



1 **Enhancing flood forecasting reliability in data-scarce regions with a distributed**  
2 **hydrology-guided neural network framework**

3 Confidence Duku<sup>1,2</sup>

4 <sup>1</sup>Wageningen Environmental Research, Team Climate Resilience, Wageningen University & Research,  
5 Wageningen, The Netherlands

6 <sup>2</sup>Earth Systems and Global Change Group, Wageningen University & Research, Wageningen, The  
7 Netherlands

8 *Correspondence to:* Confidence Duku([confidence.duku@wur.nl](mailto:confidence.duku@wur.nl))

9

10 **Abstract**

11 Flood early warning systems are critical for reducing disaster impacts, yet their effectiveness remains  
12 limited in data-scarce regions such as Africa and South America. Existing global platforms—including  
13 GloFAS and the Google Flood Hub—exhibit low reliability in these areas, particularly for rare flood events  
14 and under strict timing constraints. Here, I demonstrate the potential of a distributed, hydrology-guided  
15 neural network framework, Bakaano-Hydro, to enhance flood forecasting reliability in data-scarce regions.  
16 The proposed framework integrates process-based runoff generation, topographic routing, and a Temporal  
17 Convolutional Network for streamflow simulation. Using a hindcast-based evaluation across 470 gauging  
18 stations from 1982 to 2016, I benchmark Bakaano-Hydro’s flood detection skill against GloFAS and Google  
19 AI model across multiple return periods (1-, 2-, 5-, and 10-year) and timing tolerances (0–2 days). Results  
20 show that Bakaano-Hydro consistently achieves higher Critical Success Index (CSI), lower False Alarm Rate  
21 (FAR), and higher Probability of Detection (POD), even under exact-day (0-day) timing constraints. Its  
22 median CSI scores at 0-day tolerance exceed or match those of GloFAS and Google AI model under more  
23 lenient timing thresholds. These performance gains are statistically significant across diverse hydroclimatic  
24 regions, including arid and tropical basins, demonstrating the model’s spatial generalization capacity. By  
25 coupling physical realism with machine learning generalizability, Bakaano-Hydro provides a reliable,  
26 interpretable, and open-source tool for enhancing flood forecasting in regions most vulnerable to climate  
27 extremes and least equipped with observational infrastructure.

28 **1. Introduction**

29 Floods are the most frequent and widespread natural disaster globally, accounting for more than 40% of  
30 all weather-related hazards and affecting over 1.5 billion people between 2000 and 2020 alone (Yaghmaei,  
31 2020). In 2021, floods caused an estimated \$105 billion USD in economic losses worldwide, with Africa,



1 South-East Asia and Oceania bearing a disproportionate share of fatalities and damage due to weaker  
2 infrastructure and limited adaptive capacity (AON, 2023). The burden is particularly severe in Africa, where  
3 more than 75 million people live in high flood-risk zones and simultaneously face extreme poverty  
4 (Rentschler et al., 2022). Recent catastrophic events in Nigeria, Sudan, South Africa, and Mozambique  
5 underscore how recurrent flooding exacerbates food insecurity, damages livelihoods, and displaces  
6 communities—often with cascading socio-economic consequences.

7 As climate change intensifies extreme precipitation events, and urbanization continues to encroach upon  
8 floodplains (Mazzoleni et al., 2022; Tellman et al., 2021), the need for accurate and timely early warning  
9 systems becomes increasingly urgent. For example, the last twenty years have seen the number of major  
10 floods more than double (Yaghmaei, 2020). Flood early warning systems can reduce mortality by up to 43%  
11 (WMO, 2013) and cut economic losses by as much as 50% (Pilon, 2002; Rogers & Tsirkunov, 2011), yet  
12 these benefits are unevenly distributed. Enhancing the performance of early warning systems in  
13 developing countries to levels comparable with high-income settings could prevent an estimated 23,000  
14 deaths annually (Hallegatte, 2012). However, many of the most flood-prone and affected regions lack  
15 robust national flood early warning systems. This limits governments' ability to anticipate and respond to  
16 disasters in a timely manner. As a result, disaster response agencies—including national emergency  
17 authorities, humanitarian organizations, and international development partners—often depend on global  
18 flood forecasting platforms such as the Global Flood Awareness System (GloFAS) (Alfieri et al., 2013; Alfieri  
19 et al., 2020; Harrigan et al., 2023), part of the Copernicus Emergency Management Service and the Google  
20 Flood Hub (Nearing et al., 2024; Nevo et al., 2022). These systems offer publicly accessible, transboundary  
21 early warnings and have become essential tools for operational response in data-scarce regions. GloFAS,  
22 couples weather forecasts with the LISFLOOD model (De Roo et al., 2000), a physically based, fully  
23 distributed hydrological model in generating flood forecasts. By contrast, the Google AI model incorporates  
24 weather forecasts into a global data-driven hydrological model (Kratzert et al., 2019), which simulates river  
25 flows and potential flood events.

26 Despite recent advances, the skills of these models in data-scarce regions, are limited and consequently  
27 forecasts remain unreliable. An evaluation of GloFAS and the Google AI model performance across flood  
28 events with varying return periods (1-, 2-, 5-, and 10-year) from 1984 to 2022 revealed notably lower  
29 median performance in Africa and South America compared with Europe and North America (Nearing et  
30 al., 2024). The skill gap is especially pronounced for rarer flood events, where the models' reliability  
31 deteriorates further. The persistent low reliability is particularly concerning given the heightened



1 vulnerability of communities in these regions, where limited infrastructure, higher poverty rates, and lower  
2 adaptive capacity amplify the social and economic consequences of flooding (Fox et al., 2024; Sauer et al.,  
3 2024). A key factor underlying poor model performance is the sparse, fragmented observational networks  
4 in these regions, compounded by issues of data quality, restricted accessibility, and infrastructural  
5 limitations (Grimes et al., 2022). The process-based modelling approach employed by GloFAS typically  
6 requires extensive calibration to individual basins and struggles to generalize to ungauged areas. While  
7 Google AI model’s data-driven approach can potentially learn representations across basins, its lumped  
8 modeling framework—relying on area-weighted averages of climatic forcings, vegetation, and soil  
9 attributes—cannot fully capture localized hydrological processes.

10 Several methodological advances have been proposed to address the limitations of physically based  
11 hydrological models in data-scarce regions. Regionalization techniques, which transfer calibrated  
12 parameters from data-rich to poorly monitored catchments, have been widely adopted, but their  
13 performance remains limited in climatically heterogeneous or poorly instrumented areas (Beck et al., 2016;  
14 Guo et al., 2021; Pagliero et al., 2019). Bias correction of meteorological forcings and streamflow, has also  
15 been proposed to improve the accuracy of forecasts (Crochemore et al., 2016; Tanguy et al., 2025).  
16 However, bias correction requires long, high-quality historical observations to calibrate correction  
17 functions, limiting their applicability in regions with sparse or unreliable in-situ data. Similarly, a promising  
18 direction is data assimilation, particularly the integration of satellite-derived hydrological variables—such  
19 as soil moisture, snow cover, and river levels—into forecasting models to correct state variables and  
20 improve flood prediction in regions with limited in-situ observations. (Alfieri et al., 2022; Emerton et al.,  
21 2016; Wongchuig et al., 2024). However, in data-limited regions, the effectiveness of data assimilation is  
22 often constrained by the coarse resolution and uncertainty of satellite products, scale mismatches, and the  
23 need for careful tuning and error characterization in the absence of reliable ground-based data. In parallel,  
24 data-driven models, particularly deep learning architectures such as Long Short-Term Memory (LSTM)  
25 networks, have demonstrated strong performance in capturing complex, nonlinear rainfall–runoff  
26 relationships without requiring explicit process parameterization (Feng et al., 2020; Hunt et al., 2022;  
27 Kratzert et al., 2018). These models have contributed to advances in streamflow prediction and, by  
28 extension, flood forecasting skill. However, they are often developed using lumped or basin-aggregated  
29 inputs, lack spatially explicit representation of hydrological processes, and performance tends to degrade  
30 with increasing basin size (Hunt et al., 2022). Notably, incorporating physical constraints into data-driven  
31 models has been shown to improve predictive accuracy and enhance robustness across diverse  
32 hydrological settings (Kratzert et al., 2019). Additionally, accounting for spatial variability in inputs,



1 particularly rainfall, can further improve the performance of lumped data-driven models (Wang & Karimi,  
2 2022), highlighting the value of physically informed machine learning approaches that combine data-driven  
3 flexibility with process-based realism.

4 To address the limitations of current flood forecasting systems in data-scarce regions, here, I demonstrate  
5 the potential of a distributed, hydrology-guided neural network modelling framework to significantly  
6 enhance forecasting reliability across Africa and South America. The modelling framework uniquely  
7 integrates physically based hydrological principles with the generalization capacity of machine learning in  
8 a spatially explicit and physically meaningful way. The framework referred to as Bakaano-Hydro is  
9 benchmarked against GloFAS and Google AI model. The name Bakaano comes from Fante, a language  
10 spoken along the southern coast of Ghana. Loosely translated as "by the river side" or "stream-side", it  
11 reflects the lived reality of many vulnerable riverine communities across the Global South - those most  
12 exposed to flood risk and often least equipped to adapt. In this study, retrospective simulations (hindcasts)  
13 are used to evaluate and compare Bakaano-Hydro performance against GloFAS and Google AI model. While  
14 not true forecasts, these hindcasts serve as a standard proxy for assessing forecasting reliability, as  
15 commonly done in forecast evaluation and earlier research (e.g. Alfieri et al., 2013; Nearing et al., 2024).  
16 Accordingly, I use the term 'forecasting reliability' to refer to the skill of models in reproducing observed  
17 flood events across return periods and timing tolerances using historical input data. The Bakaano-Hydro  
18 workflow consists of three stages: (1) the estimation of spatially distributed runoff using a land surface  
19 model designed to capture spatiotemporal variability; (2) routing of this runoff using a topographic flow  
20 direction algorithm; and (3) prediction of daily streamflow using a temporal neural network trained on  
21 routed runoff. Observed streamflow data from 643 gauging stations across Africa and South America are  
22 used for training and evaluation. A single model is employed and trained jointly across Africa and South  
23 America to enhance spatial generalization. Flood return periods (1-, 2-, 5-, and 10-year events) are  
24 computed from observed streamflow at each station and consistently applied across Bakaano-Hydro,  
25 GloFAS, and Google AI model outputs. Critical Success Index (CSI), False Alarm Ratio (FAR), and Probability  
26 of Detection (POD) metrics are calculated under multiple flood timing tolerances (0, 1, and 2 days) to  
27 evaluate both detection accuracy and temporal precision.

28

29

30



## 1 **2. Methods**

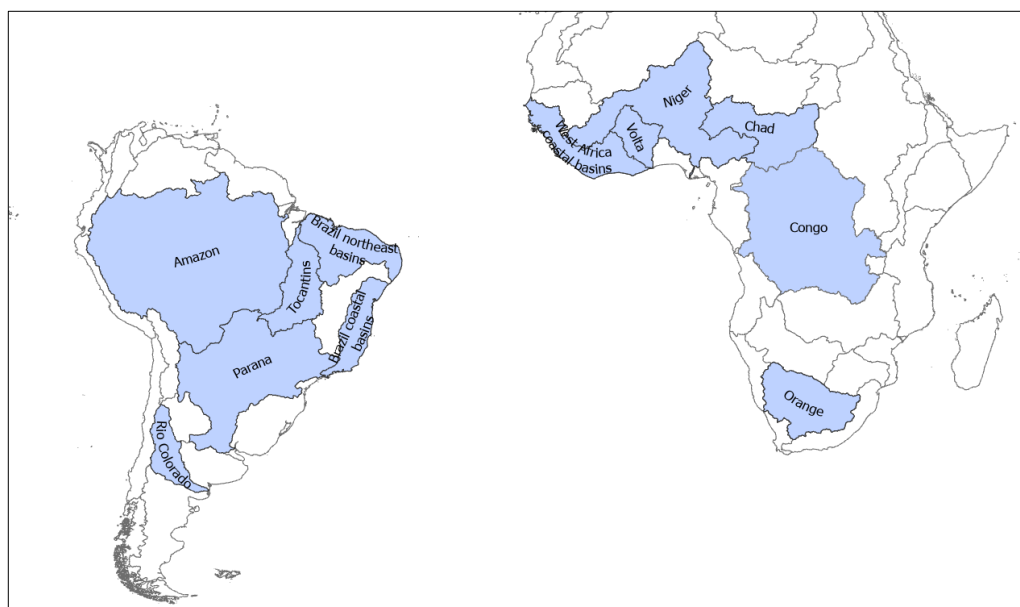
### 2 **2.1 The Bakaano-Hydro modelling framework and data**

3 The selection of river basins in this study (Figure.1) was guided by the need to evaluate flood forecasting  
4 skill across a hydroclimatically and geographically diverse set of regions, representative of key challenges  
5 in data-scarce environments. The basins span a wide range of climatic zones—from humid tropics (e.g.,  
6 Amazon, Congo, Niger) to semi-arid and arid systems (e.g., Volta, Orange, Brazil Northeast, and Rio  
7 Colorado)—and vary substantially in terms of flood regime, seasonality, land use, and topographic  
8 complexity. Several basins (e.g., Amazon, Congo, and Niger) are among the largest in the world and are  
9 subject to recurrent fluvial flooding, while others (e.g., West Africa coastal basins, Brazilian coastal and  
10 interior basins) represent smaller, flashier catchments prone to localized flood events. Importantly, all  
11 selected basins are located in regions with limited in-situ hydrometeorological infrastructure, where global  
12 forecasting systems often perform poorly due to sparse calibration data and limited representation of local  
13 hydrological processes. This selection enables a robust, spatially explicit evaluation of model generalization  
14 across diverse flood-generating mechanisms, data conditions, and socio-environmental contexts in the  
15 Global South.

16 In this implementation of the Bakaano-Hydro framework, the VegET method was first employed to  
17 estimate total runoff. VegET, primarily designed for actual evapotranspiration estimation, also incorporates  
18 key hydrological processes such as rainfall interception, runoff generation, and soil moisture storage (Senay,  
19 2008; Senay et al., 2023). The method operates on a daily time-step, where actual evapotranspiration is  
20 estimated as a function of the Normalized Difference Vegetation Index (NDVI) and reference  
21 evapotranspiration. Rainfall interception is estimated as a function of tree cover, herbaceous cover, and  
22 bare soil fractions, while runoff is computed using a saturation-excess approach, whereby excess soil water  
23 beyond the field capacity is considered unavailable for plant uptake in the root zone. Daily meteorological  
24 inputs, including precipitation and temperature (minimum, maximum, and mean), were obtained from the  
25 CHELSA-W5E5 database a bias-adjusted high resolution global climate dataset derived from CHELSA  
26 (Climatologies at High Resolution for the Earth's Land Surface Areas) and adjusted using W5E5 reference  
27 data (Karger et al., 2023; Karger, 2021). Data availability was limited to the period 1981–2016. NDVI and  
28 fractional tree and herbaceous cover data spanning 2001–2016 were obtained from Moderate Resolution  
29 Imaging Spectroradiometer (MODIS) remote sensing products (Didan, 2021; DiMiceli et al., 2015). MODIS-  
30 NDVI products are unavailable for periods before 2001. Following the VegET procedure, a daily mean



1 climatology of NDVI was established with linear interpolation. Tree cover and herbaceous cover fractions  
2 were also obtained from MODIS (DiMiceli et al., 2015) for the period 2001 – 2016. MODIS tree cover and  
3 herbaceous cover products are unavailable for periods before 2001. Annual mean tree cover and  
4 herbaceous cover fractions were then computed over this time period and were used in runoff estimation.  
5 Soil properties, including wilting point, field capacity, and saturation, were obtained from SoilGrids  
6 database, which is a global soil information system that provides predictions of standard soil properties at  
7 multiple depths and at high resolution (Poggio et al., 2021). Digital Elevation Model (DEM) was sourced  
8 from HydroSHEDS database (Lehner et al., 2008), which provides hydrologically corrected topographic  
9 data. All input datasets were resampled to 1km<sup>2</sup> resolution, matching the spatial resolution of climate data.  
10 Total runoff was computed for the 1981–2016 period, consistent with the availability of climate input data.  
11 Potential evapotranspiration as part of VegET was estimated using the Hargreaves equation (Hargreaves &  
12 Samani, 1985). Further details on the VegET model and its parameterizations are available in Senay et al.  
13 (2023).



14

15 **Figure 1.** Basins for which observed streamflow data were used in training, evaluating and computing  
16 reliability metrics of Bakaano-Hydro. Basin delineation were obtained from HydroSHEDS (Lehner et al.,  
17 2008). Brazil northeast basins, Brazil coastal basins and West Africa coastal basins are not recognized  
18 names but refer to group of smaller basins in the named locations that have been aggregated for the  
19 purpose of this study.

20



1 The second phase of Bakaano-Hydro is the flow routing phase which bridges empirical or physically based  
2 surface runoff generation with data-driven streamflow prediction. In this analysis, daily runoff was routed  
3 through the river channel network using the weighted flow accumulation method based on the multiple  
4 flow direction approach (Quinn et al., 1991). This routing scheme distributes flow from each cell to up to  
5 eight neighboring cells, with partitioning proportional to the elevation gradient between cells. For each of  
6 the 470 gauging stations, routed runoff time-series were extracted from the river network to facilitate  
7 further analysis. Routed runoff time-series are extracted based on the longitude and latitude of the gauging  
8 station. To allow for distortions in geometry and provide tolerance, the coordinates are snapped to the  
9 river network.

10 The final phase of Bakaano-Hydro involves the application of a deep learning model to simulate streamflow.  
11 To capture the complex, nonlinear relationships between routed runoff and streamflow a Temporal  
12 Convolutional Network (TCN) was employed. TCNs are a class of deep learning models designed for  
13 sequence modeling tasks and are based on one-dimensional causal convolutions with dilated filters,  
14 allowing the network to efficiently capture long-range temporal dependencies while maintaining a fixed  
15 input length and avoiding information leakage from future timesteps (Bai et al., 2018). Compared to LSTM  
16 networks, which rely on recurrent structures and hidden states passed sequentially through time, TCNs  
17 use convolutional layers that process entire sequences in parallel. This leads to faster training times, greater  
18 stability during optimization, and better scalability to long input sequences, especially in large  
19 spatiotemporal datasets. Additionally, TCNs avoid issues commonly associated with recurrent networks,  
20 such as vanishing gradients and limited interpretability of memory cells. In the context of Bakaano-Hydro,  
21 TCNs are particularly advantageous because they can efficiently model multi-scale temporal patterns in the  
22 routed runoff time series across hundreds of locations, while enabling the use of a shared model  
23 architecture for diverse river basins. Their ability to learn temporal features over fixed receptive fields,  
24 combined with their parallelizability and robustness, makes TCNs a suitable and scalable choice for  
25 simulating streamflow from spatially distributed runoff inputs in data-limited, flood-prone regions. In  
26 Bakaano-Hydro, the TCN architecture is configured to utilize a hindcast sequence of 365 days to predict  
27 daily streamflow. The architecture consists of two input branches, both employing TCN layers to process  
28 dynamic and static variables. The first input branch processes nine variables—three dynamic and six  
29 static—standardized across all 643 stations in Africa and South America using z-score normalization.  
30 Dynamic variables include routed runoff extracted at gauging stations, augmented by dividing by the  
31 upstream contributing area (in grid cells) and the depth-to-water index. Static variables include the  
32 upstream contributing area, depth-to-water index, and sine-cosine transformations of station latitude and



1 longitude to encode spatial periodicity. The static variables were repeated across the entire duration of the  
2 dynamic variables. The second branch processes station-specific routed runoff data, scaled by the  
3 maximum daily routed runoff across all stations within the basin and log-transformed to reduce skewness.  
4 TCN layers employ dilation rates of 2, 4, 8, 16, 32, 64, 128, and 256, capturing temporal dependencies  
5 across multiple timescales. The outputs from the two TCN branches are concatenated and further  
6 processed through dense layers to establish full connectivity, integrating global and station-specific runoff-  
7 streamflow relationships. Observed streamflow data were sourced from the Global Runoff Data Centre  
8 (GRDC) (GRDC, 2025) and included gauging stations across Africa and South America with at least five years  
9 of data. The training period covered 1989–2016, while model validation was conducted for the 1982–1988  
10 period. Only hydrological stations with observed streamflow record of at least five years were used in the  
11 training of the model.

## 12 **2.2 Benchmark data from GloFAS and Google AI model**

13 The Google AI model builds on previous LSTM-based nowcasting approaches by employing an encoder-  
14 decoder LSTM architecture and area-weighted averages over basin polygons over the total upstream area  
15 of each gauge or prediction point. The model was trained using data from 5680 stations across the globe  
16 (Nearing et al., 2024). The Google AI model data used here were derived from a full model run  
17 encompassing all stations, rather than the cross-validation splits reported in Nearing et al. (2024). The  
18 streamflow predictions are right-labeled, meaning the predicted value at day  $t$  corresponds to the  
19 observation at day  $t - 1$ . To ensure consistency, we relabeled their predictions to match the observation  
20 timestamps used in our dataset.

21 GloFAS data are from GloFAS version 4, which is the latest version as at the time of submission. As part of  
22 GloFAS, the LISFLOOD OS model, implemented at a 0.05-degree quasi-global resolution was calibrated  
23 using in-situ discharge data from 1,996 stations with drainage areas of at least 500 km<sup>2</sup> and observation  
24 records post-1980. These stations were calibrated using the Distributed Evolutionary Algorithm for Python  
25 (Fortin et al., 2012). Parameter values for ungauged catchments estimated through regionalization  
26 (Grimaldi, 2024).

## 27 **2.3 Estimation of flow thresholds for multiple return periods**

28 A common subset of 470 hydrological stations was identified across Bakaano-Hydro, GloFAS, and Google  
29 AI model to enable consistent model evaluation. Return periods for these stations were estimated using  
30 observed streamflow records from the GRDC database. The estimation followed a modular approach based





1 on the U.S. Geological Survey (USGS) Bulletin 17C guidelines (England Jr et al., 2018), incorporating  
2 procedures for annual peak extraction, low outlier treatment, and frequency distribution fitting. Flow  
3 thresholds corresponding to the 1-, 2-, 5-, and 10-year return periods were computed and applied  
4 uniformly across model outputs to evaluate flood detection skill.. The methodology comprises three main  
5 components: (1) extraction of peak flows, (2) identification and treatment of potentially influential low  
6 floods (PILFs), and (3) return period estimation using either parametric or empirical approaches. Peak flows  
7 were extracted from continuous daily discharge records using the annual maximum series (AMS) method.  
8 This method identifies the single highest discharge in each water year and is the standard recommended  
9 by England Jr et al. (2018). The water year is defined such that it ends in September, aligning with typical  
10 hydrological year definitions in the Northern Hemisphere. A minimum data completeness threshold of 50%  
11 of the expected daily observations per year was imposed to ensure robustness. To mitigate the influence  
12 of potentially impactful low floods (PILFs), we employed the Grubbs-Beck Test (GBT) as described in  
13 England Jr et al. (2018). This test detects low outliers by comparing observations against a lower threshold  
14 derived from sample mean and standard deviation, using critical values from a tabulated KN statistic. Data  
15 points below the threshold were flagged as PILFs and subsequently censored using a left-censoring  
16 approach during distribution fitting. The PILF threshold was defined as the midpoint between the largest  
17 outlier and the smallest non-outlier. For stations with insufficient sample size to run the GBT (i.e., fewer  
18 than 10 observations), no censoring was applied.

19 The primary method employed was the Generalized Expected Moments Algorithm (GEMA) for fitting a log-  
20 Pearson Type III distribution, following the full procedure described in England Jr et al. (2018). This iterative  
21 method jointly estimates distribution parameters (shape  $\alpha$ , scale  $\beta$ , and location  $\tau$ ) from censored and  
22 uncensored data by aligning sample and distribution moments. Convergence was achieved when the L1  
23 norm of moment differences fell below a threshold of  $(1e-10)$ . PILFs were incorporated as left-censored  
24 values below the GBT-derived threshold. This method is particularly suited to datasets with both low  
25 outliers and log-normality. In cases where the user explicitly disabled GEMA or when GEMA failed to  
26 converge, we applied standard moment-based fitting of a log-Pearson III distribution without PILF  
27 treatment. This approach estimates distribution parameters directly from the log-transformed sample  
28 mean, standard deviation, and skewness (England Jr et al., 2018). As a fallback for short time series or when  
29 moment-based fitting was numerically unstable, we estimated return periods using a log-log linear  
30 regression between observed flow magnitudes and empirical exceedance probabilities. Exceedance  
31 probabilities were computed using the Weibull plotting position formula. While this empirical method lacks



1 the theoretical rigor of the log-Pearson III distribution, it provides a practical option for very limited  
2 samples.

### 3 **2.4 Estimation of flood forecasting reliability metrics**

4 CSI, FAR and POD were used as a measure of reliability and were computed over the period 1982 – 2016.  
5 The streamflow thresholds were applied to the simulated streamflow data from GloFAS, Google AI model  
6 and Bakaano-Hydro. Based on this, a model's prediction of an event with a given return period was  
7 considered correct if both the modeled and observed hydrographs exceeded their respective return period  
8 threshold flow values within defined timing tolerance. We used three flood timing tolerance which define  
9 the permissible temporal deviation between observed and predicted flood events when evaluating model  
10 performance. In the 0-day timing tolerance, a predicted flood event must occur on the exact day of the  
11 observed event to be considered a match; the 1-day timing tolerance allows for a match if a predicted flood  
12 occurs within one day (before or after) of the observed event and the 2-day timing tolerance extends the  
13 matching period to two days. CSI is computed as  $TP/(TP + FP + FN)$ ; POD is  $TP/(TP + FN)$  and FAR is  
14 computed as  $FP/(TP + FP)$  where TP is True Positive describing correctly predicted flood events for a  
15 specified return period and flood timing tolerance; FP is False Positives describing predicted flood events  
16 that are not in observed data; and FN is False Negatives describing flood events in the observed data but  
17 were not predicted.

## 18 **3. Results**

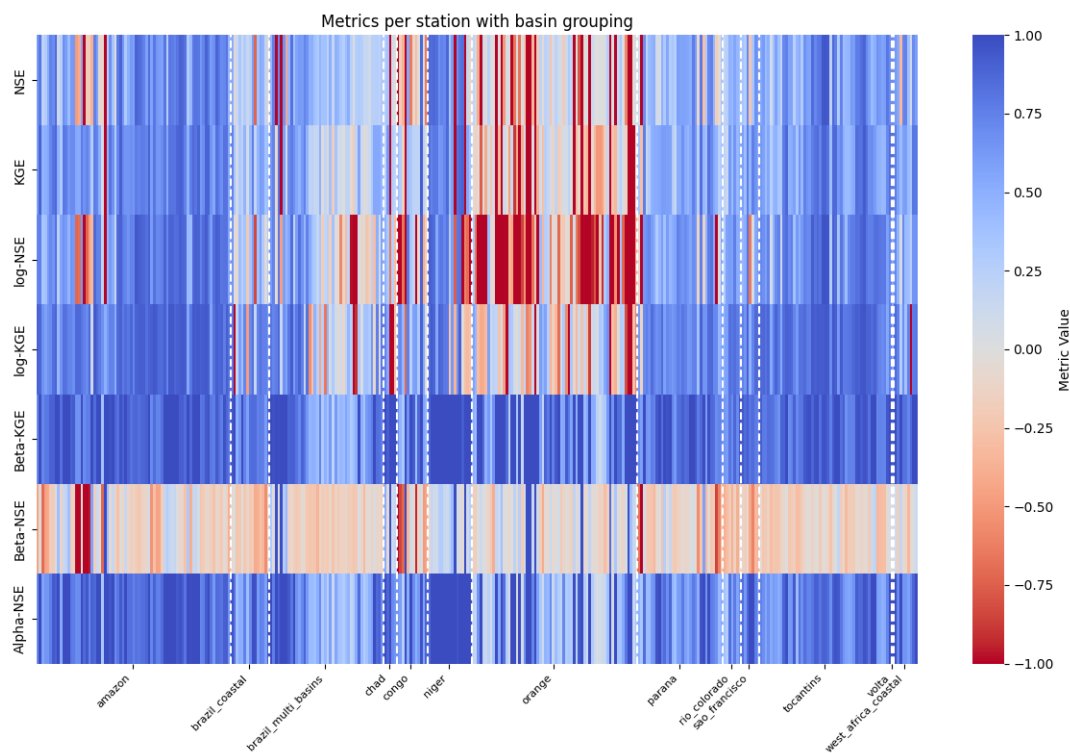
### 19 **3.1 Performance of Bakaano-Hydro**

20 Out-of-sample validation (Figure 2) demonstrates that the Bakaano-Hydro framework achieves strong and  
21 spatially consistent performance across major basins in Africa and South America. Median values of Nash-  
22 Sutcliffe Efficiency (NSE), Kling-Gupta Efficiency (KGE) exceed 0.5 in most basins, indicating robust  
23 agreement between simulated and observed streamflow. Performance remains high even in log-  
24 transformed variants of NSE and KGE, reflecting the model's ability to capture both high- and low-flow  
25 regimes. Evaluation metrics assessing bias (Alpha-NSE) and variability (Beta-KGE) are close to 1.0 across  
26 most regions, suggesting accurate representation of both amplitude and temporal dynamics of flow.  
27 Crucially, this out-of-sample validation ensures that the model's skill reflects true generalization rather than  
28 overfitting to known conditions—an essential distinction in the context of flood forecasting, where  
29 performance metrics often involve comparisons across both training or calibration and test periods due to



1 limited data availability. By validating against held-out time periods, we demonstrate that Bakaano-Hydro  
2 is not merely replicating known hydrological patterns, but can generalize in previously unseen contexts.  
3 This enhances the robustness of our benchmarking framework, enabling a fair and unbiased comparison  
4 with global systems such as GloFAS and the Google AI model. Figure A1 and A2 also compares the  
5 hydrograph of Bakaano-Hydro predicted streamflow against GloFAS and Google AI model for the period  
6 1982 – 2016, which covers both the training and testing period.

7



8

9 **Figure 2.** Out-of-sample evaluation (1982–1988) of Bakaano-Hydro performance across all 470 hydrological  
10 stations used in the forecasting reliability assessment. Stations are grouped by basin, with dashed white  
11 lines indicating basin boundaries. The heatmap displays multiple evaluation metrics, including Nash–  
12 Sutcliffe Efficiency (NSE), Kling–Gupta Efficiency (KGE), log-NSE, and extended variants such as Beta-KGE,  
13 Alpha-NSE, and Beta-NSE (Gupta et al., 2009). The color scale ranges from red (lower performance) to blue  
14 (higher performance), enabling visual comparison of skill across metrics and basins.

15

16

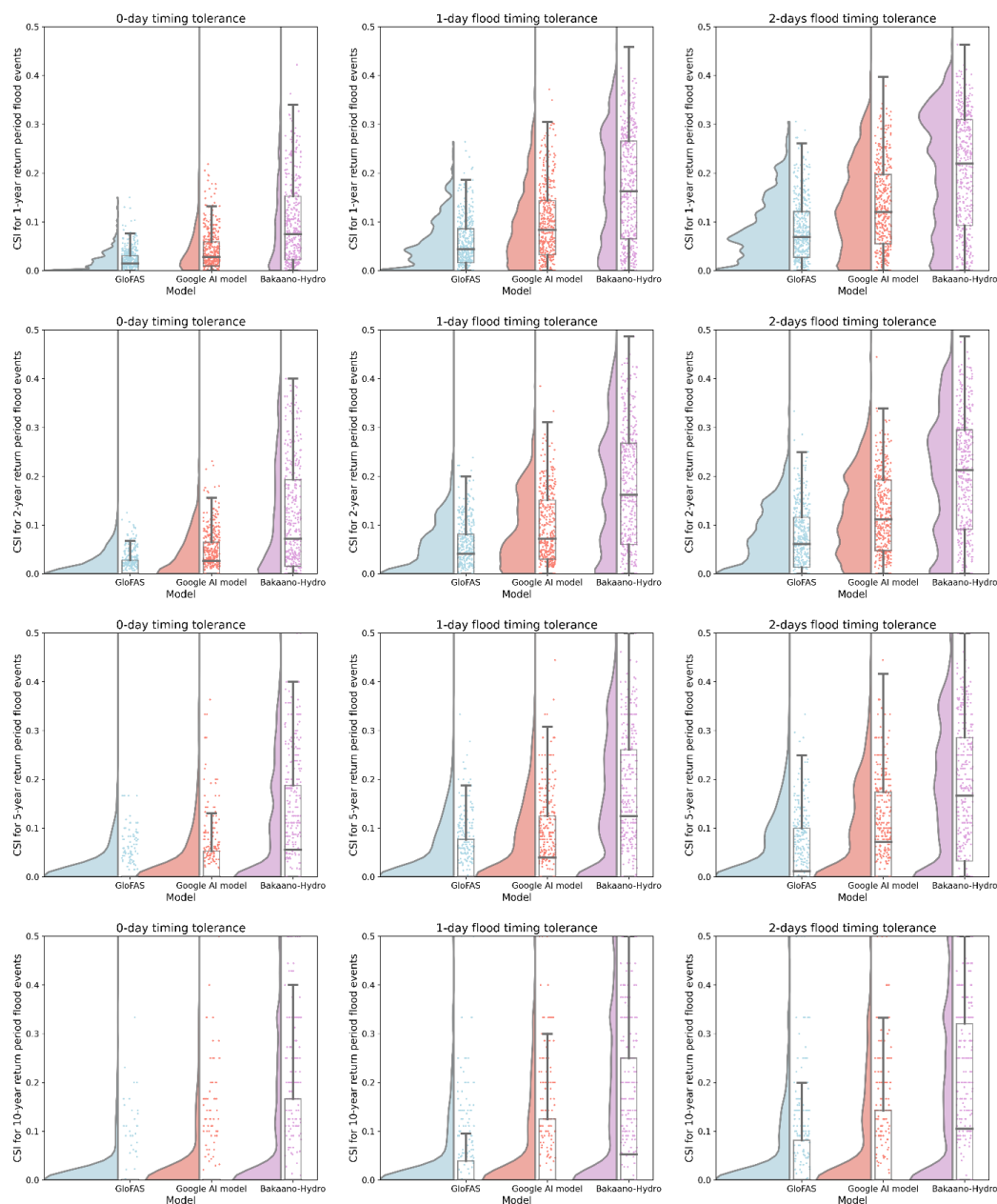


### 1 **3.2 Model intercomparison of flood forecasting reliability**

2 Three metrics—probability of detection (POD), false alarm rate (FAR), and critical success index (CSI)—were  
3 employed to assess reliability of the proposed framework. POD quantifies the proportion of observed flood  
4 events accurately predicted by the models, FAR measures the rate of false positives (incorrect flood  
5 predictions), and CSI evaluates overall prediction accuracy by accounting for true positives and false alarms.  
6 Hindcast skill, assessed through CSI, POD, and FAR, is used as a measure of forecasting reliability. This  
7 hindcasting serves as a proxy for operational forecasting reliability, helping to identify strengths and  
8 limitations in each model's performance and informing future flood predictions. The use of hindcast skill  
9 as a measure of forecast reliability is a standard practice in climate and hydrological forecasting.

10 Similar retrospective evaluations have been conducted to compare GloFAS and Google's AI model (Nearing  
11 et al., 2024), with GloFAS itself relying on historical streamflow simulations to validate its forecast skill  
12 (Alfieri et al., 2013). CSI, POD and FAR values for flood events with 1-, 2-, 5-, and 10-year return periods  
13 and 0-, 1- and 2-day timing tolerance were calculated using simulated data from Bakaano-Hydro for the  
14 period 1982 to 2016. Timing tolerances are permissible temporal deviations between observed and  
15 simulated flood events. Performance was benchmarked against GloFAS and the Google AI model; notably,  
16 the Google AI model data used here were derived from a full model run encompassing all stations, rather  
17 than the cross-validation splits reported in Nearing et al. (2024). To ensure sufficient data for calculating  
18 return-period thresholds, we included both training and out-of-sample streamflow predictions from each  
19 model similar to Nearing et al. (2024) and Alfieri et al. (2013). A common subset of stations across the  
20 three models (470 stations in total) was selected for the calculation of these metrics. Return period  
21 thresholds were computed from observed data at each station and applied uniformly to the predictions of  
22 all models.

23 The results show that the median CSI scores of Bakaano-Hydro at a 0-day timing tolerance are comparable  
24 to or exceed those of GloFAS at 1-day and 2-day timing tolerance, and are comparable to Google AI model's  
25 1-day timing tolerance CSI scores for every return period examined (Figure 3). The median scores of  
26 Bakaano-Hydro at 1-day timing tolerance are also higher than that of Google AI model at 2-day timing  
27 tolerance. These indicate substantial improvement in temporal precision of flood detection by Bakaano-  
28 Hydro. Figure 3 also reveals important differences in the reliability and variability of flood predictions from  
29 GloFAS, Google AI model, and Bakaano-Hydro, as illustrated by the shapes and ranges of their CSI  
30 distributions across Africa and South America. For GloFAS, the distribution skews markedly toward lower  
31 CSI values—particularly under stricter timing conditions (e.g., the 0-day timing tolerance) and at higher



1

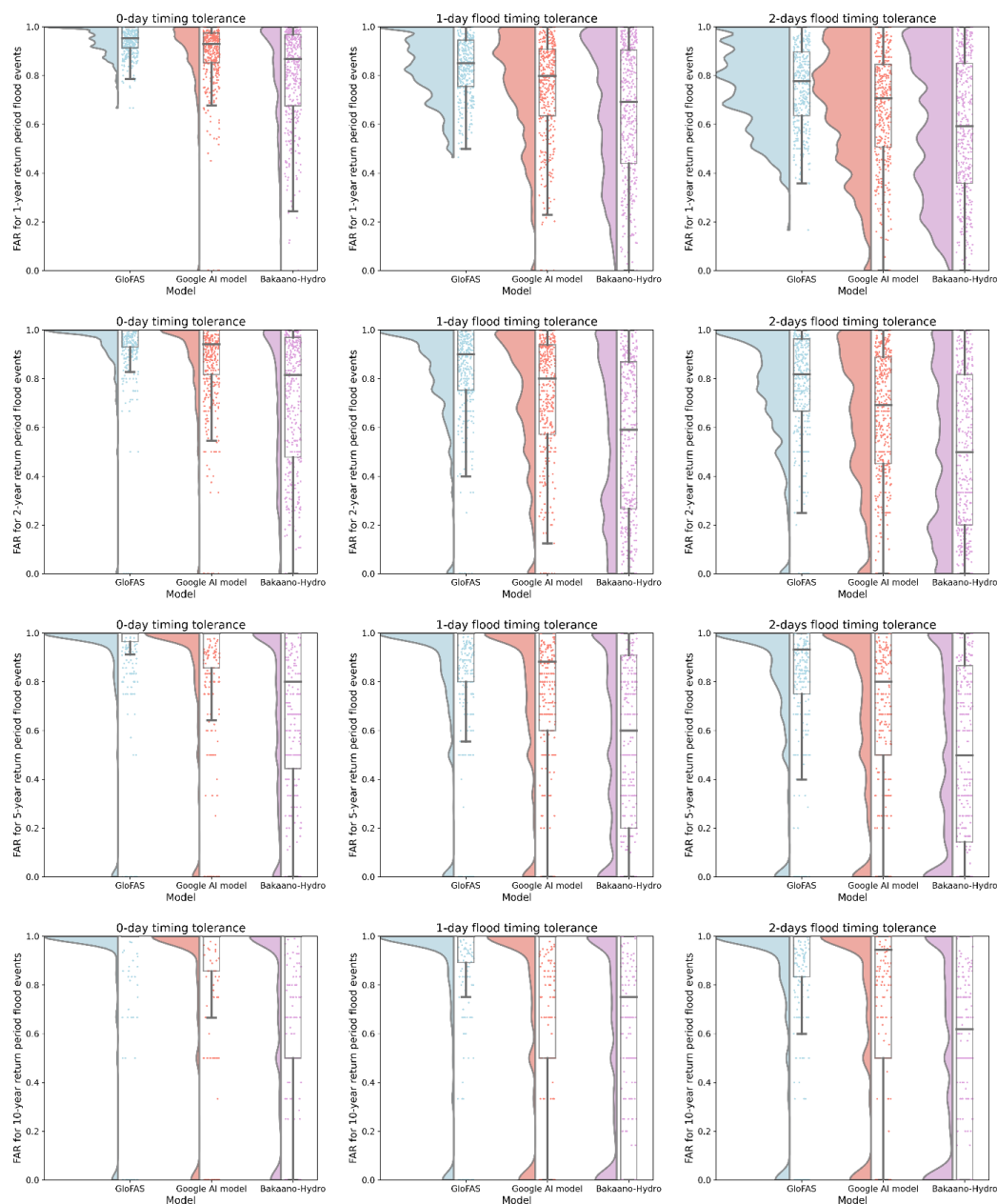
2 **Figure 3.** Comparison of the CSI distribution for multiple return period flood events across three models—  
3 GloFAS, Google, and Bakaano-Hydro—under different timing tolerances (0-day, 1-day, and 2-day) in Africa  
4 and South America. Each panel presents half violin plot, boxplots and density plots to visualize the  
5 distribution of CSI values. The boxplot shows distribution quartiles and whiskers show the full range; the  
6 points show the distribution for the gauging stations and the half-violin plots.



1 return periods—indicating consistently weak performance and a narrow range of outcomes. In contrast,  
2 the Google AI model displays broader distributions, generally clustering within a lower-to-mid CSI range.  
3 While it slightly outperforms GloFAS, its performance still declines considerably for rarer, more extreme  
4 flood events (e.g., from 1-year to 5-year return periods). Notably, Bakaano-Hydro exhibits a wider and more  
5 balanced distribution, with peaks shifting toward higher CSI values, reflecting a greater overall capacity for  
6 reliable flood prediction. Although Bakaano-Hydro’s broader spread indicates some variability—meaning  
7 its performance can occasionally align with lower CSI values—the model consistently demonstrates higher  
8 median CSI scores than GloFAS and Google AI model. As timing tolerances expand from 0 to 2 days, all  
9 three models show higher CSI values, signaling that relaxing the timing tolerance can enhance apparent  
10 predictive skill. Statistical testing show that Bakaano-Hydro’s superior CSI scores over GloFAS and Google  
11 AI model are significant in all return periods and timing tolerances. Comparing Bakaano-Hydro and GloFAS,  
12 p-values ranged from  $5.8 \times 10^{-46}$  (Cohen’s  $d = 0.64$ ) to  $7.0 \times 10^{-11}$  (Cohen’s  $d = 0.19$ ) and from  $1.28 \times 10^{-32}$   
13 (Cohen’s  $d = 0.29$ ) to  $3.5 \times 10^{-4}$  (Cohen’s  $d = 0.065$ ) for Bakaano-Hydro and Google AI model.

14 Across all return periods and timing tolerances, the patterns in FAR (Figure 4) and POD (Figure 5) highlight  
15 important trade-offs that ultimately shape the CSI for each model. For all models, FAR decreases and POD  
16 increases as the timing tolerance widens (1–2 days). GloFAS consistently exhibits the highest FAR,  
17 particularly for shorter return periods and under strict (0-day) timing tolerances, indicating frequent over-  
18 prediction (i.e., many false alarms). Simultaneously, GloFAS shows the lowest POD, suggesting relatively  
19 few correct detections (hits) and a larger number of misses. This high FAR coupled with low detection  
20 directly hinders its CSI, which accounts for both misses and false alarms. By contrast, the Google AI model  
21 demonstrates moderate FAR levels and a broader range of POD outcomes. Nonetheless, the model’s POD  
22 remains more variable and still includes lower-probability detections, especially for shorter return periods.  
23 Finally, Bakaano-Hydro attains the strongest performance across these metrics: it consistently maintains  
24 the lowest FAR and the highest POD. Even under strict 0-day timing tolerances, Bakaano-Hydro’s false  
25 alarms remain minimal, and it achieves more hits relative to misses than the other models do.

26 Statistical testing show that Bakaano-Hydro’s decrease in FAR and increase in POD scores over GloFAS and  
27 Google AI model are significant in all return periods and timing tolerances with p-values substantially lower  
28 than 0.05. Cohen’s  $d$  analyses underscore substantial effect sizes favoring Bakaano-Hydro (ranging from 0.2  
29 to 0.8), most pronounced in comparisons with GloFAS. For FAR, statistical testing of Bakaano-Hydro and  
30 GloFAS showed p-values ranging from  $8.74 \times 10^{-35}$  (Cohen’s  $d = -0.65$ ) to  $1.6 \times 10^{-9}$  (Cohen’s  $d = -0.35$ ) and  
31 p values ranging from  $2.14 \times 10^{-23}$  (Cohen’s  $d = -0.41$ ) to  $1.38 \times 10^{-5}$  (Cohen’s  $d = -0.25$ ) for Bakaano-Hydro



1

2 **Figure 4.** Comparison of the FAR distribution for multiple return period flood events across three models—  
3 GloFAS, Google, and Bakaano-Hydro—under different timing tolerance (0-day, 1-day, and 2-day) in Africa  
4 and South America. Each panel presents half violin plot, boxplots and density plots to visualize the  
5 distribution of FAR values. The boxplot shows distribution quartiles and whiskers show the full range; the  
6 points show the distribution for the gauging stations and the half-violin plots.



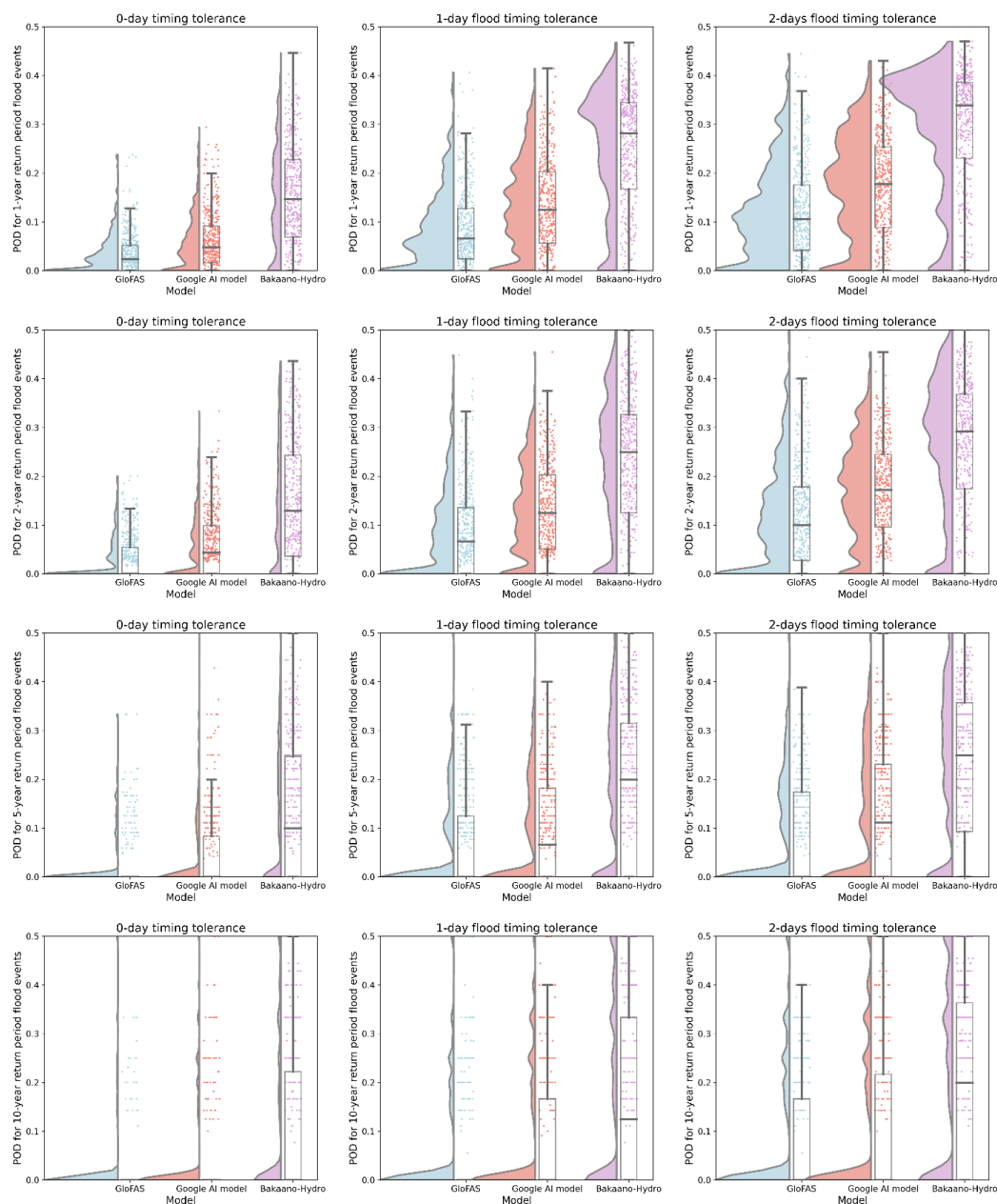
1 and Google AI model. For POD, statistical testing of Bakaano-Hydro and GloFAS showed p-values ranging  
2 from  $1.30 \times 10^{-55}$  (Cohen's  $d = 1.03$ ) to  $6.2 \times 10^{-17}$  (Cohen's  $d = 0.48$ ) and p-values ranging from  $2.2 \times 10^{-46}$   
3 (Cohen's  $d = 0.82$ ) to  $1.2 \times 10^{-10}$  (Cohen's  $d = 0.35$ ) for Bakaano-Hydro and Google AI model.

### 4 **3.3 Spatial variability in forecasting reliability across hydroclimatic regions**

5 Across diverse hydro-climatic contexts, the performance of GloFAS, the Google AI model, and Bakaano-  
6 Hydro in predicting fluvial flooding varies notably, reflecting local hydrological complexities. Overall,  
7 Bakaano-Hydro outperforms GloFAS and Google AI model across diverse hydroclimatic areas with the  
8 improved performance in most basins statistically significant (Figures 6 and 7). Google AI model also shows  
9 strong predictive skill across diverse basins, albeit under 2-day timing tolerance. GloFAS, being a physically-  
10 based model, tends to perform well in larger, more predictable river basins but struggles in basins with  
11 highly dynamic hydrological patterns.

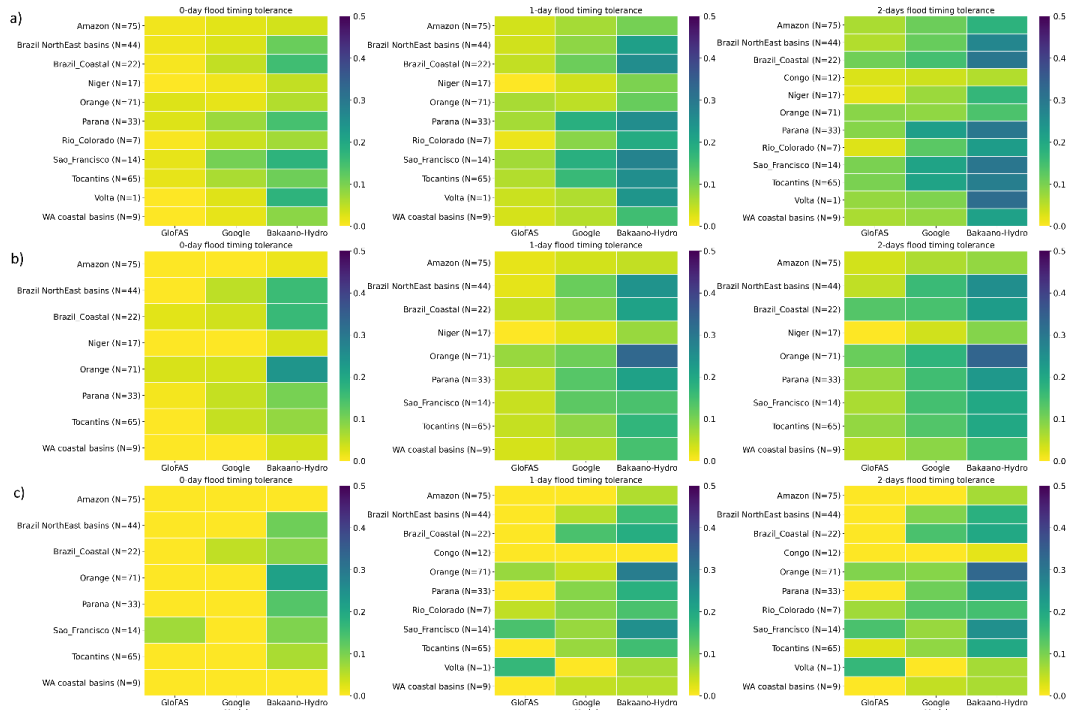
12 Generally, across all models performance is lower in semi-arid to arid areas and higher in humid and  
13 tropical areas. In semi-arid to arid regions—such as the Niger, Orange, Namibia coastal basins—  
14 performance across all models is generally lower, with CSI values indicating difficulty in capturing the  
15 intermittent and highly variable river flows, driven by extreme upstream rainfall that causes rapid river-  
16 level rises, high transmission losses, and complex flood propagation. In these settings, GloFAS exhibits the  
17 lowest CSI values mainly because of the highest false alarm rates, indicating an overestimation of flood  
18 risk, while the Google AI model offers improved detection. Bakaano-Hydro outperforms GloFAS and Google  
19 AI model in these settings with relatively higher CSI values. Figures 6 and 7 show the basins for which  
20 differences in CSI and FAR were statistically significant.





1

2 **Figure 5.** Comparison of the POD distribution for multiple return period flood events across three models—  
3 GloFAS, Google, and Bakaano-Hydro—under different timing tolerance (0-day, 1-day, and 2-day) in Africa  
4 and South America. Each panel presents half violin plot, boxplots and density plots to visualize the  
5 distribution of POD values. The boxplot shows distribution quartiles and whiskers show the full range; the  
6 points show the distribution for the gauging stations and the half-violin plots.

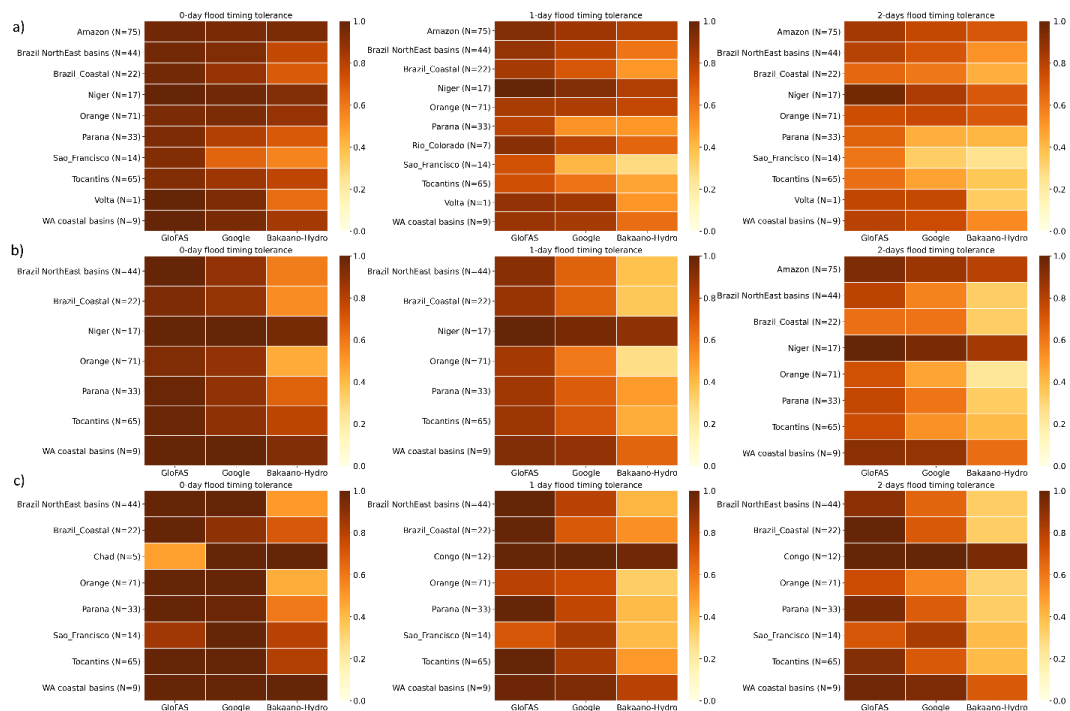


1

2 **Figure 6** Heatmaps of median CSI values for river basins across Africa and South America. a) shows  
3 heatmaps for floods with 1-year return period; b) floods with 2-year return periods; c) floods with 5-year  
4 return periods. For each plot only basins for which differences among the three models were statistically  
5 significant are shown. Brazil NorthEast basins, Brazil\_Coastal and WA coastal basins are not officially  
6 recognized basin names but are used in this paper to refer to a group of small basins in Brazil and South  
7 Africa.

8

9 In tropical and equatorial climates, such as West Africa coastal basins, the Congo, and the Amazon basins,  
10 prolonged and intense rainfall from monsoonal and convective systems introduces another layer of  
11 complexity. Spatially variable rainfall, seasonal soil saturation, and delayed flood wave propagation in large  
12 river systems often lead to overprediction by GloFAS and the Google AI model, as evidenced by higher FAR  
13 and moderate CSI. In temperate and subtropical basins such as the Paraná and Tocantins, models show  
14 consistent performance improvements with longer timing tolerances, indicating better predictability of  
15 hydrological events in these regions. GloFAS performs better in these regions compared to arid areas,  
16 reflecting its suitability for regions with relatively predictable seasonal cycles.



1

2 **Figure 7.** Heatmaps of median FAR values for river basins across Africa and South America. . a) shows  
 3 heatmaps for floods with 1-year return period; b) floods with 2-year return periods; c) floods with 5-year  
 4 return periods. For each plot, only basins for which differences among the three models were statistically  
 5 significant are shown. Brazil NorthEast basins, Brazil\_Coastal and WA coastal basins are not officially  
 6 recognized basin names but are used in this paper to refer to a group of small basins in Brazil and South  
 7 Africa.

8

9 **4. Discussion**

10 The results of this study demonstrate that the distributed hydrology-guided neural network framework  
 11 significantly improves the reliability of flood forecasts in hydrologically diverse, data-scarce regions. Across  
 12 a range of flood return periods, the model achieves higher CSI values and maintains lower FAR and higher  
 13 probability of detection (POD) than both GloFAS and the Google AI model. Notably, Bakaano-Hydro  
 14 outperforms GloFAS and Google AI model under precise timing constraints. Its performance at 0-day timing  
 15 tolerance exceeds or matches GloFAS' performance at 2-day and Google AI model performance at 1-day  
 16 timing tolerance. This enhanced temporal precision is critical for real-time flood early warning, especially  
 17 in regions with high vulnerability and limited emergency response capacity. Importantly, in this paper, I



1 adopt a consistent benchmarking methodology by using observed streamflow data to define flood  
2 thresholds across all models—unlike model-specific thresholds that can artificially inflate or deflate  
3 performance metrics (e.g. Nearing et al., 2024). This evaluation approach enables a more objective and  
4 physically meaningful comparison, ensuring that improvements in forecast skill are attributable to model  
5 behavior rather than calibration biases.

6 Bakaano-Hydro contributes a new paradigm to the field of hydrological modeling by demonstrating that  
7 fully distributed, process-informed neural networks can reliably simulate streamflow and flood dynamics.  
8 It advances the state of the art in data-driven hydrological modelling by integrating process-based runoff  
9 generation and topographic flow routing with a deep learning architecture capable of generalizing across  
10 basins, climatic zones, and hydrological regimes. This overcomes existing trade-off between physical  
11 realism and data-driven scalability, offering a unified approach that performs well even in environments  
12 with limited or fragmented observational records. A key scientific insight is that Bakaano-Hydro achieves  
13 increased reliability despite being trained only on data from Africa and South America, regions typically  
14 underrepresented in global hydrological training datasets. In contrast, the Google AI model (Nearing et al.,  
15 2024) model draws on global data, including from data-rich regions such as Europe, North America, and  
16 Australia. That Bakaano-Hydro performs better or comparably—particularly in Africa and South America—  
17 highlights not only the framework’s efficiency in learning from sparse, heterogeneous data but also the  
18 advantages of physically guided spatial representation over lumped learning schemes.

19 The development and evaluation of Bakaano-Hydro has far-reaching implications for climate resilience,  
20 disaster risk reduction, and sustainable development, particularly in the Global South. By offering a  
21 distributed, interpretable, and generalizable modeling framework, Bakaano-Hydro bridges a critical gap  
22 between global early warning systems and locally relevant flood risk information—a gap that has  
23 historically undermined emergency preparedness in data-scarce regions. For national governments and  
24 hydrometeorological agencies, Bakaano-Hydro offers a pathway toward strengthening domestic early  
25 warning capacity without requiring dense monitoring networks or extensive model calibration. Its modular  
26 and open-source architecture enables customization for local use cases—whether for national disaster  
27 response systems, basin authorities, or regional climate services. This opens up the possibility of  
28 transitioning from dependency on external, often coarse-resolution forecasts (like GloFAS or Google AI  
29 model) to context-specific, high-resolution models that are co-developed and owned by local institutions.  
30 For the humanitarian and development sector, Bakaano-Hydro provides a decision-support tool for  
31 anticipatory action, allowing earlier and more precise deployment of flood relief, social protection



1 schemes, and infrastructure safeguards. Its ability to generalize to ungauged basins also supports climate  
2 adaptation planning and nature-based solutions, such as ecosystem restoration in floodplains and  
3 watershed management. At the global scale, Bakaano-Hydro advances the science-policy interface by  
4 demonstrating that physically informed machine learning models can be both scalable and regionally  
5 accurate—paving the way for hybrid forecasting systems that serve diverse geographies, socio-economic  
6 contexts, and governance capacities. From a scientific standpoint, the framework opens new pathways for  
7 interpretable machine learning in hydrology. By using runoff as an intermediate, physically meaningful  
8 variable, Bakaano-Hydro avoids the black-box pitfalls of purely statistical models and supports diagnostic  
9 evaluation of model behavior. The architecture also enables flexible experimentation with modular  
10 components—allowing researchers to test how different land surface models or routing schemes influence  
11 predictive performance across space and time. Operationally, Bakaano-Hydro offers a practical and scalable  
12 solution for flood forecasting in regions with limited data infrastructure, where conventional physically  
13 based models struggle to calibrate, and global lumped models fail to resolve local dynamics.

14

15 A primary limitation of the current study is its use of historical observed discharge records rather than true  
16 forecast or reforecast datasets. While this retrospective analysis enables rigorous comparison across  
17 models and return periods, it does not fully capture the operational uncertainties present in real-time  
18 forecasting—such as delays in data assimilation, atmospheric forecast errors, or recent changes in land use  
19 and hydraulic infrastructure. Moreover, under ongoing climate change, stationarity assumptions may  
20 become less valid, potentially limiting the relevance of hindcasting performance as a predictor of future  
21 skill. Nonetheless, the long-term, multi-basin dataset used in this study spans a wide range of hydrological  
22 conditions and extreme events, offering a robust foundation for model validation. A model that performs  
23 reliably across these historical extremes is well-positioned to adapt to future conditions when coupled with  
24 evolving inputs such as satellite rainfall data and climate forecast ensembles.

## 25 **5. Conclusion**

26 This study demonstrates the potential of a distributed, hydrology-guided neural network framework that  
27 integrates process-based runoff generation, topographic flow routing, and temporal convolutional  
28 networks to improve flood forecasting reliability in data-scarce regions. The framework—Bakaano-Hydro—  
29 was benchmarked against GloFAS and the Google AI model using retrospective simulations across 470  
30 gauging stations in Africa and South America. Results show that Bakaano-Hydro consistently achieves  
31 higher forecasting reliability across multiple flood return periods and timing tolerances. Statistically



1 significant improvements in critical reliability metrics, including the critical success index (CSI), probability  
2 of detection (POD), and false alarm rate (FAR), indicate the model's robustness in capturing both the  
3 occurrence and timing of flood events.

4 Importantly, Bakaano-Hydro advances the state of data-driven hydrological modeling by embedding  
5 physically meaningful processes within a neural network architecture. This hybrid approach addresses key  
6 limitations of lumped data-driven models and calibration-intensive physically based systems, offering a  
7 generalizable and interpretable alternative suited to diverse hydroclimatic contexts. The ability of the  
8 framework to generalize from sparse training data, while maintaining physically plausible representations  
9 of runoff and flow routing, highlights its potential to support robust hydrological analysis in observationally  
10 limited settings. Overall, Bakaano-Hydro contributes a methodologically rigorous and scientifically  
11 grounded step forward in the development of hybrid flood forecasting models.

## 12 **Code availability**

13 Source codes for Bakaano-Hydro are available at <https://github.com/confidence-duku/bakaano-hydro>.  
14 The neural network architecture in this version was adapted for this study and can be found at  
15 <https://doi.org/10.5281/zenodo.15322955>. All codes for Bakaano-Hydro evaluation and return period  
16 estimation can be found at  
17 [https://github.com/google-research-datasets/global\\_streamflow\\_model\\_paper](https://github.com/google-research-datasets/global_streamflow_model_paper).

## 18 **Data availability**

19 Daily streamflow data for 1982 – 2016 simulated for hydrological stations across Africa and South America  
20 as well as validation and CSI, POD, FAR results are available at <https://doi.org/10.5281/zenodo.15322955>.  
21 GloFAS and Google AI model benchmark data are available at <https://doi.org/10.5281/zenodo.10397664>

## 22 **Competing interests**

23 The authors declare that they have no conflict of interest.

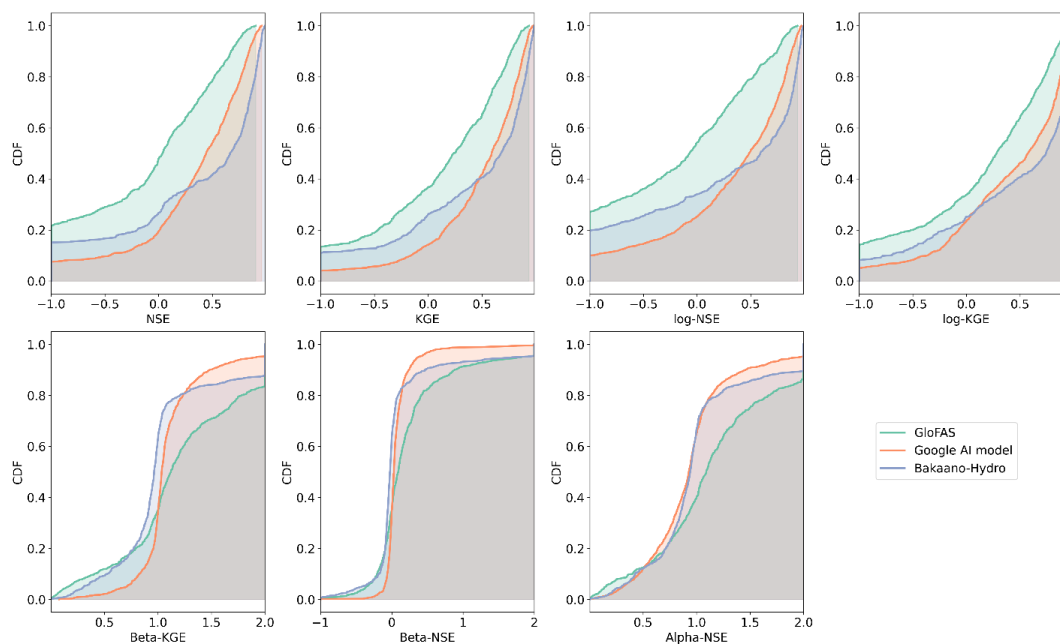
24

25

26

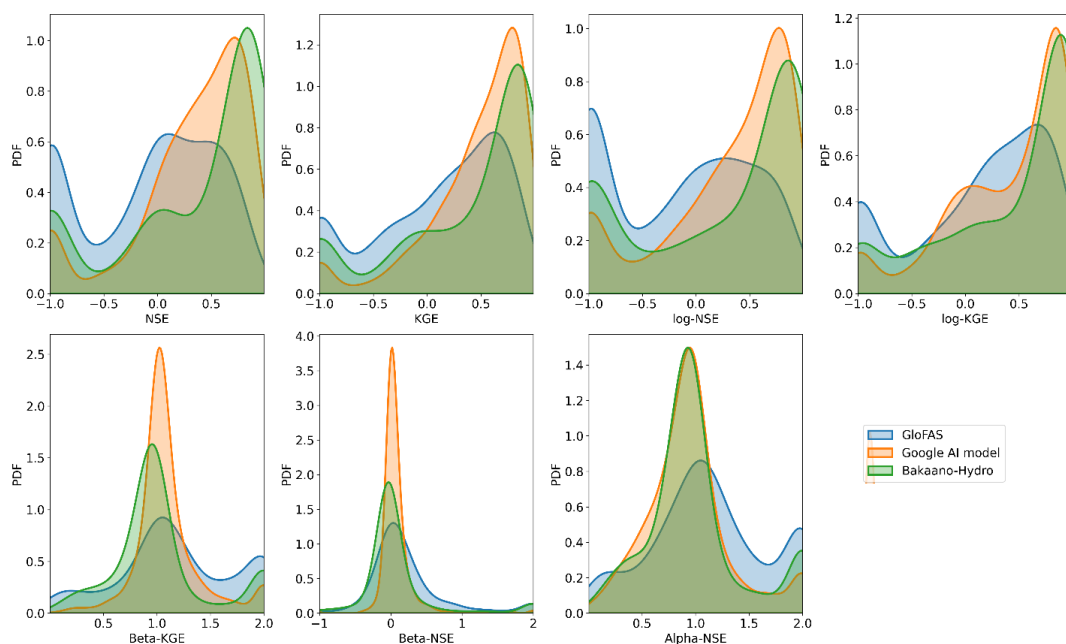


## 1 Appendix A



2

3 **Figure A1.** Comparison of cumulative distribution functions (CDFs) of hydrological performance metrics  
4 across values for the three models - GloFAS, Google AI model, and Bakaano-Hydro - for the period 1982  
5 - 2016. Metrics include Nash-Sutcliffe Efficiency (NSE), Kling-Gupta Efficiency (KGE), log-NSE and  
6 additional variants such as Beta-KGE, Alpha-NSE and Beta-NSE



1

2 Figure A2. Comparison of probability density functions (PDFs) of hydrological performance metrics across  
 3 values for the three models - GloFAS, Google AI model, and Bakaano-Hydro - for the period 1982 - 2016.  
 4 Metrics include Nash-Sutcliffe Efficiency (NSE), Kling-Gupta Efficiency (KGE), log-NSE and additional  
 5 variants such as Beta-KGE, Alpha-NSE and Beta-NSE

6

7

## 8 References

- 9 Alfieri, L., Avanzi, F., Delogu, F., Gabellani, S., Bruno, G., Campo, L., Libertino, A., Massari, C., Tarpanelli, A.,  
 10 Rains, D., Miralles, D. G., Quast, R., Vreugdenhil, M., Wu, H., & Brocca, L. (2022). High-resolution  
 11 satellite products improve hydrological modeling in northern Italy. *Hydrol. Earth Syst. Sci.*, 26(14),  
 12 3921-3939. <https://doi.org/10.5194/hess-26-3921-2022>
- 13 Alfieri, L., Burek, P., Dutra, E., Krzeminski, B., Muraro, D., Thielen, J., & Pappenberger, F. (2013). GloFAS–  
 14 global ensemble streamflow forecasting and flood early warning. *Hydrology and Earth System  
 15 Sciences*, 17(3), 1161-1175.
- 16 Alfieri, L., Lorini, V., Hirpa, F. A., Harrigan, S., Zsoter, E., Prudhomme, C., & Salamon, P. (2020). A global  
 17 streamflow reanalysis for 1980–2018. *Journal of Hydrology X*, 6, 100049.  
 18 <https://doi.org/https://doi.org/10.1016/j.hydroa.2019.100049>
- 19 AON. (2023). *Weather, Climate and Catastrophe Insight*.
- 20 Bai, S., Kolter, J. Z., & Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent  
 21 networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 10.





- 1 Beck, H. E., van Dijk, A. I. J. M., de Roo, A., Miralles, D. G., McVicar, T. R., Schellekens, J., & Bruijnzeel, L. A.  
2 (2016). Global-scale regionalization of hydrologic model parameters. *Water Resources Research*,  
3 52(5), 3599-3622. [https://doi.org/https://doi.org/10.1002/2015WR018247](https://doi.org/10.1002/2015WR018247)
- 4 Crochemore, L., Ramos, M. H., & Pappenberger, F. (2016). Bias correcting precipitation forecasts to  
5 improve the skill of seasonal streamflow forecasts. *Hydrol. Earth Syst. Sci.*, 20(9), 3601-3618.  
6 <https://doi.org/10.5194/hess-20-3601-2016>
- 7 De Roo, A. P. J., Wesseling, C. G., & Van Deursen, W. P. A. (2000). Physically based river basin modelling  
8 within a GIS: the LISFLOOD model. *Hydrological Processes*, 14(11-12), 1981-1992.  
9 [https://doi.org/https://doi.org/10.1002/1099-1085\(20000815/30\)14:11/12<1981::AID-](https://doi.org/https://doi.org/10.1002/1099-1085(20000815/30)14:11/12<1981::AID-)  
10 [HYP49>3.0.CO;2-F](https://doi.org/10.1002/1099-1085(20000815/30)14:11/12<1981::AID-HYP49>3.0.CO;2-F)
- 11 Didan, K. (2021). *MODIS/Terra Vegetation Indices 16-Day L3 Global 250m SIN Grid V061 [Data set]*. NASA  
12 *EOSDIS Land Processes Distributed Active Archive Center*.
- 13 DiMiceli, C., Carroll, M., Sohlberg, R., Kim, D.-H., Kelly, M., & Townshend, J. (2015). MOD44B MODIS/Terra  
14 vegetation continuous fields yearly L3 global 250m SIN grid V006. *NASA EOSDIS Land Processes*  
15 *DAAC*, 10.
- 16 Emerton, R. E., Stephens, E. M., Pappenberger, F., Pagano, T. C., Weerts, A. H., Wood, A. W., Salamon, P.,  
17 Brown, J. D., Hjerdt, N., Donnelly, C., Baugh, C. A., & Cloke, H. L. (2016). Continental and global  
18 scale flood forecasting systems. *WIREs Water*, 3(3), 391-418.  
19 <https://doi.org/https://doi.org/10.1002/wat2.1137>
- 20 England Jr, J. F., Cohn, T. A., Faber, B. A., Stedinger, J. R., Thomas Jr, W. O., Veilleux, A. G., Kiang, J. E., &  
21 Mason Jr, R. R. (2018). *Guidelines for determining flood flow frequency—Bulletin 17C*  
22 (14111342232).
- 23 Feng, D., Fang, K., & Shen, C. (2020). Enhancing Streamflow Forecast and Extracting Insights Using Long-  
24 Short Term Memory Networks With Data Integration at Continental Scales. *Water Resources*  
25 *Research*, 56(9), e2019WR026793. <https://doi.org/https://doi.org/10.1029/2019WR026793>
- 26 Fortin, F.-A., De Rainville, F.-M., Gardner, M.-A. G., Parizeau, M., & Gagné, C. (2012). DEAP: Evolutionary  
27 algorithms made easy. *The Journal of Machine Learning Research*, 13(1), 2171-2175.
- 28 Fox, S., Agyemang, F., Hawker, L., & Neal, J. (2024). Integrating social vulnerability into high-resolution  
29 global flood risk mapping. *Nature Communications*, 15(1), 3155. [https://doi.org/10.1038/s41467-](https://doi.org/10.1038/s41467-024-47394-2)  
30 [024-47394-2](https://doi.org/10.1038/s41467-024-47394-2)
- 31 GRDC. (2025). The Global Runoff Data Centre, 56068 Koblenz, Germany (<https://grdc.bafg.de/>). In.  
32 Grimaldi, S. (2024). *GloFAS v4 calibration methodology and parameters*.
- 33 Grimes, D. R., Rogers, D. P., Schumann, A., & Day, B. F. (2022). *Charting a Course for Sustainable*  
34 *Hydrological and Meteorological Observation Networks in Developing Countries*. World Bank.
- 35 Guo, Y., Zhang, Y., Zhang, L., & Wang, Z. (2021). Regionalization of hydrological modeling for predicting  
36 streamflow in ungauged catchments: A comprehensive review. *WIREs Water*, 8(1), e1487.  
37 <https://doi.org/https://doi.org/10.1002/wat2.1487>
- 38 Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). Decomposition of the mean squared error  
39 and NSE performance criteria: Implications for improving hydrological modelling. *Journal of*  
40 *Hydrology*, 377(1), 80-91. <https://doi.org/https://doi.org/10.1016/j.jhydrol.2009.08.003>
- 41 Hallegatte, S. (2012). A cost effective solution to reduce disaster losses in developing countries: hydro-  
42 meteorological services, early warning, and evacuation. *World Bank policy research working*  
43 *paper*(6058).
- 44 Hargreaves, G. H., & Samani, Z. A. (1985). Reference crop evapotranspiration from temperature. *Applied*  
45 *engineering in agriculture*, 1(2), 96-99.
- 46 Harrigan, S., Zsoter, E., Cloke, H., Salamon, P., & Prudhomme, C. (2023). Daily ensemble river discharge  
47 reforecasts and real-time forecasts from the operational Global Flood Awareness System. *Hydrol.*  
48 *Earth Syst. Sci.*, 27(1), 1-19. <https://doi.org/10.5194/hess-27-1-2023>



- 1 Hunt, K. M. R., Matthews, G. R., Pappenberger, F., & Prudhomme, C. (2022). Using a long short-term  
2 memory (LSTM) neural network to boost river streamflow forecasts over the western United  
3 States. *Hydrol. Earth Syst. Sci.*, 26(21), 5449-5472. <https://doi.org/10.5194/hess-26-5449-2022>
- 4 Karger, D. N., Lange, S., Hari, C., Reyer, C. P. O., Conrad, O., Zimmermann, N. E., & Frieler, K. (2023).  
5 CHELSA-W5E5: daily 1 km meteorological forcing data for climate impact studies. *Earth Syst. Sci.*  
6 *Data*, 15(6), 2445-2464. <https://doi.org/10.5194/essd-15-2445-2023>
- 7 Karger, D. N. L., S; Hari, C.; Reyer, P. O. C.; Zimmermann, E. N. (2021). *CHELSA-W5E5 v1.1: W5E5 v1.0*  
8 *downscaled with CHELSA v2.0* ISIMIP Repository.
- 9 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall–runoff modelling using  
10 Long Short-Term Memory (LSTM) networks. *Hydrol. Earth Syst. Sci.*, 22(11), 6005-6022.  
11 <https://doi.org/10.5194/hess-22-6005-2018>
- 12 Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., & Nearing, G. S. (2019). Toward  
13 improved predictions in ungauged basins: Exploiting the power of machine learning. *Water*  
14 *Resources Research*, 55(12), 11344-11354.
- 15 Lehner, B., Verdin, K., & Jarvis, A. (2008). New Global Hydrography Derived From Spaceborne Elevation  
16 Data. *Eos, Transactions American Geophysical Union*, 89(10), 93-94.  
17 <https://doi.org/https://doi.org/10.1029/2008EO100001>
- 18 Mazzoleni, M., Dottori, F., Cloke, H. L., & Di Baldassarre, G. (2022). Deciphering human influence on  
19 annual maximum flood extent at the global level. *Communications Earth & Environment*, 3(1),  
20 262. <https://doi.org/10.1038/s43247-022-00598-0>
- 21 Nearing, G., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., Hassidim, A., Klotz, D., Kratzert, F.,  
22 Metzger, A., Nevo, S., Pappenberger, F., Prudhomme, C., Shalev, G., Shenzi, S., Tekalign, T. Y.,  
23 Weitzner, D., & Matias, Y. (2024). Global prediction of extreme floods in ungauged watersheds.  
24 *Nature*, 627(8004), 559-563. <https://doi.org/10.1038/s41586-024-07145-1>
- 25 Nevo, S., Morin, E., Gerzi Rosenthal, A., Metzger, A., Barshai, C., Weitzner, D., Voloshin, D., Kratzert, F.,  
26 Elidan, G., Dror, G., Begelman, G., Nearing, G., Shalev, G., Noga, H., Shavitt, I., Yuklea, L., Royz, M.,  
27 Giladi, N., Peled Levi, N., . . . Matias, Y. (2022). Flood forecasting with machine learning models in  
28 an operational framework. *Hydrol. Earth Syst. Sci.*, 26(15), 4013-4032.  
29 <https://doi.org/10.5194/hess-26-4013-2022>
- 30 Pagliero, L., Bouraoui, F., Diels, J., Willems, P., & McIntyre, N. (2019). Investigating regionalization  
31 techniques for large-scale hydrological modelling. *Journal of Hydrology*, 570, 220-235.  
32 <https://doi.org/https://doi.org/10.1016/j.jhydrol.2018.12.071>
- 33 Pilon, P. J. (2002). *Guidelines for reducing flood losses*.
- 34 Poggio, L., de Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E., & Rossiter, D. (2021).  
35 SoilGrids 2.0: producing soil information for the globe with quantified spatial uncertainty. *SOIL*,  
36 7(1), 217-240. <https://doi.org/10.5194/soil-7-217-2021>
- 37 Quinn, P., Beven, K., Chevallier, P., & Planchon, O. (1991). The prediction of hillslope flow paths for  
38 distributed hydrological modelling using digital terrain models. *Hydrological Processes*, 5(1), 59-  
39 79.
- 40 Rentschler, J., Salhab, M., & Jafino, B. A. (2022). Flood exposure and poverty in 188 countries. *Nature*  
41 *Communications*, 13(1), 3527. <https://doi.org/10.1038/s41467-022-30727-4>
- 42 Rogers, D., & Tsirkunov, V. (2011). Costs and benefits of early warning systems. *Global assessment rep.*
- 43 Sauer, I. J., Mester, B., Frieler, K., Zimmermann, S., Schewe, J., & Otto, C. (2024). Limited progress in global  
44 reduction of vulnerability to flood impacts over the past two decades. *Communications Earth &*  
45 *Environment*, 5(1), 239. <https://doi.org/10.1038/s43247-024-01401-y>
- 46 Senay, G. B. (2008). Modeling Landscape Evapotranspiration by Integrating Land Surface Phenology and a  
47 Water Balance Algorithm. *Algorithms*, 1(2), 52-68. <https://www.mdpi.com/1999-4893/1/2/52>



- 1 Senay, G. B., Kagone, S., Parrish, G. E. L., Khand, K., Boiko, O., & Velpuri, N. M. (2023). Improvements and  
2 Evaluation of the Agro-Hydrologic VegET Model for Large-Area Water Budget Analysis and  
3 Drought Monitoring. *Hydrology*, 10(8), 168. <https://www.mdpi.com/2306-5338/10/8/168>
- 4 Tanguy, M., Eastman, M., Chevuturi, A., Magee, E., Cooper, E., Johnson, R. H. B., Facer-Childs, K., &  
5 Hannaford, J. (2025). Optimising ensemble streamflow predictions with bias correction and data  
6 assimilation techniques. *Hydrol. Earth Syst. Sci.*, 29(6), 1587-1614. [https://doi.org/10.5194/hess-](https://doi.org/10.5194/hess-29-1587-2025)  
7 [29-1587-2025](https://doi.org/10.5194/hess-29-1587-2025)
- 8 Tellman, B., Sullivan, J. A., Kuhn, C., Kettner, A. J., Doyle, C. S., Brakenridge, G. R., Erickson, T. A., &  
9 Slayback, D. A. (2021). Satellite imaging reveals increased proportion of population exposed to  
10 floods. *Nature*, 596(7870), 80-86.
- 11 Wang, Y., & Karimi, H. A. (2022). Impact of spatial distribution information of rainfall in runoff simulation  
12 using deep learning method. *Hydrol. Earth Syst. Sci.*, 26(9), 2387-2403.  
13 <https://doi.org/10.5194/hess-26-2387-2022>
- 14 WMO. (2013). The global climate 2001–2010: A decade of climate extremes. *WMO-No. 1103*, 119.
- 15 Wongchuig, S., Paiva, R., Siqueira, V., Papa, F., Fleischmann, A., Biancamaria, S., Paris, A., Parrens, M., & Al  
16 Bitar, A. (2024). Multi-Satellite Data Assimilation for Large-Scale Hydrological-Hydrodynamic  
17 Prediction: Proof of Concept in the Amazon Basin. *Water Resources Research*, 60(8),  
18 e2024WR037155. <https://doi.org/https://doi.org/10.1029/2024WR037155>
- 19 Yaghmaei, N. (2020). *Human cost of disasters: An overview of the last 20 years, 2000-2019*. UN Office for  
20 Disaster Risk Reduction.
- 21
- 22