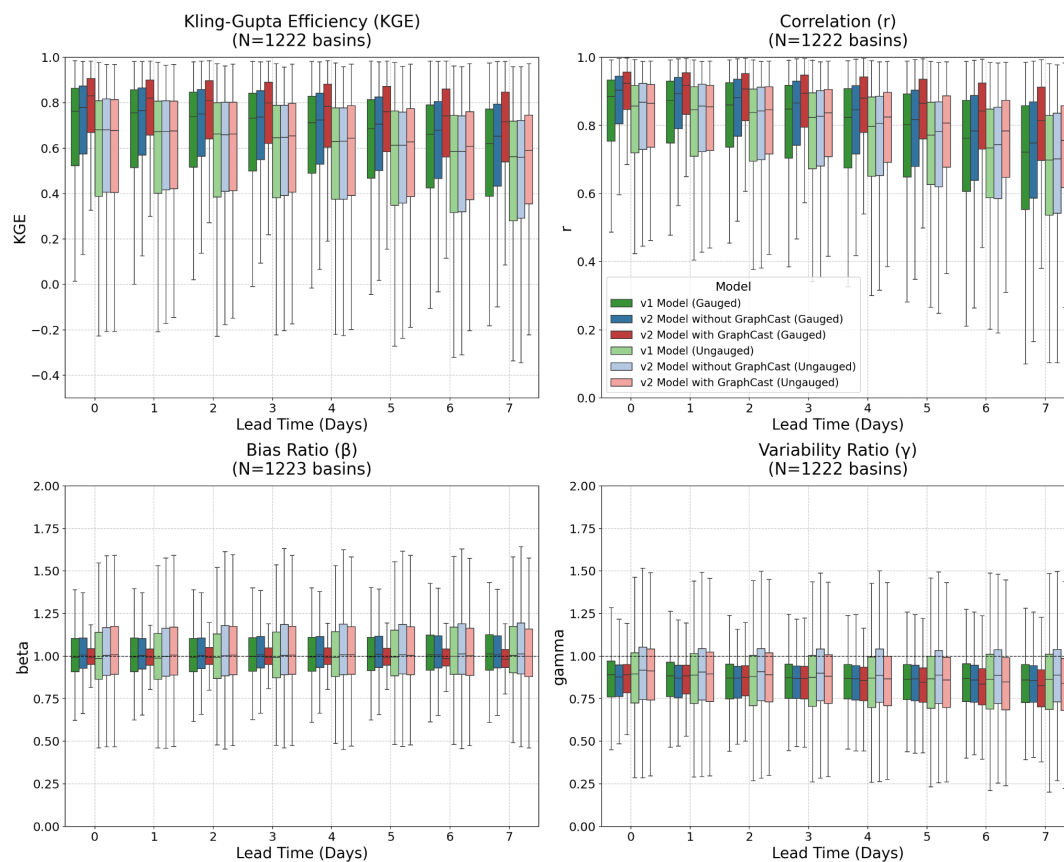


1

2 **Figure 6:** Paired treatment effects across lead times. Blue boxes represent the  $\Delta$ NSE gained by  
 3 transitioning from the v1 to v2 model architecture and expanded Caravan training data. Green  
 4 boxes represent the additional  $\Delta$ NSE gained by incorporating GraphCast.

5 It is instructive to decompose the KGE metric into three constituent parts as described in  
 6 **Section 3.2.** Figure 7 shows this KGE decomposition across lead times. The primary effect is  
 7 that GraphCast forcings improve correlation but lower forecast variance. This is conceptually  
 8 consistent with the fact that AI weather forecasts have higher accuracy than forecasts from  
 9 traditional models, but that they are typically optimized for mean squared error, which  
 10 inherently leads to spatial and temporal smoothing and an underprediction of variance at  
 11 longer lead times (Ben Bouallègue et al., 2024).

12

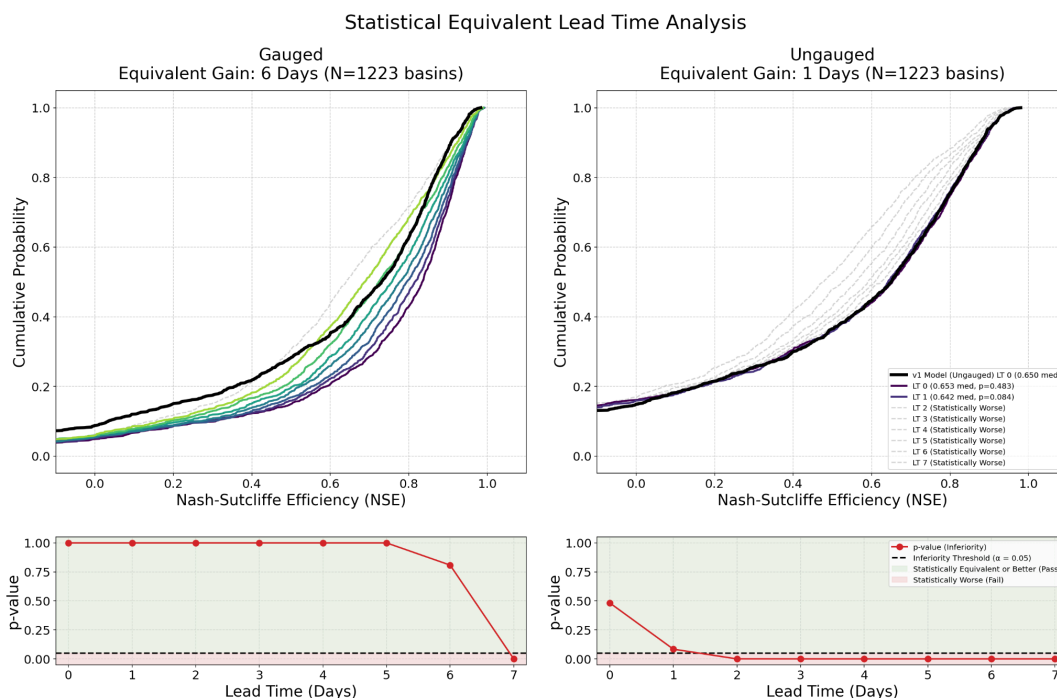


1

2 **Figure 7:** Component Decomposition of KGE (Correlation, Bias Ratio, Variability Ratio) across  
 3 lead times.

#### 4 4.3 Lead Time Extension

5 The previous analyses, and especially **Figure 5**, demonstrate that the v2 system offers an  
 6 extended predictive horizon relative to v1. We used a one-sided Wilcoxon signed-rank test to  
 7 measure how much extra lead time is gained at a statistically significant level ( $\alpha=0.05$ ). **Figure**  
 8 **8** visualizes this test. In the gauged setting, the v2 system demonstrates statistically equivalent  
 9 or superior predictive skill up to a 6-day forecast horizon relative to the v1 benchmark at a  
 10 0-day (nowcast) forecast horizon. In the ungauged setting, this statistical equivalence holds  
 11 out to 1 day of extra lead time.



1  
 2 **Figure 8:** Statistical equivalent lead time evaluation for the v2 system forecasts (lead times 0  
 3 through 7) relative to the v1 benchmark nowcast (lead time 0). The bottom row displays  
 4 p-values derived from a one-sided Wilcoxon signed-rank test comparing the NSE of the v2  
 5 model (with GraphCast) at extending lead times against the v1 model at Day-0. The p-value  
 6 (red line) being within the green area ( $p > \alpha$ ) indicates lead times where the v2 forecast remains  
 7 statistically equivalent to, or better than, the v1 nowcast. The horizontal dashed line represents  
 8 the significance threshold ( $\alpha=0.05$ ). The v2 system achieves up to 6 days of statistically  
 9 non-different lead time in the gauged setting (left panel) and 1 day in the ungauged setting  
 10 (right panel).

#### 11 4.4 Effect of Hydrological Characteristics

12 We computed Pearson correlation coefficients between static catchment attributes and both  
 13 absolute predictive skill (NSE) and the change in skill ( $\Delta$ NSE). **Figure 9** displays the ten attributes  
 14 with the highest correlations (for the v2 model) in gauged and ungauged basins.

15 Absolute predictive skill correlates positively with soil organic carbon, aridity index, and  
 16 moisture indexes, and snow cover fractions across all models. It is important to note that the  
 17 particular aridity index we are using (from HydroSheds) is lower in dryer basins and higher in  
 18 wetter basins. Forecasts are most accurate in wet catchments with more snow and more  
 19 vegetation and less accurate in arid catchments, due to the fact that streamflow in arid basins  
 20 is flashy and determined mostly by the accuracy of precipitation data. Vegetation, soil water  
 21 storage, and snow act as low-pass filters on the rainfall-runoff relationship. Prior work

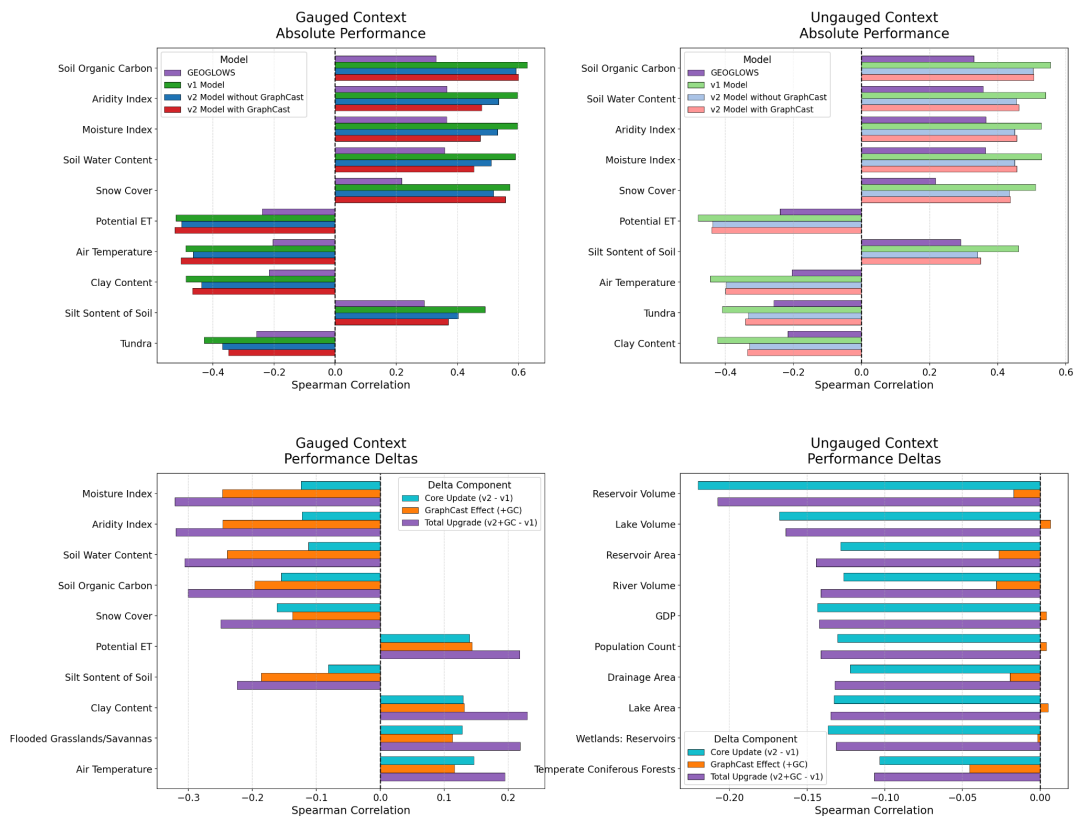


1 demonstrates how LSTM-based models are able to learn seasonal snow accumulation and melt  
 2 dynamics from precipitation and streamflow data (Kratzert et al., 2019c).

3 In gauged catchments, skill differences ( $\Delta$ NSE) between model versions correlates negatively  
 4 with humidity indexes, meaning that *improvements* are larger in more arid locations. In  
 5 ungauged catchments,  $\Delta$ NSE correlates negatively with lake and reservoir area fractions. This  
 6 indicates less improvement in managed or regulated catchments when direct streamflow  
 7 observations are unavailable. Notice that this effect in ungauged basins is not related to the  
 8 introduction of GraphCast forcings, but due to the model itself.

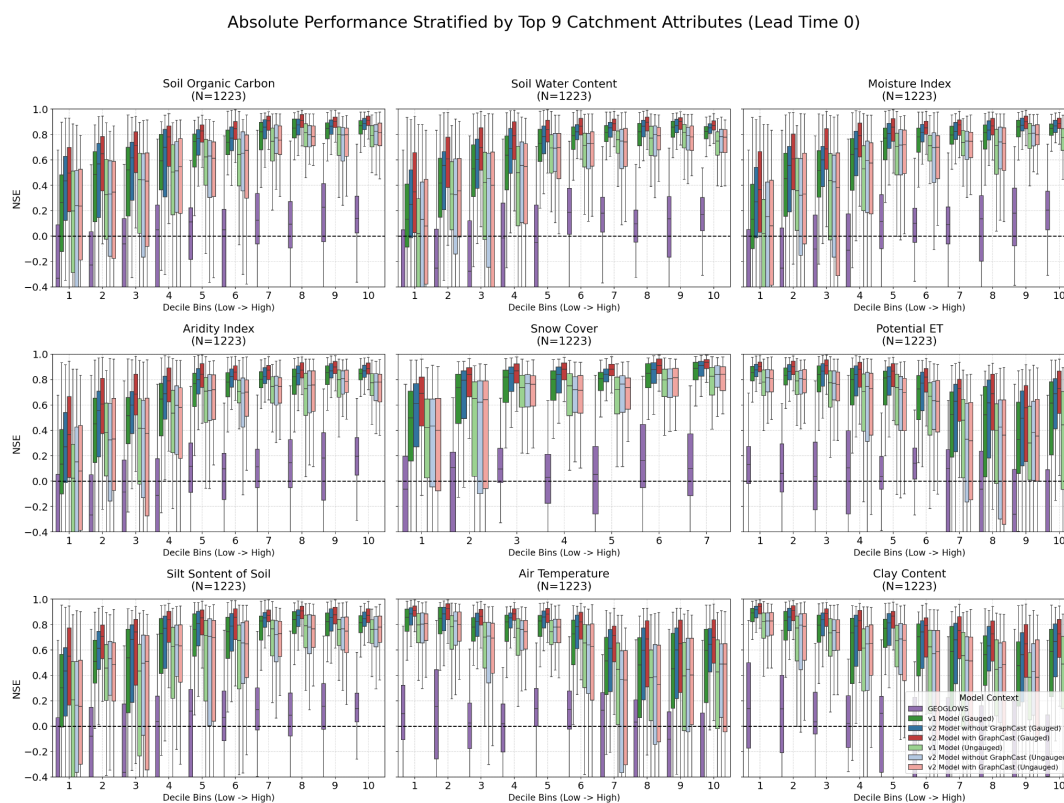
9 **Figure 10** shows the absolute (not delta) performance of each model stratified along deciles of  
 10 the 9 catchment attributes with the highest correlation with both gauged and ungauged  
 11 performance. While trends are generally the same between v1 and v2 models in both gauged  
 12 and ungauged settings, there are certain attribute deciles where the v2 model is not always  
 13 better, however these differences are small. The attribute-to-performance trends also exist for  
 14 GEOGLOWS, but the correlations are weaker (also, see **Figure 9**), however the GEOGLOWS  
 15 scores are lower overall.

Feature Correlations: Top 10 Unique Catchment Attributes Driving Sub-Contexts (Lead Time 0)





1 **Figure 9:** Pearson correlation coefficients between static catchment attributes and model  
 2 performance. The top panels show correlations with absolute NSE in gauged (left) and  
 3 ungauged (right) basins. The bottom panels show correlations with the  $\Delta$ NSE between system  
 4 versions.



6 **Figure 10:** Absolute model performance stratified by key catchment attributes. Distribution of  
 7 Nash-Sutcliffe Efficiency (NSE) at 0-day lead time for the benchmark (GEOGLWS), baseline  
 8 (v1), and updated (v2) model configurations across gauged and ungauged contexts. Basins are  
 9 divided into ten equal-sized bins (deciles) ranging from lowest (1) to highest (10) values for the  
 10 top nine catchment attributes most strongly rank-correlated with performance (**Figure 9**).  
 11 Boxplots indicate the median and interquartile range of NSE within each decile, illustrating how  
 12 specific physical and climatological characteristics impact the predictive skill of the varying  
 13 modeling frameworks.

#### 14 4.5 Gauged vs. Ungauged Performance

15 **Figure 11** shows skill differences ( $\Delta$ NSE) between gauged vs. ungauged basins in the v1  
 16 system, the v2 system without GraphCast, and the v2 system with GraphCast. The bottom  
 17 panel of **Figure 11** shows median skill differences and relative percentage NSE penalties across

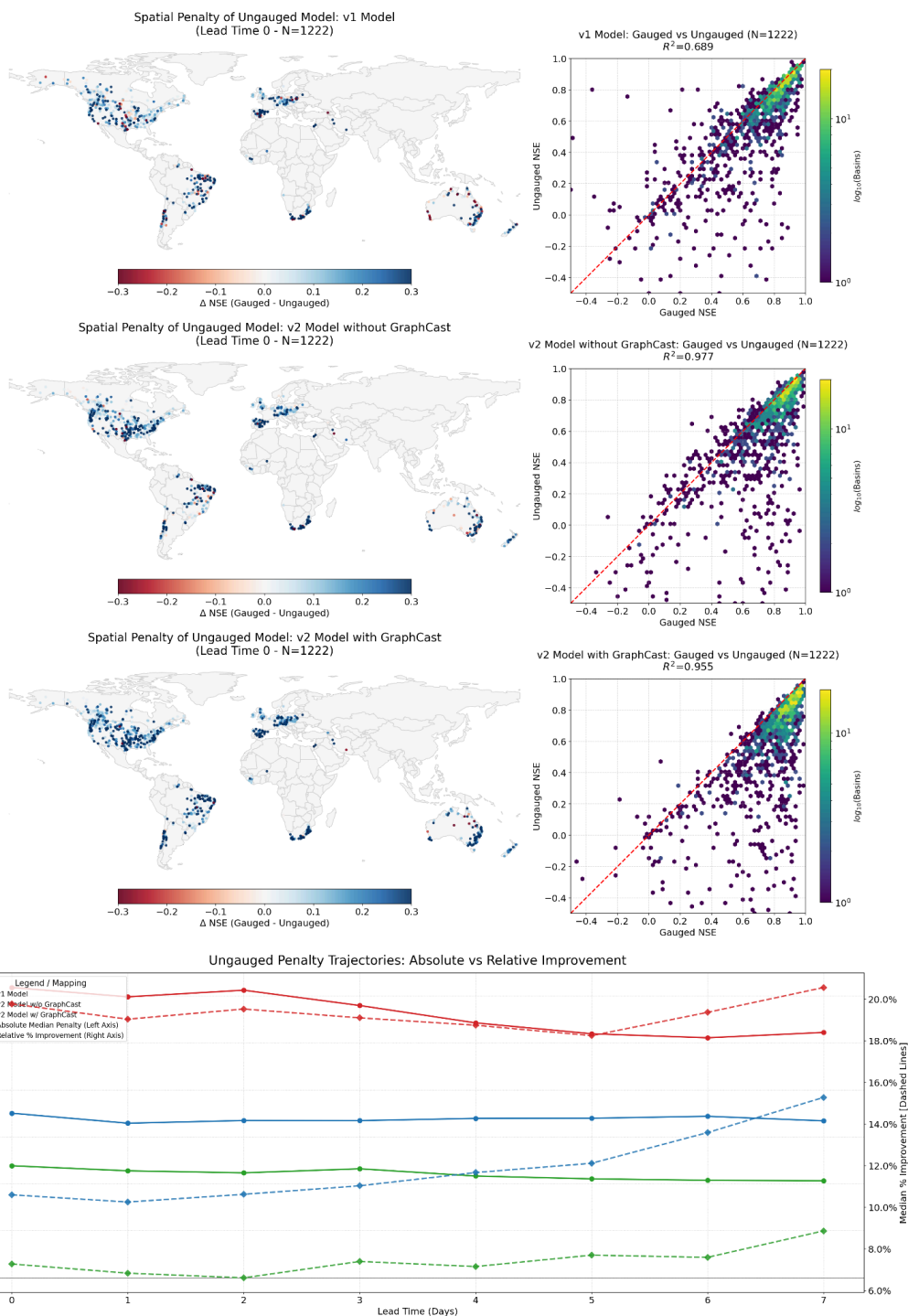


1 lead times.

2 At a 0-day lead time, the v2 system without GraphCast has a median gauged NSE of 0.78, with  
3 an absolute median penalty of 0.07 (10.6% relative decrease). The v2 system with GraphCast  
4 achieves a median gauged NSE of 0.83, with an absolute median penalty of 0.12 (19.8% relative  
5 decrease).

6 The absolute median penalty remains relatively constant over the forecast horizon. At a 7-day  
7 lead time, the v2 system without GraphCast yields a median gauged NSE of 0.4989 with an  
8 absolute penalty of 0.07 (15.3% relative decrease). The v2 system with GraphCast yields a  
9 median gauged NSE of 0.64 with an absolute penalty of 0.10 (20.6% relative decrease).

10 Although the largest skill difference between gauged and ungauged settings in both an  
11 absolute and relative sense is between the v2 systems with and without GraphCast at longer  
12 lead times, it is important to recognize that the v2 system with GraphCast is the most accurate  
13 overall (according to this median NSE metric). The ungauged median NSE at a 7-day lead time is  
14 0.48 for v2 with GraphCast and 0.37 for both v1 and v2 without Graphcast. The second row of  
15 **Figure 5** shows the median scores for all models at all lead times. Results indicate that while the  
16 v2 system yields higher absolute performance globally, it is proportionally more sensitive to the  
17 absence of local streamflow data for training.



1



1 **Figure 11:** Gauged versus ungauged performance penalty measured as  $\Delta NSE$ . The top three  
2 panels display the spatial distribution and scatter of gauged versus ungauged NSE scores for  
3 the v1 system, the v2 system without GraphCast, and the full v2 system, respectively. The  
4 bottom panel illustrates the absolute median and relative percentage NSE penalties across the  
5 forecast horizon.

## 6 5 Conclusions

### 7 5.1 Summary of System Upgrades

8 Version 2 of the Google Global Flood Forecasting system replaces the v1 encoder-decoder  
9 LSTM with a continuous Mean Embedding (ME-LSTM) architecture. The v2 system incorporates  
10 an expanded training dataset from the Caravan community dataset and integrates GraphCast  
11 AI-based meteorological forcings alongside traditional numerical weather prediction products.  
12 These architectural changes resolve certain forecast initialization artifacts observed in the v1  
13 system (**Section 2.2.1**).

14 The v2 system improves aggregate predictive skill as compared to the v1 system and two  
15 third-party benchmarks (GloFAS and GEOGLOWS). The main takeaways from our analysis are  
16 as follows:

- 17 • KGE decomposition shows these improvements are primarily driven by enhanced  
18 temporal correlation rather than reductions in volumetric bias or variance. GraphCast  
19 inputs improve correlation at longer lead times, but reduce forecast variance.
- 20 • The v2 system shows statistically significant lead time extension by 6 days in gauged  
21 basins and 1 day in ungauged basins relative to the v1 nowcast.
- 22 • Performance improvements are correlated with arid environments for gauged basins  
23 and with lake and reservoir indicators in ungauged basins, indicating less improvement  
24 in managed or regulated catchments .
- 25 • The addition of GraphCast appears to be the largest driver of skill increases in  
26 ungauged basins, especially at longer lead times.

### 27 5.2 Global Data Frameworks & Open Science

28 The technical advancements described in this paper and the associated open data and open  
29 source code releases represent an attempt to support the principles of open science in global  
30 hydrology. The future of operational AI hydrology depends on the availability of standardized,  
31 historical data archives and real-time data streams. We would like to highlight and endorse  
32 ongoing efforts by the World Meteorological Organization (WMO) to establish a unified global  
33 data infrastructure. Data collected by National Meteorological and Hydrology Services



1 (NMHSs), and global efforts to collect and standardize hydrology data are critical for the  
2 advancement of large sample hydrology (Addor et al., 2020).

3 Specifically, we highlight three initiatives that warrant the support and contributions of the  
4 hydrological science and operations communities:

- 5 • **WMO Unified Data Policy (Resolution 1, Cg-Ext(2021)):** This resolution commits  
6 Member states to the free and unrestricted exchange of Earth system data. It provides the  
7 political and legal foundation necessary for training robust, less-biased AI models (WMO,  
8 2021).
- 9 • **WMO Information System 2.0 (WIS 2.0):** By transitioning from legacy  
10 message-switching to a modern, pub/sub-based data exchange framework, WIS 2.0  
11 lowers the barrier to entry for NMHSs to share and discover data in real-time (WMO,  
12 2023a).
- 13 • **WMO Integrated Processing and Prediction System (WIPPS):** Formerly the GDPFS,  
14 WIPPS creates a standardized mechanism for cascading high-quality global model  
15 products (such as the ones presented here) down to local services. This ensures that  
16 advancements in global AI forecasting can be systematically integrated into national early  
17 warning systems (WMO, 2023b).

18 By releasing open-access datasets like the Google Runoff Reanalysis & Reforecast (GRRR)  
19 dataset (see Code and Data Availability Section) and Caravan (Kratzert et al., 2023), we aim to  
20 accelerate the collective ability of the scientific community to validate, improve, and deploy  
21 flood forecasting systems at both global and local scales.

## 22 **6 Code and Data Availability**

23 To support transparency, reproducibility, and future research in large-sample hydrology, the  
24 code, data, and model outputs associated with this study are made publicly available.

- 25 • **Codebase:** The open-source GoogleHydrology framework provides the ability to train  
26 and evaluate models that are structurally and functionally identical to the models used in  
27 this study. This is available on GitHub at:  
28 <https://github.com/google-research/flood-forecasting>.
- 29 • **Archived Release:** A persistent, citable archive containing the v2 model outputs and  
30 analysis code used in this paper are published (Nearing et al., 2026) under Zenodo with a  
31 CCA-4 license (<https://doi.org/10.5281/zenodo.19676842>).
- 32 • **Benchmark Data:** Data for the v1 and GloFAS benchmarks are published (Nearing et. al.,  
33 2023b) under Zenodo with a CCA-2 Non-Commercial license  
34 (<https://doi.org/10.5281/zenodo.813937>).



- 1 • **Additional Data Release:** The Google Runoff Reanalysis & Reforecast (GRRR) dataset  
2 contains global river discharge estimates and hydrologic predictions for approximately 1M  
3 locations is publicly accessible via the Google Flood Forecasting resources page at  
4 <https://sites.research.google/gr/floodforecasting/resources/>. This dataset is provided to  
5 support open science by allowing validation with local data almost anywhere in the world,  
6 but was not used in this publication.

## 7 Author Contributions

8 Grey Nearing wrote the paper with help from Martin Gauch, Frederik Kratzert, and  
9 Deborah Cohen.

10

11 Martin Gauch and Grey Nearing ran the v2 model runs used in this paper. Frederik Kratzert  
12 performed the Geogluws benchmarking. Grey Nearing performed the data analysis.

13

14 Guy Shalev, Martin Gauch, Frederik Kratzert, Deborah Cohen, Ben Feinstein,  
15 Asher Metzger, Daniel Klotz, Ido Zemach, Oren Gilon, Rotem Green, Rotem Mayo,  
16 Oleg Zlydenko, Grey Nearing contributed to developing the v2 model.

17

18 Rony Amira, Omri Shefi, Amit Markel, Martin Gauch, Frederik Kratzert, Gila Loike and  
19 Grey Nearing contributed to developing the GoogleHydrology open source codebase.

20

21 Aviel Niego, Hadas Fester, Rom Aschner, Dan Korenfeld, Yuval Shildan,  
22 Benny Mosheyev, Amitay Sicherman, Frederik Kratzert, Martin Gauch, Guy Shalev,  
23 Oren Gilon contributed to operationalizing the v2 model.

24

25 Stephanie Rees, Emily Reinstein, Gila Loike, Yuval Carny, Grey Nearing,  
26 Asher Metzger supported international partnerships that provided feedback to guide model  
27 development.

28

29 Yossi Matias, Avinatan Hassidim, Deborah Cohen, Oren Gilon, and Shmulik Fronman  
30 provided team leadership.

## 31 Acknowledgements

32 During the preparation of this work the authors used Gemini to edit and improve the readability  
33 of the manuscript. After using this service, the authors reviewed and edited the entirety of the  
34 text and take full responsibility for the content of the publication.

## 35 Competing Interests

36 All authors are employed by Google, the organization that developed and operates the Google



- 1 Global Flood Forecasting system and the associated open-source GoogleHydrology codebase
- 2 evaluated in this manuscript. Additionally, the authors may hold stock or share ownership in
- 3 Alphabet Inc., Google's parent company.

## 4 Financial Support

- 5 The authors received no financial support for the research, authorship, and/or publication of this
- 6 article.

## 7 References

- 8 Acuña Espinoza, Eduardo, et al. "Analyzing the generalization capabilities of a hybrid hydrological model for
- 9 extrapolation to extreme events." *Hydrology and Earth System Sciences* 29.5 (2025): 1277-1294.
- 10 Addor, Nans, et al. "The CAMELS data set: catchment attributes and meteorology for large-sample studies."
- 11 *Hydrology and Earth System Sciences* 21.10 (2017): 5293-5313.
- 12 Addor, Nans, et al. "Large-sample hydrology: recent progress, guidelines for new datasets and grand challenges."
- 13 *Hydrological Sciences Journal* 65.5 (2020): 712-725.
- 14 Alfieri, L., Burek, P., Dutra, E., Krzeminski, B., Muraro, D., Thielen, J., & Pappenberger, F. "GloFAS – global ensemble
- 15 streamflow routing and flood early warning." *Hydrology and Earth System Sciences* 17.3 (2013): 1161-1175.
- 16 Allen, Richard G., et al. "Crop evapotranspiration–Guidelines for computing crop water requirements–FAO Irrigation
- 17 and drainage paper 56." *Fao, Rome* 300.9 (1998): D05109.
- 18 Alvarez-Garreton, Camila, et al. "The CAMELS-CL dataset: catchment attributes and meteorology for large sample
- 19 studies–Chile dataset." *Hydrology and Earth System Sciences* 22.11 (2018): 5817-5846.
- 20 Arsenault, Richard, et al. "Continuous streamflow prediction in ungauged basins: Long Short-Term Memory Neural
- 21 Networks clearly outperform hydrological models." *Hydrology and Earth System Sciences Discussions* 2022 (2022):
- 22 1-29.
- 23 Arsenault, Richard, et al. "A comprehensive, multisource database for hydrometeorological modeling of 14,425 North
- 24 American watersheds." *Scientific Data* 7.1 (2020): 243.
- 25 Başağaoğlu, Hakan, et al. "A review on interpretable and explainable artificial intelligence in hydroclimatic
- 26 applications." *Water* 14.8 (2022): 1230.
- 27 Ben Bouallègue, Zied, et al. "The rise of data-driven weather forecasting: A first statistical assessment of machine
- 28 learning–based weather forecasts in an operational-like context." *Bulletin of the American Meteorological Society*
- 29 105.6 (2024): E864-E883.
- 30 Casado Rodriguez, J.: CAMELS-ES: Catchment Attributes and Meteorology for Large-Sample Studies Spain,
- 31 <https://doi.org/10.5281/zenodo.8373020>, 2023.
- 32 Chagas, Vinicius BP, et al. "CAMELS-BR: hydrometeorological time series and landscape attributes for 897
- 33 catchments in Brazil." *Earth System Science Data Discussions* 2020 (2020): 1-41.



- 1 Chaudhary, Priyanka, et al. "Flood uncertainty estimation using deep ensembles." *Water* 14.19 (2022): 2980.
- 2 Chen, Mingyue, et al. "Assessing objective techniques for gauge-based analyses of global daily precipitation."  
3 *Journal of Geophysical Research: Atmospheres* 113.D4 (2008).
- 4 Coxon, Gemma, et al. "CAMELS-GB: hydrometeorological time series and landscape attributes for 671 catchments in  
5 Great Britain." *Earth System Science Data Discussions* 2020 (2020): 1-34.
- 6 Dai, Zhihui, et al. "A hydrological data prediction model based on LSTM with attention mechanism." *Water* 15.4 (2023):  
7 670.
- 8 De la Fuente, Luis Andres, et al. "Towards interpretable LSTM-based modelling of hydrological systems." *Hydrology*  
9 *and Earth System Sciences Discussions* 2023 (2023): 1-36.
- 10 Deng, Chao, et al. "Catchment runoff simulation by coupling data assimilation and machine learning methods."  
11 *Advances in Water Science* 34.6 (2023): 839-849.
- 12 Färber, Claudia, et al. "GRDC-Caravan: extending caravan with data from the global runoff data centre." *Earth System*  
13 *Science Data Discussions* 2024 (2024): 1-17.
- 14 Feng, Dapeng, Kuai Fang, and Chaopeng Shen. "Enhancing streamflow forecast and extracting insights using  
15 long-short term memory networks with data integration at continental scales." *Water Resources Research* 56.9  
16 (2020): e2019WR026793.
- 17 Feng, Dapeng, et al. "Differentiable, learnable, regionalized process-based models with multiphysical outputs can  
18 approach state-of-the-art hydrologic prediction accuracy." *Water Resources Research* 58.10 (2022):  
19 e2022WR032404.
- 20 Feng, Dapeng, et al. "The suitability of differentiable, physics-informed machine learning hydrologic models for  
21 ungauged regions and climate change impact assessment." *Hydrology and Earth System Sciences* 27.12 (2023):  
22 2357-2373.
- 23 Fowler, Keirnan JA, et al. "CAMELS-AUS: hydrometeorological time series and landscape attributes for 222  
24 catchments in Australia." *Earth System Science Data Discussions* 2021 (2021): 1-30.
- 25 Frame, Jonathan M., et al. "On strictly enforced mass conservation constraints for modelling the Rainfall-Runoff  
26 process." *Hydrological Processes* 37.3 (2023): e14847.
- 27 Gauch, Martin, et al. "In defense of metrics: Metrics sufficiently encode typical human preferences regarding  
28 hydrological model performance." *Water Resources Research* 59.6 (2023): e2022WR033918.
- 29 Gauch, Martin, et al. "How to deal with missing input data." *Hydrology and Earth System Sciences* 29.21 (2025):  
30 6221-6235.
- 31 Gupta, Hoshin V., et al. "Decomposition of the mean squared error and NSE performance criteria: Implications for  
32 improving hydrological modelling." *Journal of hydrology* 377.1-2 (2009): 80-91.
- 33 Harrigan, S., Zsoter, E., Alfieri, L., Prudhomme, C., Salamon, P., Wetterhall, F., ... & Pappenberger, F. "GloFAS v4.0: a  
34 global hydrological ensemble forecast system." (2023).
- 35 Helgason, Hordur Bragi, and Bart Nijssen. "LamaH-Ice: LARge-SaMple data for hydrology and environmental sciences



- 1 for Iceland." *Earth System Science Data* 16.6 (2024): 2741-2771.
- 2 Hoedt, Pieter-Jan, et al. "Mc-Istm: Mass-conserving Istm." *International conference on machine learning*. PMLR,  
3 2021.
- 4 Höge, Marvin, et al. "Improving hydrologic models for predictions and process understanding using neural ODES."  
5 *Hydrology and Earth System Sciences Discussions* 2022 (2022): 1-29.
- 6 Höge, Marvin, et al. "CAMELS-CH: hydro-meteorological time series and landscape attributes for 331 catchments in  
7 hydrologic Switzerland." *Earth System Science Data Discussions* 2023 (2023): 1-46.
- 8 Hosseini Hossein Abadi, Farzad, Cristina Prieto Sierra, and César Álvarez Diaz. "An explainable AI approach for  
9 interpreting regionally optimized deep neural networks in hydrological prediction." (2025).
- 10 Huffman, George J., et al. "Integrated multi-satellite retrievals for the global precipitation measurement (GPM)  
11 mission (IMERG)." *Satellite precipitation measurement: Volume 1*. Cham: Springer International Publishing, 2020.  
12 343-353.
- 13 Jia, Xiaowei, et al. "Physics-guided machine learning for scientific discovery: An application in simulating lake  
14 temperature profiles." *ACM/IMS Transactions on Data Science* 2.3 (2021): 1-26.
- 15 Klingler, Christoph, Karsten Schulz, and Mathew Herrnegger. "Lamahl large-sample data for hydrology and  
16 environmental sciences for central europe." *Earth System Science Data Discussions* 2021 (2021): 1-46.
- 17 Klotz, Daniel, et al. "Uncertainty estimation with deep learning for rainfall-runoff modeling." *Hydrology and Earth  
18 System Sciences* 26.6 (2022): 1673-1693.
- 19 Kratzert, Frederik, et al. "NeuralHydrology-interpreting LSTMs in hydrology." *Explainable AI: Interpreting, explaining  
20 and visualizing deep learning*. Cham: Springer International Publishing, 2019. 347-362.
- 21 Kratzert, Frederik, et al. "Toward improved predictions in ungauged basins: Exploiting the power of machine  
22 learning." *Water Resources Research* 55.12 (2019): 11344-11354.
- 23 Kratzert, Frederik, et al. "Towards learning universal, regional, and local hydrological behaviors via machine learning  
24 applied to large-sample datasets." *Hydrology and Earth System Sciences* 23.12 (2019): 5089-5110.
- 25 Kratzert, Frederik, et al. "A note on leveraging synergy in multiple meteorological datasets with deep learning for  
26 rainfall-runoff modeling." *Hydrology and Earth System Sciences Discussions* 2021 (2021): 1-26.
- 27 Kratzert, Frederik, et al. "Caravan-A global community dataset for large-sample hydrology." *Scientific Data* 10.1  
28 (2023): 61.
- 29 Kratzert, Frederik, et al. "HESS Opinions: Never train a Long Short-Term Memory (LSTM) network on a single basin."  
30 *Hydrology and Earth System Sciences* 28.17 (2024): 4187-4201.
- 31 Krejčí, J. & Nearing, G. (2025). CAMELS-CZ: A catchment attributes and meteorology dataset for large-sample  
32 studies in Czechia [Dataset]. Zenodo. <https://www.google.com/search?q=https://doi.org/10.5281/zenodo.17593968>.
- 33 Lam, Remi, et al. "Learning skillful medium-range global weather forecasting." *Science* 382.6677 (2023): 1416-1421.
- 34 Linke, Simon, et al. "Global hydro-environmental sub-basin and river reach characteristics at high spatial resolution."



- 1 Scientific data 6.1 (2019): 283.
  
- 2 Liu, Jinping. "Progresses and Challenges on QPE/QPF Utilization in Hydrology." *Tropical Cyclone Research and Review* 1.2 (2012): 194-197.
  
- 4 Liu, Siyan, et al. "Uncertainty quantification of machine learning models to improve streamflow prediction under changing climate and environmental conditions." *Frontiers in Water* 5 (2023): 1150126.
  
- 6 Liu, Jun, et al. "CAMELS-DK: Hydrometeorological Time Series and Landscape Attributes for 3330 Catchments in Denmark." *Earth System Science Data Discussions* 2024 (2024): 1-30.
  
- 8 Lees, Thomas. *Deep learning for hydrological modelling: from benchmarking to concept formation*. Diss. University of Oxford, 2022.
  
- 10 Loritz, Ralf, et al. "CAMELS-DE: hydro-meteorological time series and attributes for 1555 catchments in Germany." *Earth System Science Data Discussions* 2024 (2024): 1-30.
  
- 12 Mai, Juliane, et al. "The great lakes runoff intercomparison project phase 4: the great lakes (GRIP-GL)." *Hydrology and Earth System Sciences Discussions* 2022 (2022): 1-54.
  
- 14 Morin, E. "Caravan extension Israel-Israel dataset for large-sample hydrology, Zenodo [data set]." 2023,
  
- 15 Muñoz-Sabater, Joaquin, et al. "ERA5-Land: A state-of-the-art global reanalysis dataset for land applications." *Earth system science data* 13.9 (2021): 4349-4383.
  
- 17 Nash, J. Eamonn, and Jonh V. Sutcliffe. "River flow forecasting through conceptual models part I-A discussion of principles." *Journal of hydrology* 10.3 (1970): 282-290.
  
- 19 Nearing, Grey, et al. "Data assimilation and autoregression for using near-real-time streamflow observations in long short-term memory networks." *Hydrology and earth system sciences discussions* 2021 (2021): 1-25.
  
- 21 Nearing, G., et al.: *From Hindcast to Forecast with Deep Learning Streamflow Models*, EGU General Assembly 2023, Vienna, Austria, 24–28 Apr 2023, EGU23-10020, <https://doi.org/10.5194/egusphere-egu23-10020>, 2023a.
  
- 23 Nearing, G.: *Global prediction of extreme floods in ungauged watersheds*, Zenodo [data set],  
24 <https://doi.org/10.5281/zenodo.8139379>, 2023b.
  
- 25 Nearing, Grey, et al. "Global prediction of extreme floods in ungauged watersheds." *Nature* 627.8004 (2024):  
26 559-563.
  
- 27 Nearing, G., Kratzert, F., and Gauch, M.: *Extending Medium-Range Global Flood Forecasts: The Google Global Flood Forecasting Model Version 2*, Zenodo [data set], <https://doi.org/10.5281/zenodo.19676842>, 2026.
  
- 29 Newman, Andrew J., et al. "Development of a large-sample watershed-scale hydrometeorological data set for the contiguous USA: data set characteristics and assessment of regional variability in hydrologic model performance." *Hydrology and Earth System Sciences* 19.1 (2015): 209-223.
  
- 32 Nevo, Sella, et al. "Flood forecasting with machine learning models in an operational framework." *Hydrology and Earth System Sciences* 26.15 (2022): 4013-4032.
  
- 34 Nourani, Vahid, Kasra Khodkar, and Mekonnen Gebremichael. "Uncertainty assessment of LSTM based groundwater



- 1 level predictions." *Hydrological Sciences Journal* 67.5 (2022): 773-790.
- 2 Pearson, Karl. "VII. Note on regression and inheritance in the case of two parents." *proceedings of the royal society of*  
3 *London* 58.347-352 (1895): 240-242.
- 4 Prieto, Cristina, et al. "Flow prediction in ungauged catchments using probabilistic random forests regionalization  
5 and new statistical adequacy tests." *Water Resources Research* 55.5 (2019): 4364-4392.
- 6 Rosati, Michael, et al. "Decoding LSTM Memory to Reveal Baseflow Contributions in Fractured and Sedimentary  
7 Mountain Basins: A Case Study in the Sangre de Cristo Mountains, Southwestern United States." (2025).
- 8 Sabzipour, Behmard, et al. "Comparing a long short-term memory (LSTM) neural network with a physically-based  
9 hydrological model for streamflow forecasting over a Canadian catchment." *Journal of Hydrology* 627 (2023):  
10 130380.
- 11 Shalev, Guy, and Frederik Kratzert. "Caravan MultiMet: Extending Caravan with Multiple Weather Nowcasts and  
12 Forecasts." *arXiv preprint arXiv:2411.09459* (2024).
- 13 Souffront Alcantara, M. A., Nelson, E. J., Shakya, K., Edwards, C., Roberts, W., Krewson, C., ... & Jones, N. "Hydrologic  
14 modeling as a service (HMaaS): a new approach to address hydroinformatic challenges in developing countries."  
15 *Frontiers in Environmental Science* 7 (2019): 158.
- 16 Tsai, Wen-Ping, et al. "From calibration to parameter learning: Harnessing the scaling effects of big data in  
17 geoscientific modeling." *Nature communications* 12.1 (2021): 5988.
- 18 WMO (2021). *WMO Unified Policy for the International Exchange of Earth System Data. Resolution 1 (Cg-Ext(2021)).*  
19 *World Meteorological Organization.*
- 20 WMO (2023a). *Manual on the WMO Information System (WIS). WMO-No. 1060. World Meteorological Organization.*
- 21 WMO (2023b). *Manual on the WMO Integrated Processing and Prediction System (WIPPS). WMO-No. 485. World*  
22 *Meteorological Organization.*