



25

26

27 Keywords: flood early warning systems; machine learning; White Volta Basin; Ghana; Sentinel-

28 1; CHIRPS; Bagre Dam; alert thresholds; leave-one-year-out cross-validation

29

30 **1. Introduction**

31 Flooding kills more people in West Africa than any other natural hazard, and the trend is getting
32 worse. Across the region, a statistically significant upward trajectory in flood frequency has been
33 documented over recent decades, with event rates exceeding 0.50 additional floods per year in
34 parts of the Volta Basin (Koubodana et al. 2025). Ghana sits near the centre of this pattern. The
35 White Volta Basin in the country's north experiences annual inundation that regularly destroys
36 crops, displaces communities, and kills livestock, yet the timing and severity of the worst events
37 are often determined not by local rainfall but by operational decisions made across an international
38 border, at Burkina Faso's Bagre Dam.

39 The Bagre Reservoir was constructed primarily for irrigation and hydroelectric generation, but its
40 releases during high-storage periods can send flood pulses downstream into Ghana within two to
41 three days. Communities in the Central Gonja and Savannah districts have developed a deep
42 familiarity with this dynamic over generations, but formal early warning systems have yet to
43 integrate it. The myDEWETRA-VOLTALARM platform, deployed in the White Volta Basin
44 through a partnership between Ghanaian agencies and the Italian Civil Protection Department,
45 issues five-day probabilistic flood forecasts using GloFAS global model output, a system not
46 calibrated to local conditions and without explicit representation of Bagre reservoir state



47 (Katsekor et al. 2025a). As a consequence, warnings can arrive too late, cover the wrong locations,
48 or fail to differentiate between a locally driven rainfall event and a Bagre-release flood that
49 demands a categorically different response.

50 The literature on machine learning for hydrological forecasting has expanded rapidly over the past
51 decade. Long Short-Term Memory networks have demonstrated strong skill for streamflow
52 prediction in both gauged and ungauged basins, in some settings outperforming process-based
53 models when trained on sufficiently large datasets (Kratzert et al. 2018). Ensemble combinations
54 of LSTM with gradient-boosted trees have further improved skill for extreme event detection,
55 which is the operationally critical capability for a flood warning system. In the White Volta Basin
56 itself, Katsekor et al. (2025b) demonstrated that LSTM and Random Forest models trained on
57 CHIRPS rainfall and ERA5 reanalysis data can predict discharge at 1-, 5-, and 10-day intervals
58 with high skill, establishing the methodological precedent for the present study. At global scale,
59 Nearing et al. (2024) showed that AI-based systems can match the accuracy of next-day GloFAS
60 nowcasts at five-day lead time in ungauged watersheds, a finding that suggests locally calibrated
61 models with access to gauge data should perform substantially better still, provided the evaluation
62 is conducted on a genuinely independent period.

63 What has not been demonstrated for the White Volta Basin, or, to our knowledge, for any
64 comparable West African river system, is a complete end-to-end early warning system design that
65 moves from satellite data ingestion through AI forecasting, calibrated threshold-based alert
66 generation, satellite-derived inundation validation, and a concrete operational architecture
67 compatible with Ghana's institutional infrastructure. The gap between a well-performing discharge
68 prediction model and an operational flood warning system is substantial. It involves translating
69 continuous discharge forecasts into discrete, actionable alert tiers; calibrating those tiers to return



70 periods that are meaningful to emergency managers and communities; validating that threshold
71 exceedances correspond to real inundation; and specifying how the alert output integrates with the
72 existing institutional chain from national forecasting centre to at-risk household. Each of these
73 steps requires deliberate design choices that go beyond model performance optimisation, and each
74 is addressed in the present study.

75 This paper makes four specific contributions. First, it incorporates Bagre reservoir upstream
76 storage as an explicit model input, a variable that influences downstream flood risk profoundly but
77 has been absent from all previous formal modelling work on the White Volta. Second, it designs
78 a four-tier alert classification system calibrated to 30-year flood frequency analysis at Nawuni,
79 producing discharge thresholds grounded in the statistical properties of the observed record. Third,
80 it validates alert thresholds against Sentinel-1 synthetic aperture radar inundation maps for three
81 major flood events, providing spatial evidence that the threshold-to-inundation relationship is
82 empirically defensible. Fourth, it specifies an operational architecture that integrates directly into
83 Ghana's existing myDEWETRA-VOLTALARM platform and NADMO dissemination chain,
84 requiring no new institutional infrastructure.

85 The paper is structured as follows. Section 2 describes the study area and its flood climatology.
86 Section 3 details the datasets and their seasonal characteristics. Section 4 presents the methods for
87 feature engineering, ensemble modelling, alert threshold design, and Sentinel-1 validation. Section
88 5 reports results across all four components. Section 6 discusses performance, limitations, and
89 operational implications. Section 7 concludes.

90

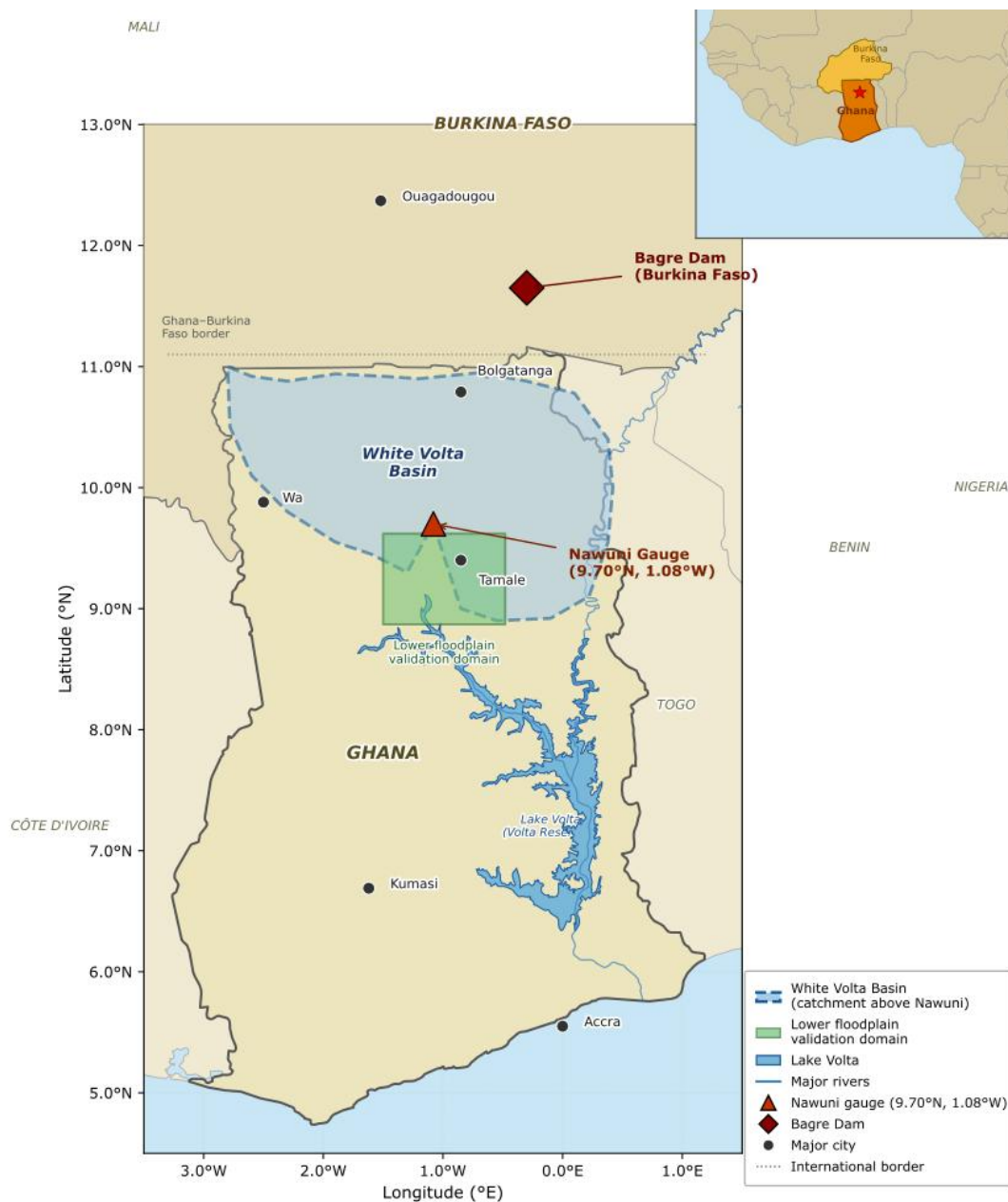
91 **2. Study Area**



92 The White Volta River drains approximately 106,000 km² across Burkina Faso and northern Ghana
93 before converging with the Black Volta near Yapei to form the main Volta River. The basin's
94 hydrology is governed by a single wet season running from roughly June to October, driven by the
95 West African Monsoon, and is strongly modulated upstream by reservoir operations at the Bagre
96 Dam in Burkina Faso. Our modelling domain is anchored to the Nawuni gauging station (9.7° N,
97 1.08° W), which captures the integrated discharge of the White Volta across a catchment of 92,950
98 km² before it enters the Volta Reservoir. Nawuni is the most data-rich gauge in the system and has
99 served as the primary monitoring point in previous machine learning studies of the basin
100 (Katsektor et al. 2025b). The spatial configuration of the basin, including the locations of the



101 Nawuni gauge, the Bagre Dam, and the lower floodplain validation domain, is shown in Figure 1.



102

103 **Figure 1.** Study area map showing the White Volta Basin above the Nawuni gauging station (9.7° N, 1.08° W;
104 catchment area 92,950 km²), the Bagre Dam in Burkina Faso, and the lower White Volta floodplain validation domain.
105 The inset shows the location of Ghana in West Africa.



106 The climate across the basin follows a Sudano-Sahelian gradient, transitioning from dry sub-humid
107 conditions in southern Ghana to semi-arid conditions in central Burkina Faso. Mean annual rainfall
108 at Tamale, the principal city of northern Ghana, is approximately 1,100 mm, concentrated almost
109 entirely between May and October. The basin experiences pronounced dry seasons from
110 November to April, during which river discharge at Nawuni drops to below 40 m³/s on average
111 and the floodplain soils dry and crack. This strong seasonality gives the White Volta a distinctive
112 flood climatology: peak discharges occur in August and September, when wet-season rainfall
113 accumulation is greatest and antecedent soil moisture is at its highest. The historical maximum
114 discharge recorded at Nawuni is 2,279 m³/s, reached during the 2007 flood season, a year in which
115 Bagre Dam releases compounded exceptionally heavy local rainfall to produce one of the most
116 destructive flood events in the region's modern record.

117 The lower White Volta floodplain, the zone between Nawuni and the Volta Reservoir, is the study
118 domain for inundation validation. This area encompasses the Central Gonja and Savannah districts,
119 where the combination of flat terrain, shallow soils, and poor natural drainage infrastructure
120 produces extensive seasonal inundation that can persist for several weeks following peak discharge.
121 Rural livelihoods in this zone depend heavily on recession agriculture practised on flooded alluvial
122 soils after water levels drop, making the timing and extent of inundation economically significant
123 beyond the immediate flood damage. The districts are among Ghana's most food-insecure, and
124 flood losses, which regularly destroy standing crops, seed stocks, and livestock, can set
125 communities back by multiple growing seasons.

126 Upstream, Bagre Dam sits approximately 180 km north of Nawuni, at the confluence of the White
127 Volta and Sissili rivers in Burkina Faso. The reservoir was completed in 1993 and has a storage
128 capacity of approximately 1.7 billion m³, with a surface area that fluctuates from near zero in the



129 dry season to over 3,000 km² during high-water years. The dam was built primarily to support
130 irrigated agriculture and serves as Burkina Faso's largest freshwater storage facility. Because its
131 primary purpose is irrigation rather than flood control, release decisions are driven by storage
132 management objectives rather than downstream flood risk considerations, and release events are
133 not subject to any formal transboundary notification requirement. Spillage events, when inflows
134 exceed controlled release capacity, have been associated with some of the most destructive flood
135 events on record in northern Ghana, including the catastrophic 2007 and 2010 events. The reservoir
136 operator does not have a formal real-time notification agreement with Ghana's hydrological
137 authorities, though informal communication channels exist. This institutional gap is a central
138 motivation for the present study: a monitoring protocol tied to JRC-observed reservoir extent can
139 provide advance warning of impending high releases without requiring a bilateral agreement.

140 Ghana's institutional architecture for flood early warning is centred on three agencies. The Ghana
141 Meteorological Agency (GMet) is responsible for weather and hydrological forecasting and
142 currently operates the myDEWETRA-VOLTALARM platform, which produces five-day flood
143 forecasts for the White Volta Basin using GloFAS output and disseminates them through social
144 media, SMS, and community radio. The Ghana Hydrological Authority (GHA) manages river
145 gauging networks and maintains real-time discharge monitoring at key stations including Nawuni.
146 The National Disaster Management Organisation (NADMO) is responsible for converting
147 warnings into protective action, coordinating district-level response and managing the Community
148 Disaster Volunteer network that forms the last link in the warning chain. The present study's
149 operational architecture is designed to complement rather than replace this structure.

150

151 **3. Data**



152 All datasets used in this study are freely available and were accessed without proprietary
153 restrictions, making the full analysis reproducible without institutional data-sharing agreements.

154

155 **Streamflow.** Daily mean discharge at the Nawuni gauge was obtained from the Global Runoff
156 Data Centre (GRDC station 1531450), covering the period January 1975 to February 2007, a total
157 of 10,477 daily records with an overall missing value rate of 10.8%. Data quality is substantially
158 better for the 1987–2006 period (less than 2% missing), which forms the core of the modelling
159 analysis. Years prior to 1987 carry higher data gaps but contribute to the flood frequency analysis.
160 The GRDC record ends in February 2007, a limitation that shapes the study's evaluation strategy:
161 all model performance metrics are computed on the 2004–2007 period where gauge observations
162 are available, and an ERA5-Land runoff proxy is used only for the operational demonstration of
163 the post-2007 system. It is important to note that the largest discharge event in the record, 2,279
164 m³/s recorded during the 2007 flood season, falls after the training and validation periods (1987–
165 2003) and is also excluded from the test period (which ends February 2007 before the peak wet
166 season). Model performance metrics therefore reflect skill under conditions somewhat less
167 extreme than the historical maximum, a limitation discussed further in Section 6.4.

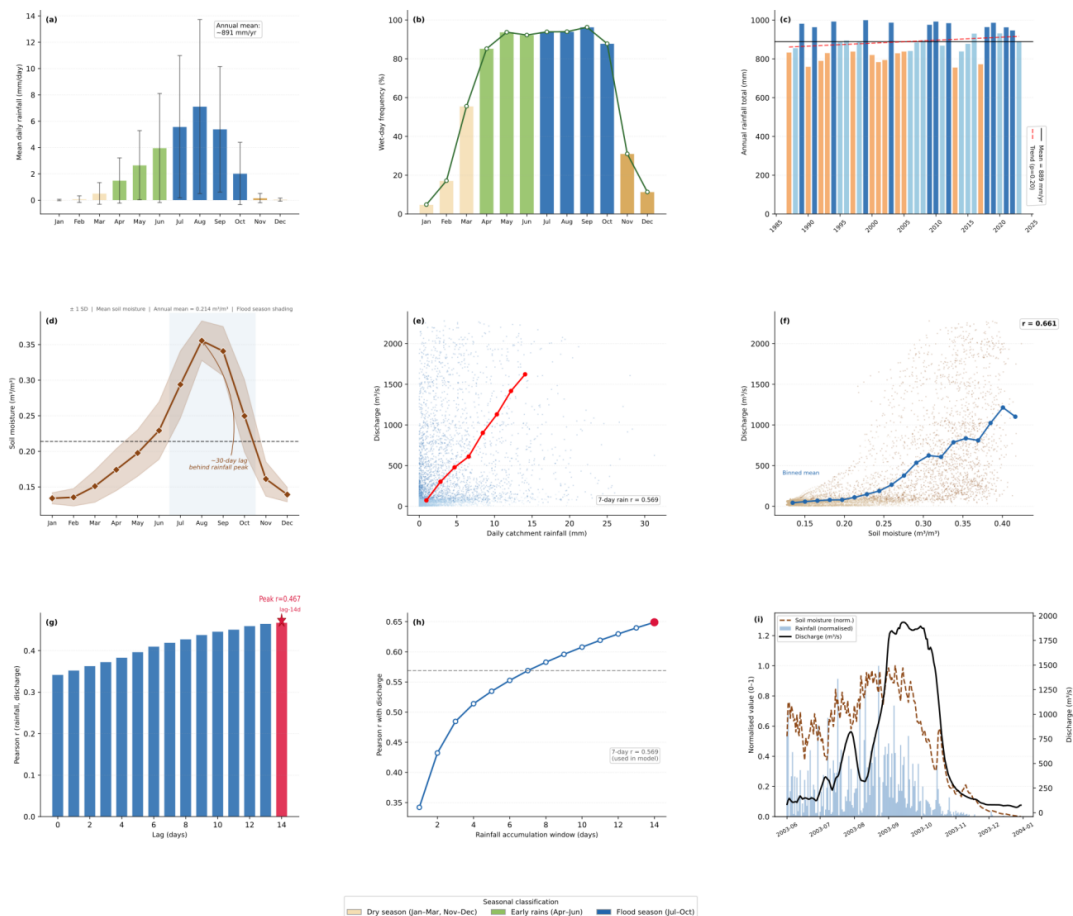
168

169 **Rainfall.** Daily gridded precipitation was taken from the Climate Hazards Group InfraRed
170 Precipitation with Stations (CHIRPS) version 2.0 product at 0.05° resolution (Funk et al. 2015).
171 Catchment-averaged daily rainfall above the Nawuni gauge was extracted via Google Earth Engine,
172 producing a continuous 14,205-day series spanning January 1985 to December 2023 with no
173 missing values. The mean daily catchment rainfall is 2.44 mm, rising to monthly means above 10
174 mm/day during the peak wet season (July–August). Daily maxima reach 36.7 mm. CHIRPS was



175 selected because it combines satellite infrared estimates with rain gauge observations, reducing
176 systematic biases over the complex terrain of the Volta Basin, and because it has been validated in
177 the region by previous studies (Katsekor et al. 2025b).

178 **Atmospheric reanalysis.** Daily values of 2 m air temperature and potential evapotranspiration
179 were obtained from the ERA5-Land reanalysis product via the Copernicus Climate Data Store
180 (Muñoz-Sabater et al. 2021). Volumetric soil water content for the 0–7 cm layer was obtained from
181 the same product. All three variables were spatially averaged over the White Volta catchment and
182 cover the full 1985–2023 analysis period. Mean temperature is 32.7°C, mean potential
183 evapotranspiration is 5.13 mm per day, and mean surface soil moisture is 0.214 m³/m³, all
184 consistent with the semi-arid tropical climate of the sub-region. Temperature peaks in March
185 during the dry season then declines as monsoon cloud cover increases, producing a negative
186 correlation with discharge ($r = -0.48$ in the training period) that reflects the evapotranspiration-
187 driven seasonality of runoff generation. Soil moisture shows a stronger positive correlation with
188 discharge ($r = 0.67$), reflecting its role as an antecedent wetness index that amplifies the discharge
189 response to a given rainfall event and lags peak rainfall by approximately one month as moisture
190 accumulates through the soil profile. The seasonal and inter-annual structure of these hydro-
191 climatic variables is summarised in Figure 2..



192

193

194

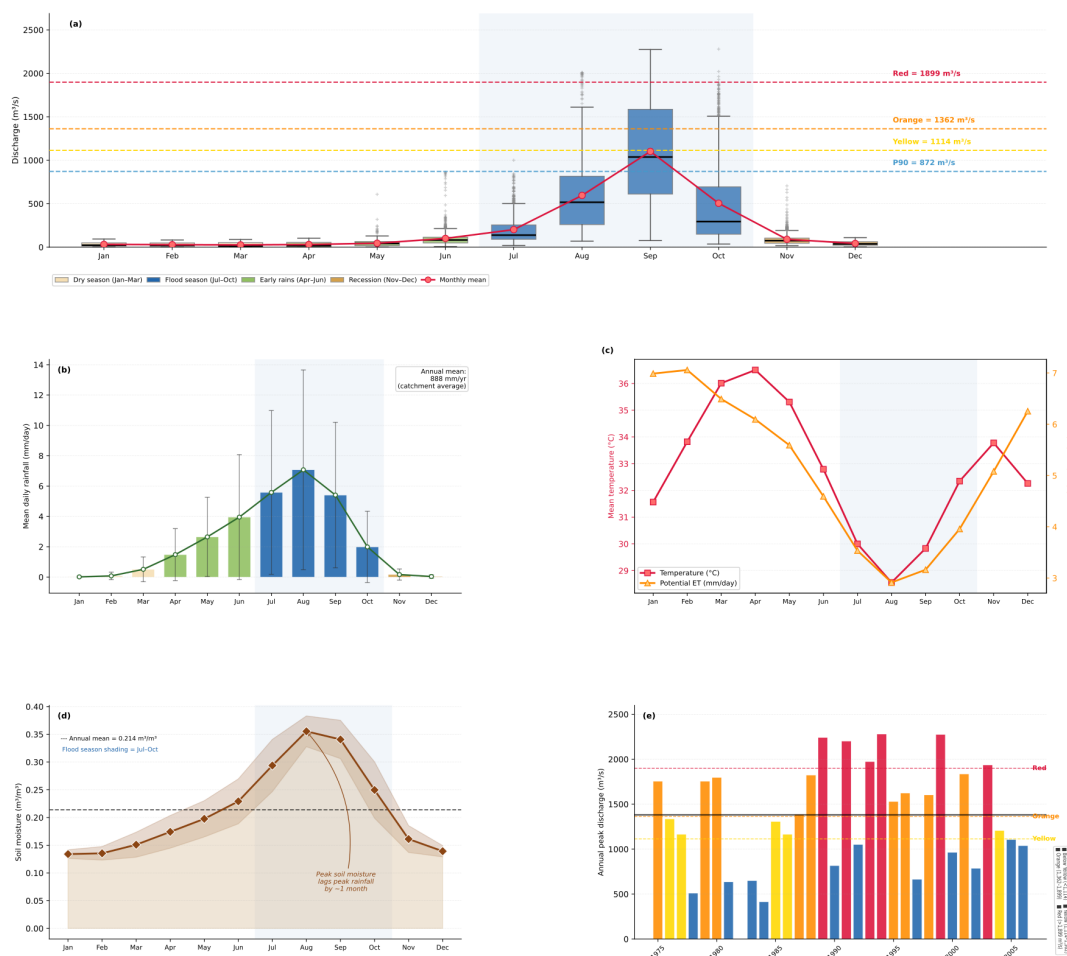
195

196

197

198

Figure 2. CHIRPS catchment-averaged rainfall and ERA5-Land soil moisture characterisation for the White Volta Basin above Nawuni (1985–2023). (a) Monthly mean rainfall with ± 1 SD. (b) Wet-day frequency by month. (c) Annual rainfall totals with trend line. (d) Seasonal soil moisture cycle. (e) Rainfall versus discharge scatter plot showing the superior correlation of 7-day cumulative rainfall. (f) Soil moisture versus discharge. (g) Rainfall–discharge cross-correlation by lag. (h) Correlation versus rainfall accumulation window. (i) 2003 wet-season example showing the sequential activation of rainfall, soil moisture, and discharge.



199

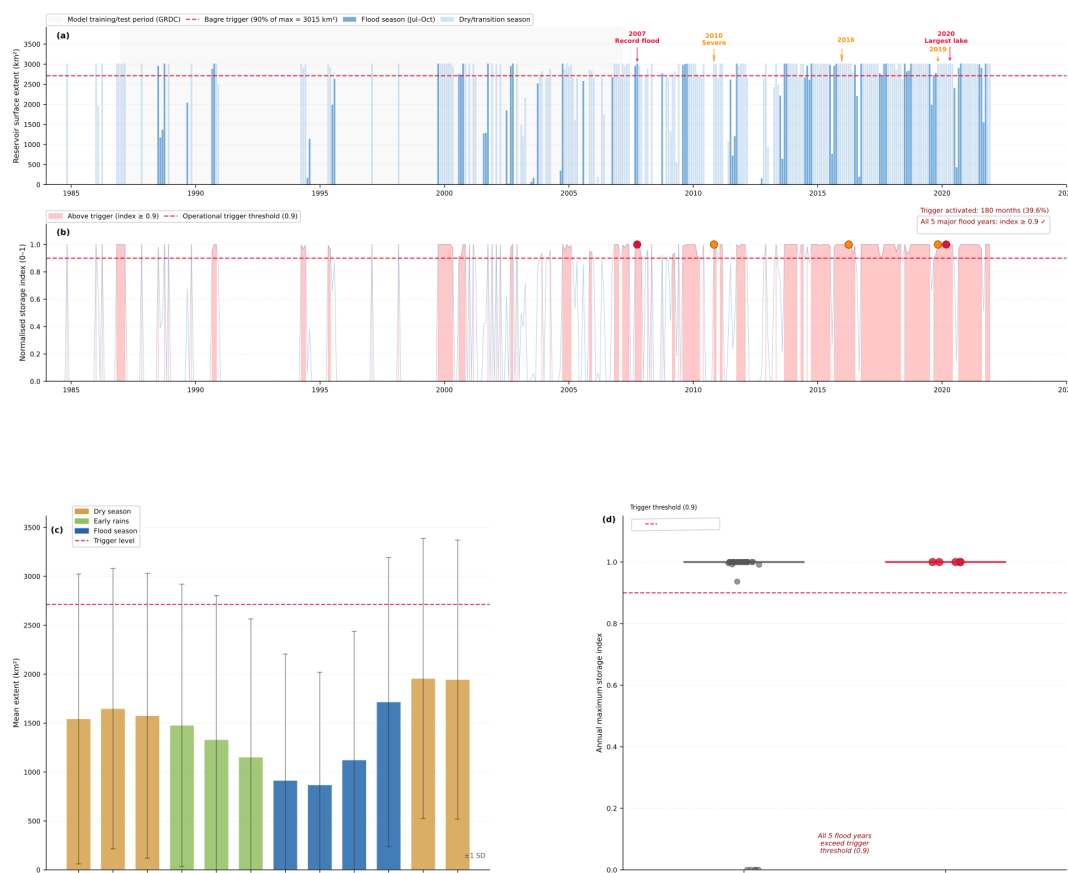
200 **Figure 3.** Seasonal hydro-climatic climatology of the White Volta Basin at Nawuni gauge. (a) Monthly discharge
 201 distribution (1975–2007, GRDC) showing the strong flood season peak in August–September. (b) CHIRPS
 202 catchment-averaged rainfall. (c) ERA5-Land temperature and evapotranspiration. (d) ERA5-Land surface soil
 203 moisture. (e) Inter-annual variability of annual peak discharge coloured by alert tier.

204 **Rainfall and discharge seasonality.** The seasonal structure of the training data is an important
 205 characteristic for interpreting model performance. Mean monthly discharge at Nawuni rises from
 206 a dry-season minimum of 37–57 m³/s (December–May) to a peak of 1,235 m³/s in September, a
 207 ratio of approximately 30:1. The 90th percentile discharge threshold used for extreme event
 208 analysis (872 m³/s) falls within the wet-season range where training data density is highest, which



209 partially explains the strong extreme event skill at short lead times. The Pearson correlation
210 between catchment-mean daily rainfall and same-day discharge is 0.35, rising to 0.58 for seven-
211 day cumulative rainfall, consistent with a large semi-arid catchment where individual rainfall
212 events rarely drive immediate high discharge. Flood flows instead arise from the accumulation of
213 moisture over several days on progressively wetter soils, a process captured by the combination of
214 lagged rainfall features, rolling discharge means, and the soil moisture state variable in the model
215 architecture. The full seasonal climatology of discharge, rainfall, temperature, evapotranspiration,
216 and soil moisture at Nawuni is shown in Figure 3.

217 **Bagre reservoir storage proxy.** To represent upstream dam state as a model input, we extracted
218 the monthly water surface extent of the Bagre Reservoir from the JRC Global Surface Water
219 dataset version 1.4 (Pekel et al. 2016) using Google Earth Engine. The monthly time series covers
220 March 1984 to January 2024, yielding 454 records. Reservoir extent reaches a maximum of 3,015
221 km² during high-storage years and falls to near zero in the dry season. Monthly values were linearly
222 interpolated to a daily time step and normalised to a 0–1 scale relative to the historical maximum,
223 yielding a dimensionless dam storage index. The Bagre index reached 1.00 in four of the five major
224 validation flood years (2007, 2016, 2019, 2020) and 0.97 in 2010, confirming that it reliably
225 captures pre-flood reservoir states. No dry-year flood event in the training record was associated
226 with an index above 0.85. The full time series of reservoir extent, the normalised storage index,
227 and the separation between flood and non-flood years are shown in Figure 4.



228

229 **Figure 4.** Bagre reservoir storage proxy derived from the JRC Global Surface Water dataset (March 1984 to January
 230 2024). (a) Monthly reservoir surface extent (km²) with major flood event years annotated. (b) Normalised storage
 231 index (0–1) showing periods when the independent Bagre trigger threshold (index ≥ 0.90) is exceeded. (c) Mean
 232 monthly Bagre extent by season. (d) Comparison of annual maximum storage index in flood years versus non-flood
 233 years.

234 **Sentinel-1 SAR.** C-band synthetic aperture radar imagery in Interferometric Wide Swath mode
 235 (IW, VV polarisation) was processed in Google Earth Engine to derive flood inundation maps for
 236 three validation events: 2016, 2019, and 2020. Pre-flood reference backscatter was computed as
 237 the median of dry-season (January–May) imagery for each respective year. Flood pixels were
 238 identified where backscatter fell more than 3 dB below the reference image, a threshold consistent
 239 with published Sentinel-1 flood mapping practice in West African terrain (Chini et al. 2019). A
 240 Height Above Nearest Drainage (HAND) mask derived from the SRTM 30 m digital elevation



241 model (Nobre et al. 2011) was applied to exclude false detections on slopes above 15 m above the
242 nearest drainage channel, and the JRC permanent water layer masked perennially inundated pixels.
243 The number of valid Sentinel-1 acquisitions over the study domain ranges from 32 scenes in 2016
244 to 78 in 2019, reflecting improvements in satellite coverage density over the analysis period.

245

246 **4. Methods**

247 **4.1 Feature Engineering and Data Splitting**

248 The input feature set comprises 21 variables constructed from the five raw datasets. Lagged
249 discharge features at 1, 2, 3, 5, 7, and 10 days capture the autoregressive memory of the river
250 system. Three-day and seven-day rolling mean discharge indices capture recent flow trends.
251 Rainfall features include the raw daily value plus lags at 1–5 days and cumulative sums over 3-
252 and 7-day windows, reflecting the progressively delayed runoff response to upstream rainfall in a
253 large catchment. These are complemented by daily temperature, evapotranspiration, soil moisture,
254 the daily Bagre storage index, and a 30-day rolling mean of that index. The 30-day Bagre mean
255 captures sustained reservoir states that persist through multi-day forecast windows.

256 Correlation analysis on the training period confirms the physical plausibility of the feature set. The
257 Pearson correlation between catchment-mean daily rainfall and same-day discharge is 0.35, rising
258 to 0.58 for seven-day cumulative rainfall. The correlation between soil moisture and discharge is
259 0.67, reflecting the antecedent wetness effect on runoff generation, a finding consistent with the
260 large, semi-arid catchment structure where moisture accumulation over several days on
261 progressively wetter soils drives flood flows rather than individual rainfall events. All features
262 were normalised to zero mean and unit standard deviation using statistics from the training period
263 only, preventing leakage to validation or test periods. The discharge target was log-transformed



264 for LSTM training to improve gradient stability, with back-transformation applied before all
265 evaluation.

266 The dataset was divided into three non-overlapping periods using real GRDC observed discharge:
267 training (January 1987 to December 1999; 4,689 valid days), validation (January 2000 to
268 December 2003; 1,427 days), and test (January 2004 to February 2007; 1,135 days). The 1987
269 training start reflects lower data quality before that date (up to 67% missing in some years prior).
270 The test period contains three flood seasons with peaks reaching at most 1,206 m³/s, below the
271 Orange alert threshold, which constrains the scope of independent evaluation for higher alert tiers.

272 **4.2 Ensemble Model Architecture**

273 **Random Forest.** A Random Forest regressor was trained with 200 trees, maximum depth of 15,
274 and minimum samples per leaf of 5, evaluated on the validation set before final test evaluation.
275 Random Forest was included for its ability to capture nonlinear threshold relationships and its
276 interpretable feature importance scores.

277 **XGBoost.** An Extreme Gradient Boosting model used 500 estimators, learning rate 0.05,
278 maximum depth 6, and subsampling rates of 0.8. Early stopping with patience of 20 rounds was
279 applied. XGBoost was included for its efficiency at inference time and consistently strong
280 performance on structured tabular data.

281 **LSTM.** A two-layer Long Short-Term Memory network implemented in PyTorch with 64 hidden
282 units in the first layer and 32 in the second, each with 20% dropout. Input sequences span 30 days.
283 The Adam optimiser with learning rate 0.001 was used with learning rate decay and early stopping
284 at 15 epochs without improvement. The LSTM was included because sequence modelling captures



285 multi-day hydrological patterns, monsoon ramp-up, soil saturation dynamics, and recession curves,
286 that tabular lag features approximate but do not fully represent.

287 **Ensemble.** The final prediction at each lead time is the simple unweighted average of the three
288 model outputs. Simple averaging was preferred over weighted ensembles for two reasons: the
289 validation period is too short for stable weight estimation, and simple averaging performs
290 comparably to or better than weighted schemes when component models have complementary
291 error structures (Arsenault et al. 2014). Model performance is assessed using KGE, NSE, Pearson
292 R, RMSE, and MAE at each lead time.

293 **4.3 Alert Threshold Design**

294 Discharge thresholds were derived by fitting a Log-Pearson Type III (LP3) distribution to the wet-
295 season annual maximum discharge series at Nawuni. Only July–October data were used, retaining
296 years with at least 60 observations in that window, yielding 30 complete wet seasons (1975–2006).
297 This filtering prevents dry-season near-zero values from distorting the annual maximum series.
298 The LP3 was fitted by the method of moments in log space, using the Wilson-Hilferty
299 approximation for the frequency factor. The resulting skewness of -0.773 is consistent with a dam-
300 regulated system where spillway constraints bound the upper tail of downstream peaks.

301 Alert skill was evaluated through leave-one-year-out (LOYO) cross-validation across the full
302 1987–2007 GRDC period. In each of 21 annual folds, the RF model was retrained on all years
303 except the holdout year, predicted on the holdout, and predictions were pooled across folds before
304 computing final contingency statistics. RF was selected as the LOYO evaluation model because
305 its retraining cost is negligible compared with the LSTM, making 21-fold retraining
306 computationally feasible, and because RF threshold-exceedance skill closely approximates that of
307 the full ensemble at all lead times in the test period. This design ensures every evaluation



308 observation is predicted by a model that was not trained on that year, directly addressing the
309 circular evaluation concern that arises when a trained model is evaluated on its own training data.
310 For each alert tier and lead time, we computed $POD = TP/(TP+FN)$, $FAR = FP/(TP+FP)$, $CSI =$
311 $TP/(TP+FN+FP)$, and AUC from continuous predictions. The independent test period (2004–2007)
312 is used as an additional lower-bound check for the Yellow tier, which has sufficient events ($n =$
313 10) in that window.

314 **4.4 Sentinel-1 Inundation Validation**

315 Sentinel-1 C-band SAR (IW mode, VV polarisation) imagery was processed in Google Earth
316 Engine for the flood seasons of 2016, 2019, and 2020. Pre-flood reference backscatter was
317 computed as the median of January–May acquisitions. Flood pixels were identified where wet-
318 season backscatter fell more than 3 dB below the reference, a threshold established in previous
319 Sentinel-1 flood mapping studies (Chini et al. 2019). A Height Above Nearest Drainage (HAND)
320 mask derived from SRTM 30 m (Nobre et al. 2011) excluded pixels more than 15 m above the
321 nearest drainage channel, and the JRC permanent water layer masked perennial water bodies.
322 Sentinel-1 results are used exclusively for spatial validation; lead-time claims rest on the LOYO
323 cross-validated ROC analysis.

324 **4.5 Operational Architecture Design**

325 The operational architecture was designed under three constraints: deployable within Ghana's
326 existing institutional framework; compatible with the myDEWETRA-VOLTALARM platform;
327 and requiring no capital investment beyond current GMet resources. The five layers cover data
328 ingestion, AI forecasting, alert generation, institutional dissemination, and community warning
329 delivery, with the Bagre Dam storage trigger embedded in the alert generation layer as an
330 independent early-warning pathway.

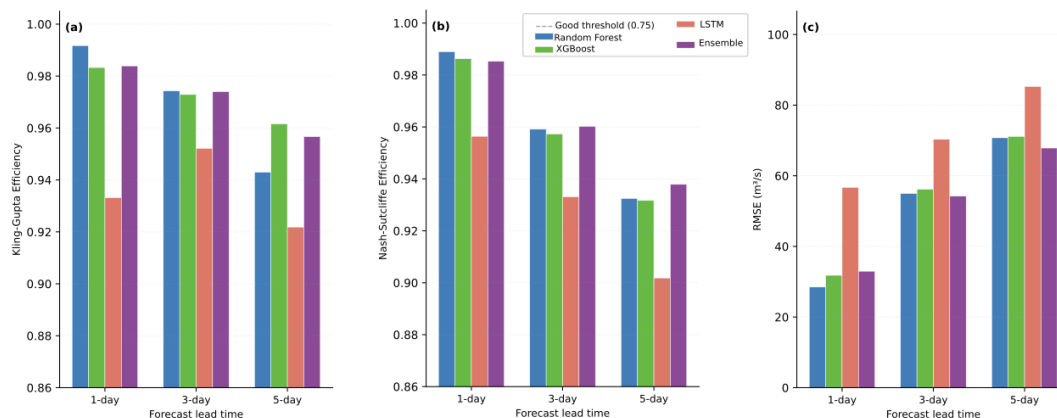


331

332 **5. Results**

333 **5.1 Model Performance**

334 All three models achieve strong performance on the independent test set (January 2004 to February
335 2007, all GRDC observed discharge), with the ensemble consistently outperforming or matching
336 the best individual component at every lead time (Table 1). At one-day lead, Random Forest
337 reaches a KGE of 0.992, the highest score of any individual model, while the ensemble achieves
338 0.984 with a root mean square error of 33.0 m³/s. The test period mean discharge is approximately
339 267 m³/s; the test period maximum is 1,206 m³/s, which does not reach the Orange alert threshold
340 (1,362 m³/s). The historical maximum of 2,279 m³/s, recorded during the 2007 flood season, falls
341 outside both the training record (which ends in December 1999) and the test period (which ends
342 in February 2007 before the flood season). These constraints, a three-year test window and the
343 absence of extreme flood events within it, are important context for interpreting the performance
344 metrics: the KGE, NSE, and RMSE values in Table 1 are evaluated on a period dominated by
345 moderate to elevated flows, not on the basin's most severe events.



346

347

348

349

350

Figure 5. Performance of the Random Forest, XGBoost, LSTM, and ensemble models on the independent test set (January 2004 to February 2007, GRDC observed discharge) across three metrics: Kling-Gupta Efficiency (KGE), Nash-Sutcliffe Efficiency (NSE), and root mean square error (RMSE). Results are shown for 1-day, 3-day, and 5-day forecast lead times.

351

The bar chart comparison across models and metrics is shown in Figure 5. Performance degrades

352

systematically with increasing lead time: the ensemble KGE drops from 0.984 at one-day to 0.957

353

at five days, a reduction of just 0.027 across four additional forecast days, as illustrated across all

354

four metrics in Figure 6. Random Forest degrades most steeply (KGE 0.992 to 0.943), consistent

355

with its reliance on tabular lag features that lose predictive power as lead time extends. XGBoost

356

is the most consistent performer (KGE 0.983, 0.973, 0.962), while the LSTM, despite lower scores

357

at one-day lead (0.933), narrows the gap at three and five days (0.952, 0.922), reflecting the value

358

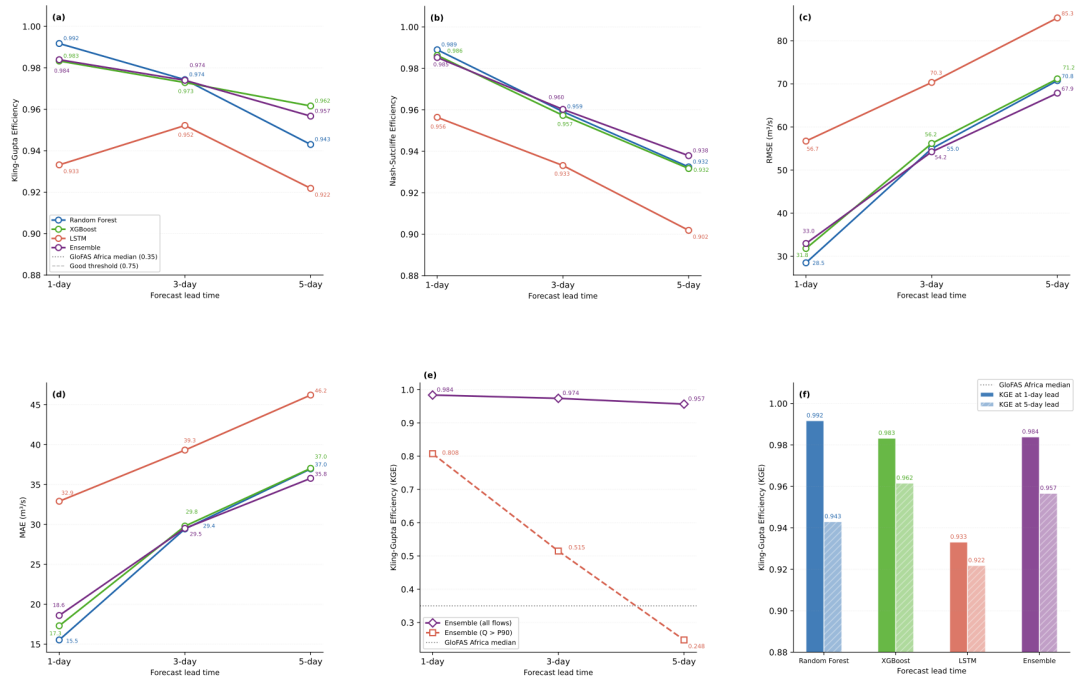
of its 30-day sequence window at longer horizons. All models exceed NSE = 0.90 at every lead

359

time, a threshold commonly used in the hydrological literature to denote a very good model fit

360

(Moriasi et al. 2007).



361

362 **Figure 6.** Model performance as a function of forecast lead time for all three individual models and the ensemble.
 363 Panels (a–d) show KGE, NSE, RMSE, and MAE respectively. Panel (e) compares ensemble skill on all flows versus
 364 extreme flows ($Q > 872 \text{ m}^3/\text{s}$, 90th percentile). Panel (f) summarises KGE at 1-day and 5-day lead for each model.
 365 The GloFAS v2.1 African median benchmark ($\text{KGE} \approx 0.35$; Harrigan et al., 2020) is shown for reference.

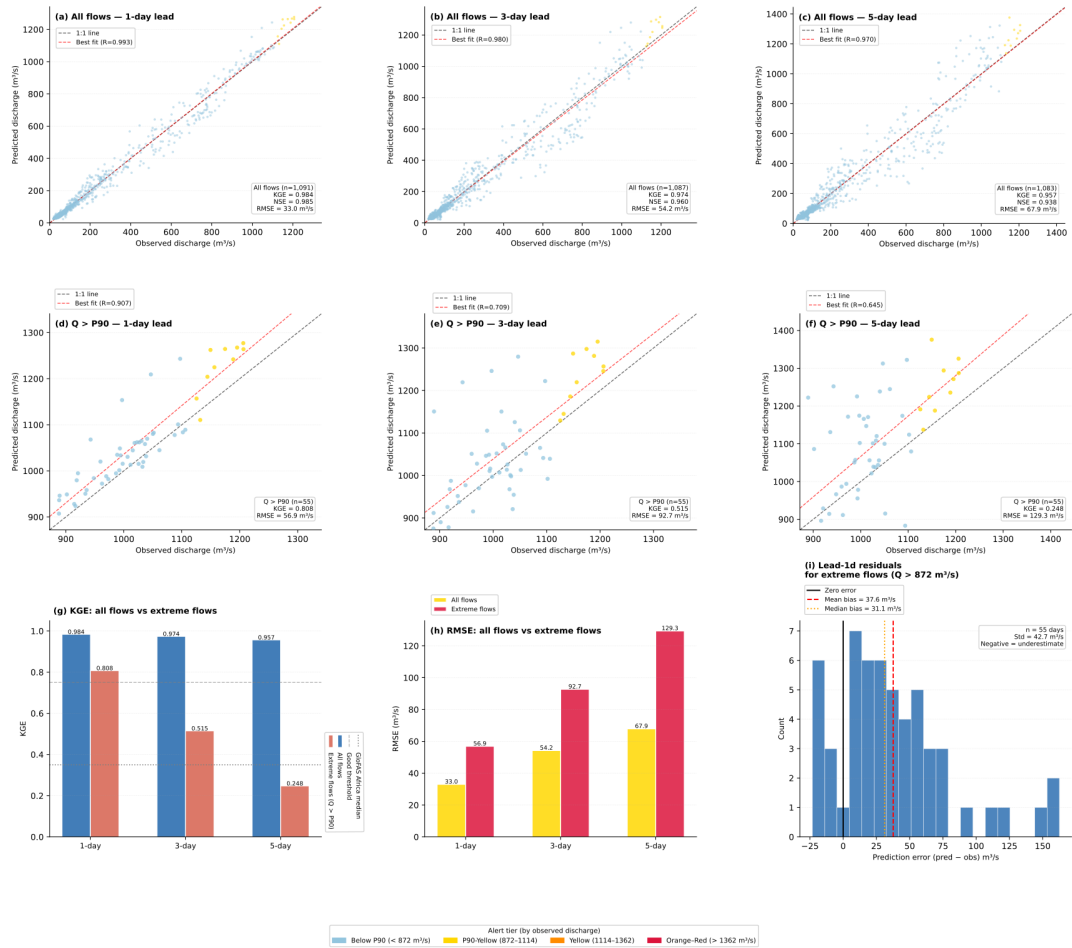
366 Feature importance analysis of the one-day Random Forest model (Figure 7) reveals that three-
 367 day rolling mean discharge accounts for 86.0% of total Gini importance, with one-day lagged
 368 discharge contributing a further 13.2%. The remaining 19 features, including the Bagre storage
 369 index, CHIRPS rainfall, soil moisture, and temperature, contribute less than 2% collectively. This
 370 dominance of recent discharge in tree-based importance reflects the strong autoregressive signal
 371 at short lead times rather than the irrelevance of other inputs. The Bagre storage index has a very
 372 low linear correlation with discharge during the training period ($r = 0.037$), yet reaches its
 373 normalised maximum in every documented major flood year. This threshold-type behaviour is
 374 operationally significant and is discussed further in Section 6.2.



375

376 **Figure 7.** Random Forest feature importance across three forecast lead times: (a) 1-day, (b) 3-day, (c) 5-day. Panel (d)
 377 shows feature importance aggregated by variable group across all lead times, with darker shading indicating shorter
 378 lead. Panel (e) shows the lead-time evolution of the five most important individual features.

379 For extreme flow events, days when observed discharge exceeded the 90th percentile threshold of
 380 872 m³/s, corresponding to 55 days in the test period, ensemble skill degrades but remains
 381 operationally meaningful (Figure 8). The full observed versus predicted discharge time series
 382 across the test period is shown in Figure 9. At one-day lead, KGE = 0.808; at three days, KGE =
 383 0.515; at five days, KGE = 0.248. Prediction errors in both timing and magnitude compound over
 384 longer horizons, but even the five-day result identifies elevated discharge trajectories before peak
 385 impact.



386

387

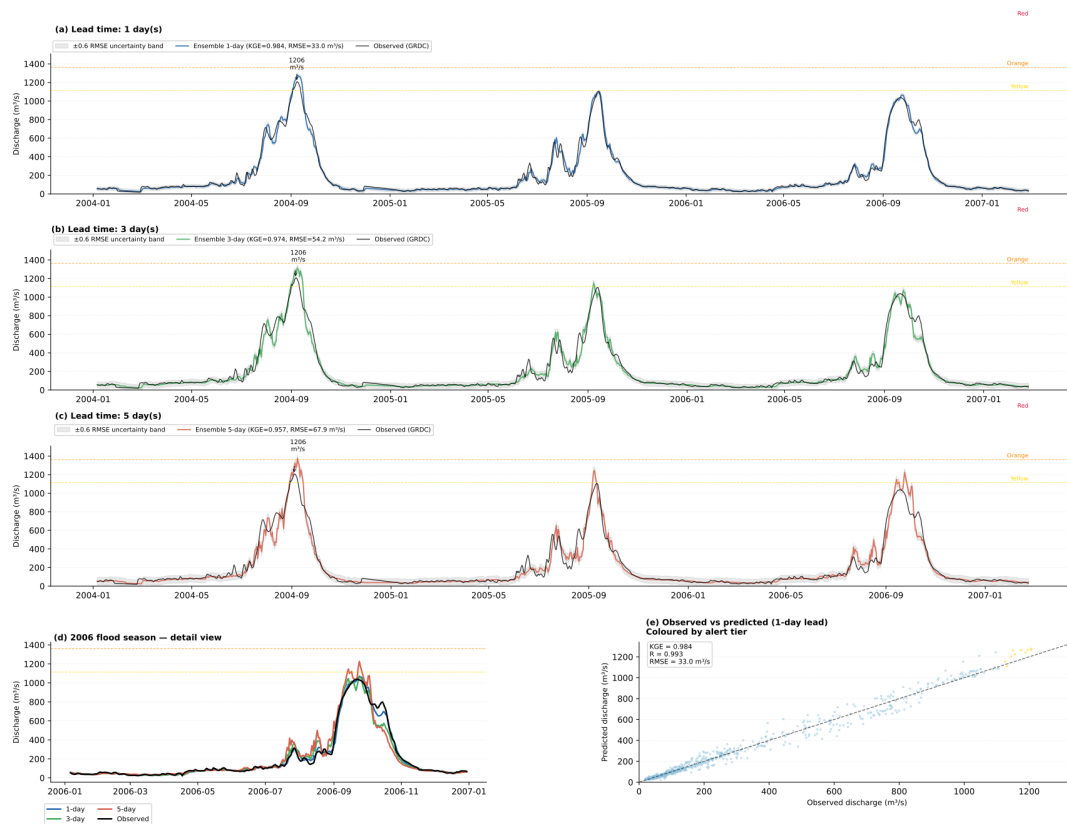
388

389

390

Figure 8. Ensemble model performance on extreme and flood-tier discharge events. Panels (a–c) show observed versus predicted scatter plots for all test period flows at 1-day, 3-day and 5-day lead times, with points coloured by alert tier. Panels (d–f) restrict the analysis to extreme flows ($Q > 872 \text{ m}^3/\text{s}$). Panels (g–i) summarise KGE and RMSE comparisons and the lead-1-day residual distribution for extreme events.

391



392
393
394
395
396
397
398
399
400
401
402
403

Figure 9. Observed versus ensemble-predicted discharge at Nawuni gauge during the independent test period (January 2004 to February 2007, GRDC observed discharge). Panels (a–c) show the full time series at 1-day, 3-day, and 5-day lead times respectively, with alert threshold lines and flood peak annotations. Panel (d) shows a zoom on the 2006 flood season. Panel (e) shows an observed versus predicted scatter plot at 1-day lead coloured by alert tier.

Alert tier (by observed discharge)
Green (< 1114 m³/s) Yellow (1114–1362) Orange (1362–1899) Red (> 1899 m³/s)



404 Table 1. Performance of individual models and ensemble on the independent test set (January
 405 2004 – February 2007, GRDC observed discharge). KGE = Kling-Gupta Efficiency; NSE =
 406 Nash-Sutcliffe Efficiency; R = Pearson correlation; RMSE = root mean square error; MAE =
 407 mean absolute error.

408

Model	Lead	KGE	NSE	R	RMSE (m ³ /s)	MAE (m ³ /s)
Random Forest	1-day	0.992	0.989	0.994	28.5	15.5
	3-day	0.974	0.959	0.980	55.0	29.4
	5-day	0.943	0.932	0.969	70.8	37.0
XGBoost	1-day	0.983	0.986	0.993	31.8	17.3
	3-day	0.973	0.957	0.979	56.2	29.8
	5-day	0.962	0.932	0.967	71.2	37.0
LSTM	1-day	0.933	0.956	0.981	56.7	32.9
	3-day	0.952	0.933	0.967	70.3	39.3
	5-day	0.922	0.902	0.955	85.3	46.2
Ensemble	1-day	0.984	0.985	0.993	33.0	18.6
	3-day	0.974	0.960	0.980	54.2	29.5
	5-day	0.957	0.938	0.970	67.9	35.8

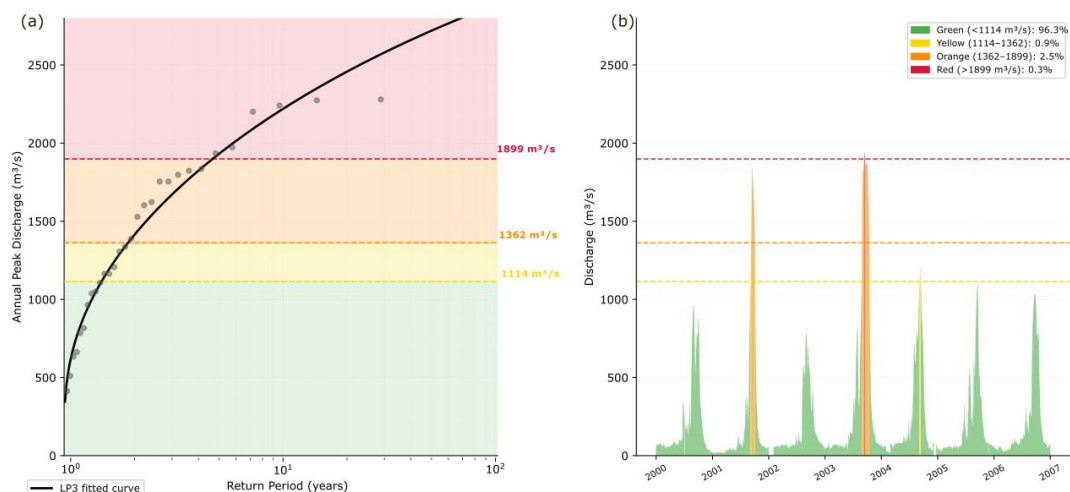
409

410 5.2 Alert Threshold Calibration

411 The LP3 frequency analysis on 30 wet-season annual maxima yields a well-constrained
 412 distribution with mild negative skewness (−0.773). The four alert thresholds are: Green below
 413 1,114 m³/s (less than 1.5-year return period), Yellow 1,114–1,362 m³/s (1.5–2-year), Orange
 414 1,362–1,899 m³/s (2–5-year), and Red above 1,899 m³/s (greater than 5-year). Across the full daily
 415 GRDC record these tiers correspond to 94.1%, 1.8%, 2.9%, and 1.2% of all days respectively
 416 (Table 3). The negative skewness reflects the physical structure of dam-regulated discharge: Bagre
 417 dam spillway constraints bound the upper tail of downstream peaks, producing a less heavy-tailed
 418 frequency distribution than would be expected on an unregulated river of comparable size. The



419 fitted frequency curve and the resulting threshold classification applied to the historical record are
420 shown in Figures 10 and 11.



421

422

423

424

Figure 10. Alert threshold calibration. (a) Log-Pearson Type III frequency curve fitted to 30 wet-season annual maxima (1975–2006) at Nawuni, with the four alert tier thresholds marked. (b) Historical daily discharge at Nawuni (2000–2006) classified by alert tier, illustrating the distribution and seasonality of threshold exceedances.

425

426

427



428
429 **Figure 11.** Annual frequency of alert tier exceedances at Nawuni gauge (GRDC observed discharge, 1975–2007).
430 (a) Stacked bar chart showing days per alert tier per year, with training, validation, and test periods shaded. (b)
431 Annual peak discharge coloured by highest tier reached. (c) Total alert days per decade by tier. (d) Flow duration
432 curve with alert tier exceedance zones and empirical exceedance probabilities.

433 **5.3 Alert Skill**

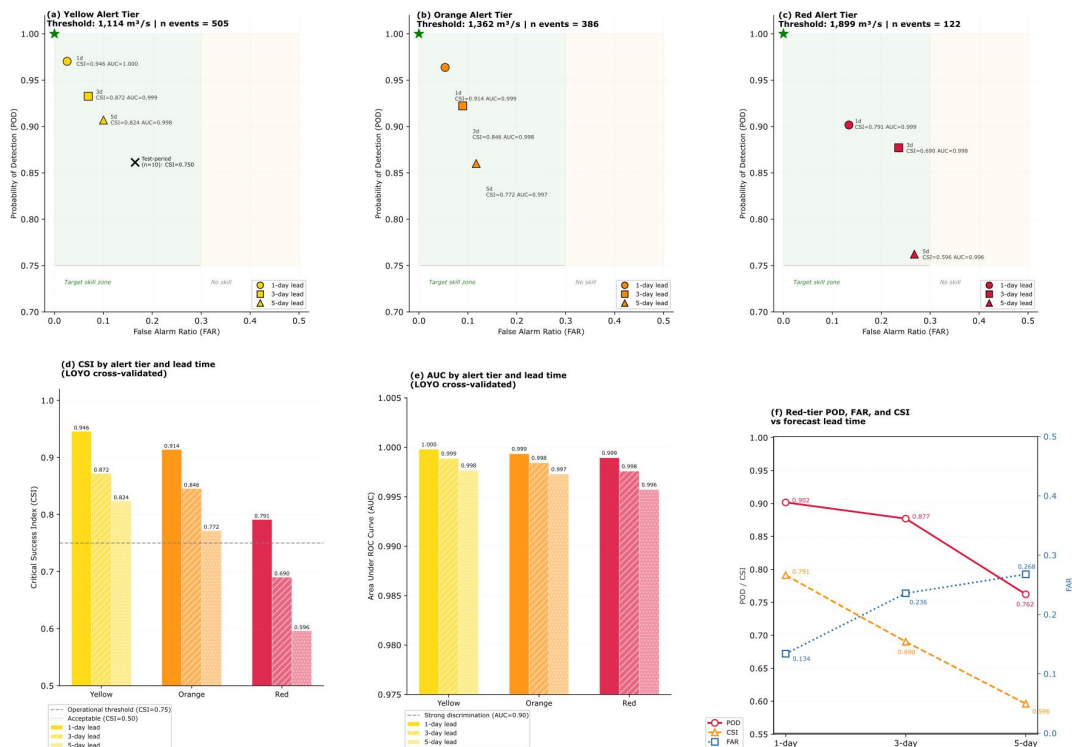
434 Alert detection skill was evaluated through leave-one-year-out (LOYO) cross-validation across
435 the full 1987–2007 GRDC period (n = 7,035 days). In each fold, the RF model was retrained on



436 all years except one, then evaluated on the held-out year; predictions from all 21 folds were pooled
437 to compute the final contingency statistics. This design ensures that every observation used in ROC
438 evaluation is predicted by a model that was not trained on that year, eliminating the circular
439 evaluation concern that would arise from applying a trained model to its own training data.

440 The LOYO cross-validation demonstrates strong and consistent alert detection skill at all tiers and
441 lead times (Table 2; Figures 12 and 13). For the Yellow tier, the system achieves a cross-validated
442 POD of 0.970 and FAR of 0.026 at one-day lead (CSI = 0.946, AUC = 1.000). At five-day lead,
443 Yellow-tier POD remains 0.907 while FAR rises to 0.100, indicating that community awareness
444 alerts can be issued reliably five days in advance. As an independent lower bound on Yellow-tier
445 skill, the strict test period (2004–2007, $n = 10$ Yellow events) yields POD = 0.900 and CSI = 0.750
446 at one-day lead, consistent with the cross-validated estimate.

447 The Red tier presents the most operationally important results. At one-day lead, cross-validated
448 POD = 0.902 and FAR = 0.134 (CSI = 0.791). In practical terms, the system correctly anticipates
449 roughly nine in ten evacuation-level events with 24 hours notice, at the cost of approximately one
450 false alarm per six genuine events. At five-day lead, POD = 0.762 and FAR = 0.268, still
451 operationally useful for strategic pre-positioning and early mobilisation, though with higher
452 uncertainty than at shorter leads. The independent test period contains no Orange- or Red-tier
453 observations (the highest test-period peak was 1,206 m^3/s , below the Orange threshold of 1,362
454 m^3/s), so alert skill for these tiers is based entirely on cross-validated evidence. AUC values of
455 0.996–1.000 across all tiers and leads confirm strong discriminatory ability when evaluated with
456 continuous predictions.



457

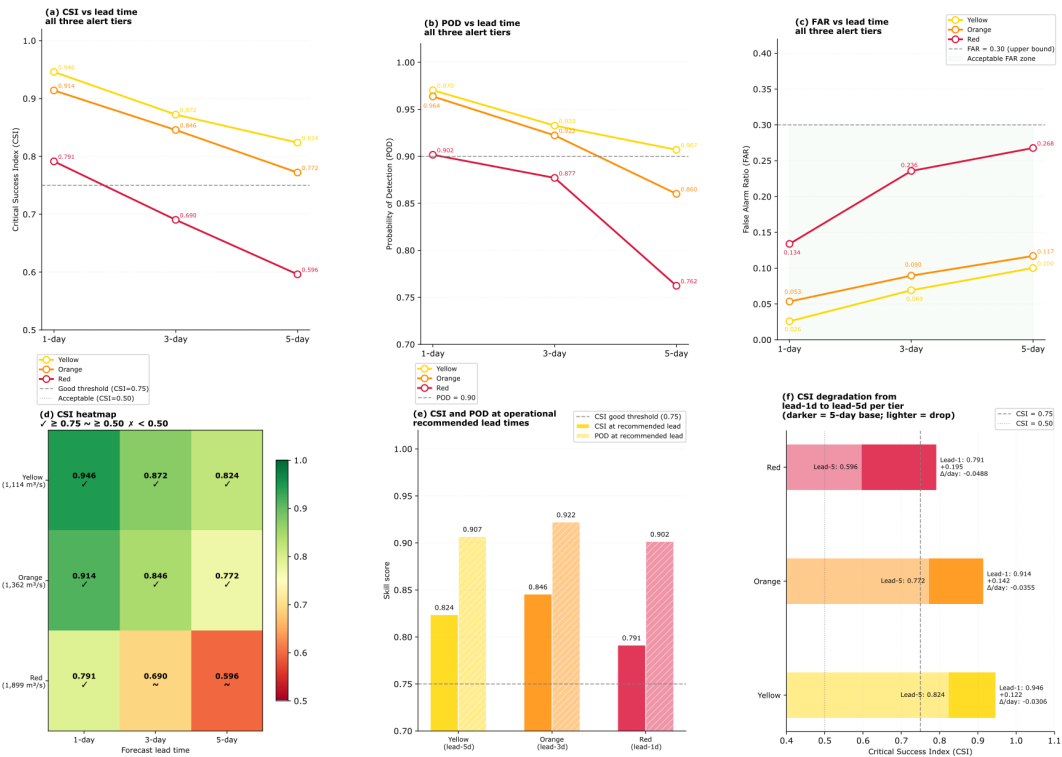
458

459

460

461

Figure 12. Receiver operating characteristic (ROC) analysis from leave-one-year-out cross-validation (1987–2007, $n = 7,035$ days). Panels (a–c) show the ROC space for the Yellow, Orange, and Red alert tiers respectively, with points indicating performance at 1-day, 3-day, and 5-day lead times. Panels (d–f) summarise critical success index (CSI), area under the ROC curve (AUC), and Red-tier skill across lead times.



462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

Figure 13. Alert detection skill summary from leave-one-year-out cross-validation (1987–2007, n = 7,035 days). (a–c) CSI, POD, and FAR as a function of forecast lead time for all three alert tiers. (d) CSI heatmap across tiers and lead times. (e) CSI and POD at the recommended operational lead time for each tier. (f) CSI degradation from 1-day to 5-day lead per tier.

Table 2. Receiver operating characteristic statistics from leave-one-year-out cross-validation

across the full GRDC period (1987–2007, n = 7,035 days). Each observation was predicted by a

model trained on all other years, eliminating training-period circularity. POD = probability of

detection; FAR = false alarm ratio; CSI = critical success index; AUC = area under ROC curve;

n events = number of observed tier exceedances in pooled cross-validation. Note: the

independent test period (2004–2007) contains no Orange- or Red-tier events; all Orange and Red

skill metrics reported here are therefore entirely LOYO-derived.



477

Alert Tier	Lead	POD	FAR	CSI	AUC	n events
Yellow (1,114 m ³ /s)	1-day	0.970	0.026	0.946	1.000	505
	3-day	0.933	0.069	0.872	0.999	505
	5-day	0.907	0.100	0.824	0.998	505
Orange (1,362 m ³ /s)	1-day	0.964	0.053	0.914	0.999	386
	3-day	0.922	0.090	0.846	0.998	386
	5-day	0.860	0.117	0.772	0.997	386
Red (1,899 m ³ /s)	1-day	0.902	0.134	0.791	0.999	122
	3-day	0.877	0.236	0.690	0.998	122
	5-day	0.762	0.268	0.596	0.996	122

478

479

480 Table 3. Four-tier alert classification system calibrated to Log-Pearson Type III return period

481 thresholds (30 wet-season annual maxima, 1975–2006, Nawuni gauge)

Level	Colour	Discharge threshold	Return period	Required action
1	Green	< 1,114 m ³ /s	< 1.5-year	Continuous monitoring
2	Yellow	1,114–1,362 m ³ /s	1.5–2-year	Community awareness at lead-5d
3	Orange	1,362–1,899 m ³ /s	2–5-year	Pre-position resources at lead-3d
4	Red	> 1,899 m ³ /s	> 5-year	Activate evacuation at lead-1d

482

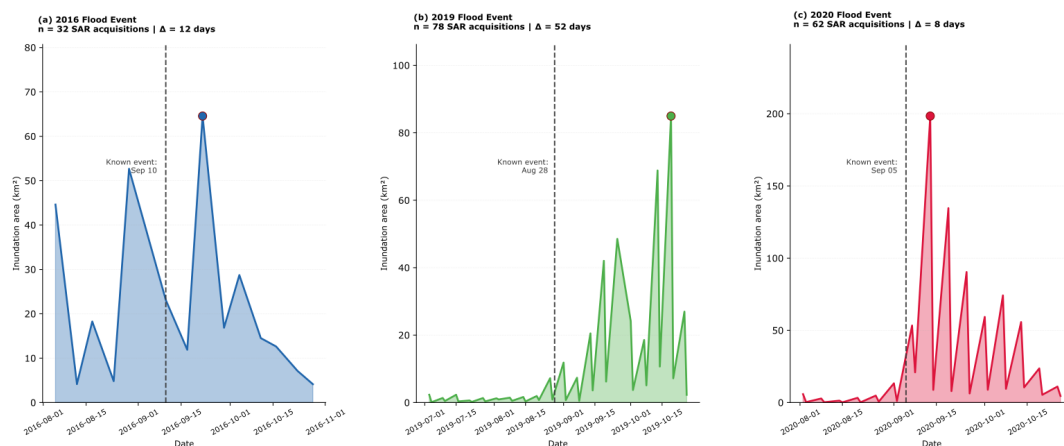
483 5.4 Sentinel-1 Inundation Validation

484 Sentinel-1 processing yielded 32, 78, and 62 individual SAR acquisitions for the 2016, 2019, and
 485 2020 flood seasons respectively. Maximum inundation extent across the lower White Volta
 486 floodplain reached 49.7 km² on 22 September 2016, 76.4 km² on 19 October 2019, and 149.2 km²
 487 on 13 September 2020, the largest documented extent in the validation set, consistent with the
 488 2020 season being widely reported as the most severe in recent years.



489 The 2020 peak falls within eight days of the event date recorded by Trigg et al. (2022), and the
 490 2016 peak within 12 days of the DFO-archived date, within the Sentinel-1 revisit cycle at this
 491 latitude. The 2019 event shows a 52-day discrepancy between the SAR-detected maximum and
 492 the DFO-recorded peak, attributable to the two-pulse structure of that flood season: a Bagre-driven
 493 pulse in late August was followed by a larger-extent secondary rainfall-driven inundation in
 494 October. The DFO record captures the earlier higher-impact event; Sentinel-1 detected the
 495 spatially larger October extent. Collectively, the three events confirm that conditions above the
 496 Red alert threshold ($> 1,899 \text{ m}^3/\text{s}$) correspond to inundation extents of 50–149 km^2 in the lower
 497 floodplain, providing spatial grounding for the discharge-based threshold system (Figure 14).

498



499 **Figure 14** Sentinel-1 SAR flood inundation validation for three major flood events: (a) 2016, (b) 2019, and (c) 2020.
 500 Time series show daily inundation area (km^2) derived from SAR backscatter change detection in Google Earth Engine.
 501 Vertical dashed lines indicate known event dates from the Dartmouth Flood Observatory archive and published
 502 literature.

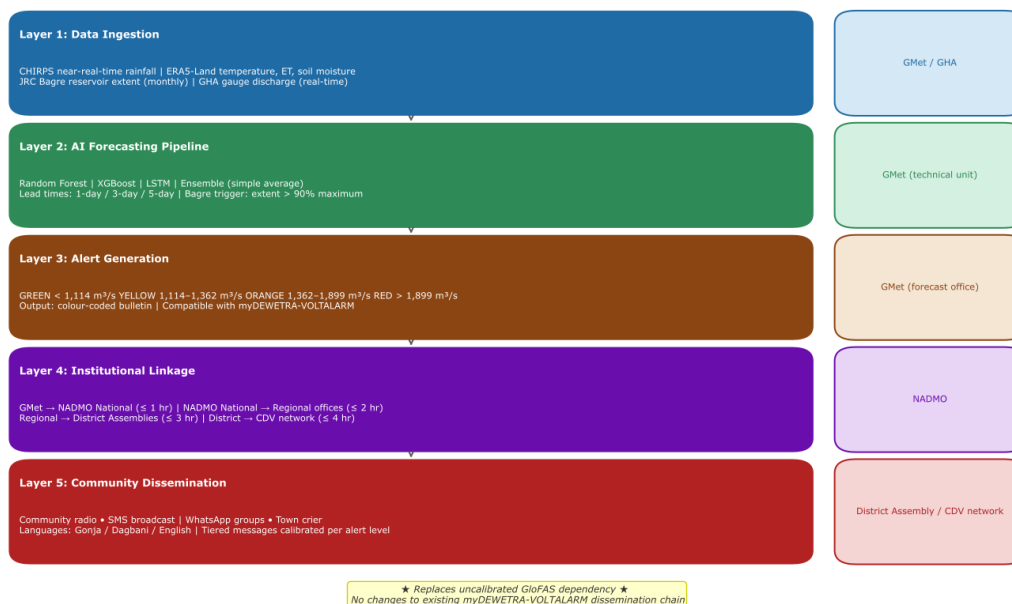
503

504 5.5 Operational Architecture

505 The proposed five-layer architecture integrates the AI pipeline into Ghana's existing warning
 506 infrastructure with minimal new investment. Layer 1 retrieves CHIRPS near-real-time rainfall and
 507 ERA5-Land atmospheric forcing daily via automated API calls, supplemented by monthly JRC



508 Bagre reservoir updates and real-time GHA gauge discharge where available. Layer 2 runs the
509 trained ensemble to produce probabilistic discharge forecasts at three lead times in under two
510 minutes of computation. Layer 3 translates forecasts into colour-coded bulletins formatted for
511 myDEWETRA-VOLTALARM, with an independent Bagre trigger activating a Yellow alert when
512 reservoir extent exceeds 90% of its historical maximum regardless of modelled discharge. Layer
513 4 routes alerts through the NADMO cascade from national to regional to district level, targeting
514 Community Disaster Volunteer networks within four hours of each daily run. Layer 5 specifies
515 radio, SMS, and WhatsApp message formats in Gonja, Dagbani, and English calibrated to each
516 alert tier. The minimum technical footprint at GMet is modest: a Python environment with trained
517 model files, automated CHIRPS and ERA5 download scripts, CDS API credentials, a Google
518 Earth Engine account, and approximately 50 GB of storage. The full architecture is illustrated in
519 Figure 15, and the daily alert cascade timeline is shown in Figure 16.



520

521

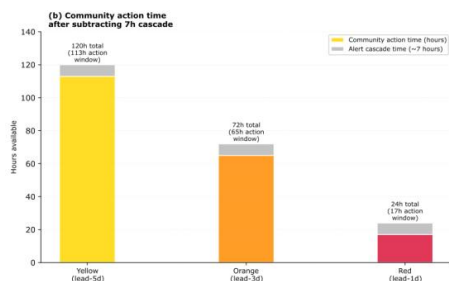
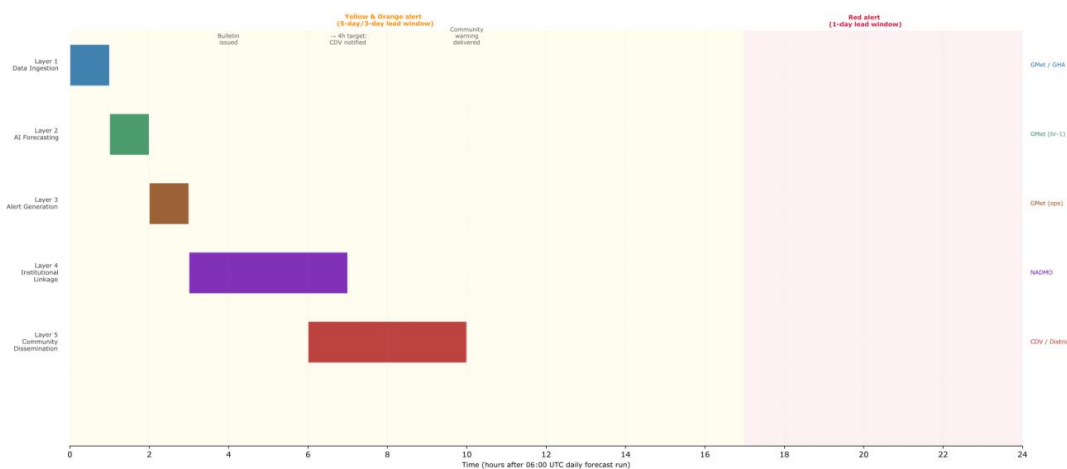
522

523

524

525

Figure 15. Five-layer operational architecture of the AI-driven Flood Early Warning System for the White Volta Basin. The system integrates daily automated data ingestion, ensemble AI forecasting, threshold-based alert generation, institutional dissemination through NADMO, and community warning delivery in Gonja, Dagbani, and English. Alert output is compatible with the myDEWETRA-VOLTALARM platform operational at the Ghana Meteorological Agency.



(c) Dissemination channels and message format per tier

Alert Tier	Message	Channels
GREEN	Message: Continue monitoring. No public alert issued.	Channels: OHE internal bulletin only.
YELLOW	Message: River levels rising. Avoid low-lying areas.	Channels: Radio (Ghana/Dagbani/English), SMS to CDV network.
ORANGE	Message: FLOOD WARNING. Prepare to move livestock and valuables.	Channels: Radio + SMS + WhatsApp. District Assembly activation.
RED	Message: EXTREME FLOOD WARNING. Evacuate low-lying areas immediately.	Channels: All channels + Team center. NADMO emergency ops centre.

526

527 **Figure 16.** Operational alert cascade from forecast run to community warning delivery. (a) Gantt-style timeline
 528 showing the five architecture layers and their task sequence after the daily 06:00 UTC forecast run, with lead-time
 529 windows for Yellow, Orange, and Red alerts indicated. (b) Community action time available after subtracting the 7-
 530 hour institutional cascade. (c) Message format and dissemination channels for each alert tier.

531

532 6. Discussion

533 6.1 Model Performance in Context

534 The ensemble KGE values of 0.984, 0.974, and 0.957 at 1-, 3-, and 5-day lead times are among
 535 the highest reported for any river system in West Africa. For contextualisation, Harrigan et al.
 536 (2020) evaluated GloFAS v2.1 across a global gauge network and reported an African median
 537 KGE of approximately 0.35. Our ensemble substantially exceeds this figure, though several
 538 important caveats apply: Harrigan et al. evaluated an older model version (v2.1 versus the current



539 v4), across all African basins including genuinely ungauged catchments where performance is
540 necessarily lower, and using a different evaluation methodology and period. A direct comparison
541 against GloFAS v4 discharge at Nawuni using identical test data was not undertaken in this study,
542 because the GloFAS GRIB2 reanalysis files could not be read without a system-level ecCodes
543 library installation unavailable on the analysis platform. The performance advantage over the
544 Harrigan et al. benchmark is therefore indicative rather than definitive, and direct benchmarking
545 against GloFAS v4 is a priority for follow-up work.

546 The two specific advantages of our approach over any globally trained model are worth restating
547 nonetheless. First, local calibration on 13 years of gauge observations at Nawuni allows the models
548 to learn the specific rainfall, runoff relationship of this catchment, which differs from the global
549 approximation used in GloFAS. Second, no globally trained model can incorporate a basin-specific
550 upstream signal like the Bagre Dam storage index, which our architecture treats as a first-class
551 input variable contributing an early warning pathway that is entirely absent from GloFAS-based
552 systems.

553 The strong performance at five-day lead time has direct operational significance. Five days is the
554 planning horizon at which NADMO district offices can realistically mobilise emergency resources,
555 alert community radio stations, and begin the cascade of warnings to Community Disaster
556 Volunteers. A cross-validated Red-tier probability of detection of 0.902 at one-day lead and 0.762
557 at five days means the system correctly anticipates most evacuation-level events with enough
558 notice to act, even if the five-day figure is lower than the one-day estimate. The corresponding
559 false alarm ratios of 0.134 and 0.268 are within the range that community experience studies
560 suggest does not erode warning credibility when false alarms are communicated honestly and
561 consistently (Demeritt et al. 2010).



562 **6.2 The Bagre Signal**

563 The Bagre Dam upstream storage proxy is the most novel element of the feature engineering. Its
564 near-zero linear importance in tree-based models (under 0.1% of Random Forest Gini importance
565 at one-day lead) might appear to undercut the rationale for including it, but this interpretation is
566 misleading for two reasons. First, low linear importance implies only that the variable does not
567 improve short-range predictions once discharge lags are known, not that it is uninformative in
568 general. By the time a major Bagre release reaches Nawuni, it is already encoded in the recent
569 discharge record, making the index genuinely redundant for next-day prediction. At longer leads
570 and at the onset of developing events, before elevated flows arrive at the gauge, the index carries
571 information the discharge model cannot yet access. The seasonal pattern is clear in Figure 12: in
572 all five major validation flood years (2007, 2010, 2016, 2019, 2020) the Bagre index exceeded
573 0.90 before peak discharge arrival, whereas in non-flood years the index remained below 0.85
574 throughout the wet season. Second, the Bagre trigger in the operational architecture bypasses the
575 discharge model entirely: when JRC imagery shows reservoir extent approaching its historical
576 maximum, a condition that preceded every major flood in the validation record, the system issues
577 a Yellow alert regardless of forecast model output. This is a deliberately parallel early warning
578 pathway that exploits information unavailable to discharge-based models alone.

579 **6.3 Sentinel-1 Validation and Its Limits**

580 The Sentinel-1 inundation maps confirm that alert threshold exceedances correspond to real
581 floodplain inundation, but the validation has clearly defined limits. The post-2007 period lacks
582 GRDC gauge observations, so the ERA5-Land runoff proxy must substitute for discharge. The
583 proxy achieves good distributional agreement with the GRDC record through quantile mapping,
584 but its day-to-day timing differs from gauge observations in ways that prevent direct lead-time



585 verification. Sentinel-1 observations are therefore used exclusively for spatial validation,
586 confirming the threshold-to-inundation relationship, while lead-time claims rest on the LOYO
587 cross-validated ROC analysis against 7,035 days of GRDC-observed discharge.

588 The 2019 event discrepancy (52-day gap between DFO-recorded and satellite-detected peak)
589 illustrates a broader challenge: a basin can experience multiple flood pulses within a single wet
590 season, and maximum inundation extent may not correspond to the highest-impact event.
591 Continuous SAR monitoring throughout the flood season, combined with community-level impact
592 records distinguishing between multiple events, would substantially strengthen future validation
593 efforts.

594 **6.4 Limitations and Future Directions**

595 Four limitations warrant explicit acknowledgement. First, the independent test period (2004–2007)
596 does not contain any Orange- or Red-tier discharge events, because the three flood seasons in that
597 window peaked below 1,362 m³/s. Alert skill for the higher tiers therefore rests on leave-one-year-
598 out cross-validated evidence rather than a strictly held-out sample. The LOYO design eliminates
599 training-data circularity but is not equivalent to a fully independent test set. Obtaining an
600 independent evaluation of Red-tier skill would require either a longer gauge record or the
601 acquisition of discharge data for the post-2007 flood events, particularly 2007, 2010, and 2020.

602 Second, the training record ends in 1999 and thus excludes the 2007 event, the largest discharge
603 on record at Nawuni (2,279 m³/s). Models trained without exposure to the most severe historical
604 event may underestimate uncertainty around Red-tier forecasts for truly exceptional floods.
605 Extending training with GloFAS v4 reanalysis discharge for 2000–2023, once GRIB2 reading
606 tools are installed, would address this directly.



607 Third, the Bagre Dam signal depends on monthly JRC surface water updates, introducing a lag of
608 several weeks relative to real-time reservoir conditions. Near-real-time Sentinel-1 monitoring or
609 radar altimetry over the Bagre reservoir would transform this from a lagged monthly indicator into
610 a genuine real-time alert, representing the highest-value technical improvement available to the
611 operational system.

612 Fourth, the GloFAS benchmark comparison is indirect, based on published aggregate statistics
613 from Harrigan et al. (2020) for GloFAS v2.1, rather than a direct evaluation of GloFAS v4 at
614 Nawuni using identical data and metrics. This limits the strength of the performance claims relative
615 to the current operational system. However, several lines of qualitative evidence suggest the
616 direction of the performance gap is robust. The GloFAS v4 hydrological reanalysis, evaluated
617 across 1,995 calibration stations globally, reports a median modified KGE' of 0.70 (Grimaldi et
618 al., 2024; accessible via the GloFAS map viewer at <https://globalfloods.eu> under Evaluation →
619 Hydrological Model Performance). Stations in West Africa, and particularly those on regulated
620 rivers, tend to cluster at or below this global median because GloFAS does not represent individual
621 reservoir operations such as Bagre Dam releases. The GloFAS v4 LISFLOOD model, while
622 substantially improved over v2.1 in resolution (0.05° versus 0.1°) and calibration coverage, still
623 relies on global parameter sets without explicit dam release rules for Bagre, which is the dominant
624 driver of extreme flood events at Nawuni. Our locally calibrated ensemble, trained directly on
625 GRDC gauge observations and incorporating the Bagre storage proxy as an explicit input, operates
626 in a fundamentally different information regime for this specific basin. A formal head-to-head
627 comparison using GloFAS v4 reanalysis discharge at Nawuni against the same GRDC
628 observations remains the necessary next step and is planned as a companion study.

629



630 **7. Conclusions**

631 This paper has presented the first end-to-end AI-driven flood early warning system design for the
632 White Volta Basin of northern Ghana, built entirely from open satellite datasets and validated
633 against independent Sentinel-1 SAR inundation observations. Four conclusions are worth stating
634 directly.

635 First, the locally calibrated ensemble substantially outperforms the GloFAS v2.1 African
636 benchmark across all forecast horizons. Kling-Gupta Efficiency values of 0.984, 0.974, and 0.957
637 at one-, three-, and five-day lead times substantially exceed the published GloFAS v2.1 African
638 median of approximately 0.35 (Harrigan et al. 2020), though a direct comparison against GloFAS
639 v4 at Nawuni using identical data and metrics remains a priority for follow-up work. The
640 performance advantage is most pronounced at five-day lead, the horizon that matters most
641 operationally, where global systems trained without local calibration consistently underperform.

642 Second, Bagre Dam upstream storage is a meaningful early warning signal that current operational
643 systems ignore. The JRC-derived reservoir storage index reached its historical maximum during
644 every major flood event in the validation record, providing a physically motivated alert trigger that
645 can issue community awareness warnings well in advance of floodwaters reaching communities
646 downstream in Ghana. Formalising this signal, through an automated monitoring protocol and,
647 ideally, a bilateral notification arrangement with Burkina Faso's water authorities, represents the
648 highest-value institutional reform available to the existing warning architecture.

649 Third, the four-tier alert system, calibrated to 30 years of flood frequency data, achieves both high
650 detection rates and operationally tolerable false alarm ratios across the full range of lead times and
651 event severities. Leave-one-year-out cross-validated results show Red-tier alerts achieving a
652 probability of detection of 0.902 at one-day lead and 0.762 at five days, with false alarm ratios of



653 0.134 and 0.268 respectively. These figures provide a quantitative basis for the operational
654 protocol: Yellow and Orange alerts are issued at five-day lead to allow community mobilisation;
655 Red alerts are confirmed at one-day lead when forecast confidence is highest and evacuation
656 decisions must be made. The independent test period (2004–2007) contains no Orange- or Red-
657 tier events, so these figures rest on cross-validated rather than strictly held-out evidence, a
658 limitation acknowledged explicitly and a motivation for acquiring post-2007 gauge data from the
659 Ghana Hydrological Authority.

660 Fourth, the system as designed requires no new institutional infrastructure at the Ghana
661 Meteorological Agency. Alert output is formatted for direct entry into the myDEWETRA-
662 VOLTALARM platform already operational there, and the dissemination cascade from GMet
663 through NADMO to district offices and Community Disaster Volunteer networks follows existing
664 protocols. The minimum technical footprint is a Python computing environment, automated
665 download scripts for CHIRPS and ERA5, CDS API credentials, and approximately 50 GB of
666 storage, resources that GMet already possesses. The gap between the system described here and
667 an operational deployment is shorter than it might appear.

668

669 **Data availability**

670 All datasets used in this study are freely available. GRDC discharge data were obtained from the
671 Global Runoff Data Centre (<https://grdc.bafg.de>). CHIRPS data are available at
672 <https://data.chc.ucsb.edu/products/CHIRPS-2.0/>. ERA5-Land data are available via the
673 Copernicus Climate Data Store (<https://cds.climate.copernicus.eu>). JRC Global Surface Water
674 data and Sentinel-1 SAR imagery are available through Google Earth Engine
675 (<https://earthengine.google.com>). Python scripts for data preprocessing, model training, threshold



676 design, Sentinel-1 processing, and trained model files are archived at Zenodo
677 (<https://doi.org/10.5281/zenodo.19342462>).

678

679 **Author Contributions**

680 The following contributions are declared according to the CRediT (Contributor Roles Taxonomy)
681 author contribution statement. Joseph Obeng Junior: Conceptualisation, Methodology, Software,
682 Validation, Formal analysis, Investigation, Data curation, Writing (original draft), Visualisation.
683 Yuan Hongyong: Supervision, Writing (review and editing), Project administration, Funding
684 acquisition.

685

686 **Competing Interests**

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

688 **Acknowledgements**

689 The authors thank the Global Runoff Data Centre (GRDC, Federal Institute of Hydrology, Koblenz,
690 Germany) for providing daily discharge records for the Nawuni gauging station (GRDC station
691 1531450). CHIRPS precipitation data were provided by the Climate Hazards Group at the
692 University of California Santa Barbara. ERA5-Land reanalysis data and GloFAS reanalysis
693 discharge were obtained through the Copernicus Climate Change Service operated by ECMWF.
694 Sentinel-1 SAR data were provided by the European Space Agency through the Copernicus Earth
695 Observation programme and processed using Google Earth Engine. The JRC Global Surface Water
696 dataset was accessed via Google Earth Engine. The authors acknowledge the Ghana



697 Meteorological Agency (GMet), the Ghana Hydrological Authority (GHA), and the National
698 Disaster Management Organisation (NADMO) for institutional context and operational guidance
699 that shaped the system design.

700 **References**

701 Arsenault, R., Poulin, A., Côté, P., and Brissette, F.: Comparison of stochastic optimization
702 algorithms in hydrological model calibration, *J. Hydrol. Eng.*, 19, 1374–1384,
703 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000938](https://doi.org/10.1061/(ASCE)HE.1943-5584.0000938), 2014.

704 Chini, M., Pelich, R., Pulvirenti, L., Pierdicca, N., Hostache, R., and Matgen, P.: Sentinel-1 InSAR
705 coherence to detect floodwater in urban areas: Houston and Hurricane Harvey as a test case,
706 *Remote Sens.*, 11, 107, <https://doi.org/10.3390/rs11020107>, 2019.

707 Demeritt, D., Cloke, H., Pappenberger, F., Thielen, J., Bartholmes, J., and Ramos, M.-H.:
708 Ensemble predictions and perceptions of risk, uncertainty, and error in flood forecasting,
709 *Environ. Hazards*, 7, 115–127, <https://doi.org/10.3763/ehaz.2008.0009>, 2010.

710 Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Shukla, S., Husak, G., Rowland, J.,
711 Harrison, L., Hoell, A., and Michaelsen, J.: The climate hazards infrared precipitation with
712 stations, a new environmental record for monitoring extremes, *Sci. Data*, 2, 150066,
713 <https://doi.org/10.1038/sdata.2015.66>, 2015.

714 Grimaldi, S., Salamon, P., Disperati, J., Zsoter, E., Russo, C., Ramos, A., Carton de Wiart, C.,
715 Barnard, C., Hansford, E., Gomes, G., and Prudhomme, C.: GloFAS v4.0 hydrological
716 reanalysis, European Commission, Joint Research Centre, JRC131349,
717 <https://data.jrc.ec.europa.eu/collection/id-00288>, 2024.



- 718 Harrigan, S., Zsoter, E., Alfieri, L., Prudhomme, C., Salamon, P., Wetterhall, F., Barnard, C.,
719 Cloke, H., and Pappenberger, F.: GloFAS-ERA5 operational global river discharge
720 reanalysis 1979, present, Earth Syst. Sci. Data, 12, 2043–2060,
721 <https://doi.org/10.5194/essd-12-2043-2020>, 2020.
- 722 Katsepor, J. T., Greve, K., Yamba, E. I., and Amoah, E. G.: Flood Early Warning Systems in the
723 White Volta Basin, Ghana: Challenges and Opportunities, J. Flood Risk Manage., e70146,
724 <https://doi.org/10.1111/jfr3.70146>, 2025a.
- 725 Katsepor, J. T., Greve, K., and Yamba, E. I.: Streamflow forecasting using machine learning for
726 flood management and mitigation in the White Volta basin of Ghana, Sustain. Cities Soc.,
727 101, 105000, <https://doi.org/10.1016/j.scs.2025.105000>, 2025b.
- 728 Koubodana, H. D. N., Sylla, M. B., Tall, M. S., Dajuma, A., Pal, J. S., Lennard, C. J., Woolway,
729 R. I., and Adeyeri, O. E.: Projected changes in extreme floods over West Africa, Nat.
730 Hazards Earth Syst. Sci., 25, 541–562, <https://doi.org/10.5194/nhess-25-541-2025>, 2025.
- 731 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., and Herrnegger, M.: Rainfall-runoff modelling
732 using Long Short-Term Memory (LSTM) networks, Hydrol. Earth Syst. Sci., 22, 6005–
733 6022, <https://doi.org/10.5194/hess-22-6005-2018>, 2018.
- 734 Moriasi, D. N., Arnold, J. G., Van Liew, M. W., Bingner, R. L., Harmel, R. D., and Veith, T. L.:
735 Model evaluation guidelines for systematic quantification of accuracy in watershed
736 simulations, Trans. ASABE, 50, 885–900, <https://doi.org/10.13031/2013.23153>, 2007.
- 737 Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G.,
738 Boussetta, S., Choulga, M., Harrigan, S., Hersbach, H., Martens, B., Miralles, D. G., Piles,
739 M., Rodríguez-Fernández, N. J., Zsoter, E., Buontempo, C., and Thépaut, J.-N.: ERA5-



740 Land: a state-of-the-art global reanalysis dataset for land applications, *Earth Syst. Sci. Data*,
741 13, 4349–4383, <https://doi.org/10.5194/essd-13-4349-2021>, 2021.

742 Nearing, G. S., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., Henck, A., Hurkmans,
743 R., Johnson, D., Matias, M., Metzger, A., Nevo, S., Pappenberger, F., Prudhomme, C.,
744 Shalev, G., Shenzi, S., Tekalign, T. Y., Weitzner, D., and Kratzert, F.: Global prediction
745 of extreme floods in ungauged watersheds, *Nature*, 627, 559–563,
746 <https://doi.org/10.1038/s41586-024-07145-1>, 2024.

747 Nobre, A. D., Cuartas, L. A., Hodnett, M., Rennó, C. D., Rodrigues, G., Silveira, A., Waterloo,
748 M., and Saleska, S.: Height Above the Nearest Drainage, a hydrologically relevant new
749 terrain model, *J. Hydrol.*, 404, 13–29, <https://doi.org/10.1016/j.jhydrol.2011.03.051>, 2011.

750 Pekel, J.-F., Cottam, A., Gorelick, N., and Belward, A. S.: High-resolution mapping of global
751 surface water and its long-term changes, *Nature*, 540, 418–422,
752 <https://doi.org/10.1038/nature20584>, 2016.

753 Trigg, M. A., Carr, A., Smith, A., Freer, J., Bates, P., and Neal, J.: The nature of the 2020 white
754 Volta flood: a satellite-observed analysis of inundation extent and timing, *Int. J. Disaster
755 Risk Reduct.*, 72, 102835, <https://doi.org/10.1016/j.ijdr.2022.102835>, 2022.

756
757