# Postprocessing East African rainfall forecasts using a generative machine learning model

Bobby Antonio<sup>1</sup>, Andrew T T McRae<sup>1</sup>, David MacLeod<sup>2</sup>, Fenwick C Cooper<sup>1</sup>, John Marsham<sup>3</sup>, Laurence Aitchison<sup>4</sup>, Tim N Palmer<sup>5</sup>, and Peter A G Watson<sup>4</sup>

<sup>1</sup>University of Oxford <sup>2</sup>Cardiff University <sup>3</sup>University of Leeds <sup>4</sup>University of Bristol <sup>5</sup>Oxford University

March 05, 2024

## Abstract

Existing weather models are known to have poor skill at forecasting rainfall over East Africa, where there are regular threats of drought and floods. Improved forecasts could reduce the effects of these extreme weather events and provide significant socioeconomic benefits to the region. We present a novel machine learning-based method to improve precipitation forecasts in East Africa, using postprocessing based on a conditional generative adversarial network (cGAN). This addresses the challenge of realistically representing tropical rainfall, where convection dominates and is poorly simulated in conventional global forecast models. We postprocess hourly forecasts made by the European Centre for Medium-Range Weather Forecasts Integrated Forecast System at 6-18h lead times, at  $0.1^{(circ)}$  resolution. We combine the cGAN predictions with a novel neighbourhood version of quantile mapping, to integrate the strengths of machine learning and conventional postprocessing. Our results indicate that the cGAN substantially improves the diurnal cycle of rainfall, and improves predictions up to the  $99.9^{(text{th})}$  percentile ( $\frac{1}{\sin 10} \frac{1}{\cot 1}$ ). This improvement extends to the 2018 March–May season, which had extremely high rainfall, indicating that the approach has some ability to generalise to more extreme conditions. We explore the potential for the cGAN to produce probabilistic forecasts and find that the spread of this ensemble broadly reflects the predictability of the observations, but is also characterised by a mixture of under- and over-dispersion. Overall our results demonstrate how the strengths of machine learning approaches can bring to this region.



















radar data in the area. In Ageet et al. (2022) a range of satellite rainfall estimating prod-160 ucts, including the IMERG product, were compared with rain gauges in an area around 161 Uganda (including parts of Kenya, Tanzania, Sudan and the Democratic Republic of Congo) 162 over 17 years. Based on a combined assessment of quantile-quantile plots, correlation, 163 and skill scores such as hit rate and false alarms, they identi ed the IMERG product (VO6B) 164 as the best performing at daily resolution. However, it still has biases; for example it has 165 a tendency to underestimate the rainfall rate, and over-predict the frequency of rainfall. 166 There are also known issues with similar satellite products in mountainous areas (Dinku 167 et al., 2010), which means these observations may be more unreliable over areas such as 168 the Ethiopian Highlands and parts of the Rift Valley. A dry bias has also been observed 169 in several studies (e.g. Vogel et al. (2018)). Overall, though, it provides a good source 170 of data with a high temporal and spatial resolution over our target region, and it has been 171 used in other studies in this area (e.g. Woodhams et al. (2018); Finney et al. (2019); Ca-172 faro et al. (2021)). 173

The forecast dataset used is the ECMWF IFS HRES deterministic hourly forecast 174 (ECMWF, 2023) as this tends to perform amongst the best compared to similar mod-175 els (Haiden et al., 2012). IFS forecasts are provided at 00h and 12h and we use lead times 176 within a 6-18h window, corresponding to short-range weather prediction (however it is 177 expected that the method we use could also equally apply to longer lead times). The data 178 is interpolated from 9km 9km resolution to  $\mathbf{O}$  =  $\mathbf{O}^{1\circ}$  to match the grid points of 179 the IMERG precipitation. The data starts at March 2016, after the increase in horizon-180 tal resolution for the IFS with the release of Cycle 41r2. To ensure the precipitation fore-181 casts are reasonably consistent, we use data up until September 2021 before the upgrade 182 to the convection parameterization scheme with the release of Cycle 47r3 in October 2021 183 (ECMWF, 2023). 184

185

2.3 Machine Learning Model

Our model architecture uses the same architecture and code that L. Harris et al. (2022) used to postprocess UK rainfall forecasts. This is itself based on Leinonen et al. (2020) and a variant was developed for downscaling tropical cyclone rainfall by Vosper et al. (2023). A conditional Wasserstein GAN is trained to predict realistic rainfall patterns conditioned on several meteorological inputs together with constant inputs such as orography, using the IMERG data as ground truth. We use the same approach to test

{7{

whether it will transfer to also perform well at postprocessing forecasts in a tropical do-main.

Both the generator and discriminator of the GAN are deep neural networks, pri-194 marily made up of residual blocks, where each residual block contains two convolution 195 layers that use square convolutional kernels of width 3 pixels (see e.g. Goodfellow et al. 196 (2016) for background on convolutional neural networks). The generator is composed 197 of 7 residual blocks (each with lters), with a nal softplus activation function, giv-198 ing a total of 27  $f_q$  intermediate arrays each of size 2**200**. The discriminator is 199 made up of 3 residual blocks (each withlters), and two dense layers, giving a total 200 of 2 3  $f_d$  intermediate arrays each of size 20200. Excluding the output layers, 201 PReLU activation functions were used, where we set the preLU to 202 0.2 following L. Harris et al. (2022). The number of noise channels was set to 4, and the 203 learning rates for the generator and discriminator were set equal 000, 1 with the 204 discriminator being trained for 5 steps for every 1 step of generator training. The batch 205 size was set to 2 based on hardware memory constraints, and the Adam optimiser was 206 used for training. 207

Following L. Harris et al. (2022) we use a Wasserstein GAN (Arjovsky et al., 2017), 208 which has been demonstrated to improve training stability in many cases (Creswell et 209 al., 2018). This modi es the GAN discriminator to output low numbers for real samples 210 and high numbers for fake samples, rather than producing a number in the range [0 211 and modi es the loss function to approximate the Wasserstein distance between the gen-212 erated and true distributions (Gulrajani et al., 2017). This approximation is parame-213 terised by a gradient penalty parameter which we set to 10 in line with Gulrajani et 214 al. (2017). 215

In order to perform shorter experiments to tune hyperparameters, smaller models with $f_g = 32$ ,  $f_d = 128$  were trained for46 10<sup>4</sup> iterations, and then nally larger models with $f_g = 64$ ,  $f_d = 256$  were trained for23 10<sup>5</sup> steps, and used for evaluation; thus our largest model was smaller than the model in L. Harris et al. (2022) that had  $f_g = 128$ ,  $f_d = 512$ . However, since the model is not being used for downscaling, and because we use a larger domain, the model dimensions scale di erently, and so using a smaller number of channels was required to achieve a reasonable training time.

{8{

As inputs, the model uses IFS forecast variables from the same time as the target 223 rainfall forecasts. On top of the 9 variables used in L. Harris et al. (2022), we used 11 224 extra elds, including temperature, convective precipitation, vertical velocity, and rel-225 ative humidity (some of which are at several pressure levels; see appendix Appendix A 226 for a full table of inputs). Convective inhibition was included, with null values set to 0. 227 We included these extra variables as they contain important information about convec-228 tive processes, which are critical for forecasting in East Africa. Based on the transfor-229 mations applied in L. Harris et al. (2022), we normalise the input variables; precipita-230 tion variables are log-normalised via!  $\log_{10}(1 + x)$ , whilst others were either di-231 vided by the maximum value, or normalised to fall within the maximum and minimum 232 values (see Appendix A). 233

Model checkpoints were saved every 3200 steps, and the best model in the last 1/3rd of checkpoints was selected based on judgement of the combined performance on CRPS, RAPSD and mean squared error, plus visual evaluation of the samples produced. Our batch size was limited to 2 due to the need to generate an ensemble to calculate part of the loss function (discussed in the next paragraph). All models were trained on a single Nvidia A100 GPU.

One notable addition by L. Harris et al. (2022) is the inclusion of a 'content loss' term, inspired by Ravuri et al. (2021), which penalises GAN predictions that do not have an ensemble mean close to the observed value. Speci cally, at each training step the generator produces an ensemble of predictions (set to 8 in this work), and the generator loss function includes a mean-squared error term between the observed image and the ensemble mean of the generated samples.

During validation of the models, we observed that using log normalisation of the 246 output precipitation predictions, as done by L. Harris et al. (2022), produced a distri-247 bution of rainfall that tended to greatly overestimate the observations at the extreme 248 rainfall values. Removing the log normalisation of the output rainfall values remedied 249 this, and also removed the need to clip the predicted rainfall values to a given maximum, 250 as done in L. Harris et al. (2022). This also required modifying the content loss param-251 eter $\lambda$ , with  $\lambda$  = 100 appearing to produce the best results according to a joint assess-252 ment of quantile-quantile plots, CRPS and RALSD (see Sec. 2.6). 253

{9{

In L. Harris et al. (2022), samples are grouped into prede ned bins based on the fraction of grid points exceeding a set threshold, and then at training time samples are drawn more frequently from the high rainfall bins, in order to oversample higher rainfall values and improve performance in these cases. However for our data, the threshold used by L. Harris et al. (2022) for postprocessing UK rainfall was not appropriate, and our attempts at using a similar approach did not give any improvements on the validation set. Therefore we did not apply this oversampling.

To increase the variation in the samples seen during training, we randomly cropped the 270 265 images to smaller images of 22000, as this has been demonstrated to improve the generalisability of deep learning models (Goodfellow et al., 2016), and produces output similar in size to that in L. Harris et al. (2022).

Similarly to the results in L. Harris et al. (2022), we observed that the model skill can vary considerably between training steps (as measured by CRPS, RAPSD and visual inspection), so it was necessary to have a validation dataset set aside in order to choose the nal model, and this data was not used in the nal training.

## 269 2.3.1 Training and evaluation strategy

For training and evaluating the model, the dataset was split up as follows:

Training set: March 2016 { February 2018 and July 2018 { Sept 2020 (excluding
 validation months)

<sup>273</sup> Validation set: Jun 2018, Oct 2018, Jan 2019, March 2019

Test sets: October 2020 - September 2021, March - May 2018

We used the nal year as a primary test dataset, and the 2018 long rains (March-May) as an extreme test set, since this was a season of particularly heavy rainfall (Kilavi et al., 2018) for which the March{May rainfall was signi cantly higher than in any other season in the full IMERG dataset.

The purpose of the validation dataset is to guide choice of the model structure and hyperparameters. The development process was to train di erent versions of the model on the training dataset, then evaluate these on the validation dataset to select the best version. This avoids over tting on the test data by selecting a model that performs well by chance. The standard choice of validation set would be the period October 2019{September

{10{

2020, which spans a full year and would sit between the training and test periods. How-284 ever, the rains of October-December 2019 were exceptionally high (Wainwright et al., 285 2021), as were the long rains of March{May 2020 (Palmer et al., 2023). So to avoid val-286 idation over an atypical year, which may have given an inaccurate assessment of the model's 287 general performance, we chose to validate over the period of June 2018{May 2019. Rather 288 than use a full year for validation, we also chose to maximise the amount of training data 289 by including a month from each of the di erent seasons in the validation period. This 290 sampling variability observed for this size of validation data also indicated there was no 291 additional bene t from using a whole year. 292

All evaluation results reported in Secs. 3.1 and 3.2 are evaluated on 4000 unique hours randomly sampled from the unseen test period October 2020 - September 2021, with 20 ensemble members used in the example plots. The ensemble calibration results in Sec. 3.3 are assessed over 500 unique hours sampled uniformly from the same period with an ensemble size of 100.

For the extreme rainfall evaluation in Sec. 3.4, we analyse all of the hours from March-May 2018. Since much of the anomalous rainfall in this season was concentrated over Kenya, we restrict our analysis to this regio**6**° §4 5.1°N, 332° 431°E, see Fig. 1).

## 2.4 Quantile mapping

Since it is known that IFS forecasts with postprocessing outperform those without in this region (Vogel et al., 2020), and IFS forecasts are not speci cally tuned to reproduce the properties of IMERG observations, we applied quantile mapping to the IFS forecasts (see e.g. Maraun and Widmann (2017)) to provide a stronger baseline.

Additionally, since GAN predictions are not guaranteed to precisely capture the 306 rainfall distribution, and we observed that our GAN predictions tended to under-predict 307 high rainfall values, we produced a variant of our model with quantile mapping applied 308 to the output. In doing so we aimed to combine the strengths of both postprocessing meth-309 ods to achieve an overall more accurate and realistic forecast. The GAN could be ex-310 pected to perform well at producing predictions with realistic spatial structure, but not 311 necessarily with realistic point frequency distributions. Quantile mapping can greatly 312 improve the latter. 313

{11{

We used empirical quantile mapping rather than a distribution-based quantile map-314 ping approach, since it has been demonstrated to work well (Gudmundsson et al., 2012), 315 and does not require a parametric distribution. Our method is based on the well-used 316 methods outlined in Boe et al. (2007), Deque (2007), and Maraun and Widmann (2017), 317 in which empirical cumulative density functions are calculated over the training period 318 and used to create a mapping between the forecast quantiles and the observed quantiles. 319 In general this means creating an estimate of the cumulative density functions 320  $F_{o}$  of the forecast and observations respectively, and mapping the forecast  $t_{f}$  values 321 an adjusted value  $r_{f}$  according to: 322

$$x_{f} = F_{0}^{-1}(F_{f}(x_{f}))$$
(1)

In Bœ et al. (2007), percentiles at 1% spacing are rst calculated on the training set to nd an approximation to the quantile distributions. Forecast values are then mapped into quantile values relative to the training data, and then converted into adjusted forecast values using the observed quantile values (using linear interpolation when the quantile falls between the known quantile values).

Since the East African precipitation is low, there can be multiple quantiles that are O; therefore for a forecast of Omm/hr there is no way to tell which quantile it belongs to. We follow the method in Bœ et al. (2007) and pick one of the O-valued forecast quantiles at random, then assign the value of the matching observational quantiles. In practise, this can lead to low level noise on the corrected forecast, but replicates the high level statistics.

In order to better match the tail of the frequency distribution in our work, the step size between the quantiles was decreased towards the higher quantiles; so a step size of 0.01 was used up to 0.99, a step size of 0.001 used from 0.99 to 0.999, and so on up to the 999999 percentile, above which we observed signi cant sampling variability. Note that these percentiles are calculated at the grid box level, so that the number of samples available to calculate these quantiles is 4000 265 grid boxes for the validation and test datasets.

For data greater than the maximum value observed in training, we follow the additive uplift method (Boe et al., 2007; Deque, 2007) and add the uplift of the highest quan-

{12{



Figure 2: a) Illustration of the general method for how grid cells are grouped together in order to estimate quantiles. In this example, the spatial domain is split into squares of 4 4 grid cells, giving 9 separate large square regions, and in each region the quantiles are calculated. b) To calculate a particular quantile for a grid point, lled in black, we perform a weighted sum of the value of this quantile calculated in the square regions nearest to that point. The weighting for each large region is proportional to the number of small squares inside the red dashed square. In this example, the red dashed square covers 25 grid cells and the weightings would  $\frac{12}{29}e \frac{3}{25}, \frac{8}{25}, \frac{2}{25}$ .

tile for the IMERG and IFS data. For example,  $ib_{max}$ ,  $f_{max}$  are the highest values seen in the training set for the observations and forecast respectively, then for any forecast value in the test data greater than we add the uplift  $o_{max}$  ( $f_{max}$ ) to it.

From experiments we found that the typical approach of quantile mapping the GAN at each grid point individually was not a robust approach for the highest values. Therefore we aggregated the data into square regions to calculate quantiles (Fig. 2(a)). The intuition is that nearby points will have similar distributions and so we can gain accuracy by grouping nearby points together.

To avoid any artefacts due to the edges of these domains, the quantiles for a given 353 grid cell were calculated as a weighted average of the nearest square regions; speci cally, 354 the quantiles used to update the values at grid  $\alpha \in M$  are calculated as a weighted 355 sum, where the weighting is calculated by drawing a square around and count-356 ing the number of grid points that fall into each quantile grouping (Fig. 2(b)). This is 357 partly motivated by the ease of implementation, as this can be easily done by broadcast-358 ing the grouped quantiles to the same dimensions as the original grid, and using square 359 convolutions with re ective padding to calculate the weighted versions of each quantiles. 360

{13{

The length scale of the weighting window was chosen to be the same as the length scale of the quantile groupings, as this was empirically observed to produce reasonably smoothed values.

To decide on the optimal grouping in the spatial domain, the quantile mapping ap-364 proach described above was trained on the same training data as the cGAN, and then 365 used to perform quantile mapping on forecasts in the validation set. The best param-366 eters were chosen by calculating the quantiles over the whole domain after quantile map-367 ping, and comparing these to the quantiles of the IMERG data over the whole domain 368 using mean-square error (MSE) up to the precentile. Using this method, the 369 cGAN performed best when split into 4 square regions (each region having width 66 grid 370 boxes), whilst the IFS forecast performed best when split into 9 regions (each region hav-371 ing width 30 grid boxes). 372

<sup>373</sup> We denote these quantile mapped models as cGAN-qm and IFS-qm.

## 2.5 Assessing Sample Variability

For many diagnostics, particularly those concerned with high rainfall events, the 375 results can be swayed by the presence or absence of a small number of high rainfall events 376 particular to the test year. To estimate the uncertainty due to these e ects, we use boot-377 strapping along the time dimension (Efron & Tibshirani, 1986). To perform bootstrap-378 ping for a property of interestalculated over a set of hourly samples, we sample 379 with replacement times from the samples, and repeat this protestimes, result-380 ing in M sets of samples of size Then the mean and standard error  $\infty$  can then be 381 estimated from the mean and standard deviation conficulated on the bootstrap sam-382 ples. Since the hours are sampled uniformly at random, this method does not take into 383 account the correlation between adjacent hours, and so it is likely that the standard er-384 ror calculated from this method is an underestimate. 385

386

2.6 Forecast veri cation measures

387

## 2.6.1 Radially Averaged Power Spectral Density (RAPSD)

In order to assess the spatial realism of the generated forecasts, we use the Radi ally Averaged Power Spectral Density (RAPSD) (Sinclair & Pegram, 2005; D. Harris et

{14{

al., 2001). This is calculated by taking the 2D Fourier transform of the precipitation image, and averaging the power spectrum over the total wavenumber magnitude, yielding
a one-dimensional series showing the distribution of weights given to di erent frequencies. For assessing multiple forecast images we take the mean RAPSD over the images,
and for situations where we need to summarise the overall similarity of two RAPSD curves
we use the Radially Averaged Log Spectral Distance (RALSD) de ned in L. Harris et
al. (2022).

#### 397

## 2.6.2 Equitable Threat Score (ETS)

The Equitable Threat Score (ETS) measures the balance between the hit rate and false alarm rate, whilst accounting for the probability of random events (Schaefer, 1990; Wilks, 2019a). This score is used in operational forecast veri cation (Mittermaier et al., 2013) and in Manzato and Jolli e (2017) was shown to be one of the most robust metrics with respect to random (unskillful) changes in the forecast. It is de ned as:

403 404

$$ETS := \frac{TP TP_r}{TP + FP + FN TP_r}$$
(2)

where TP, FP, FN, are the number of true positives, false positives and false negatives. TP<sub>r</sub> accounts for the number of true positives we would expect to achieve by guessing at random, and is often estimated from the data using the formula:

$$TP_{r} = \frac{(TP + FP)(TP + FN)}{N}$$
(3)

410

408 409

## 2.6.3 Fractions Skill Score (FSS)

Many forecast veri cation scores, such as mean square error or the ETS, do not always align with human forecasters' subjective evaluations. There have been many different approaches employed to try and remedy this problem (see e.g. Gilleland et al. (2009)). One approach, usually called the neighbourhood approach, is to smooth the forecasts and observations by averaging the forecast around each grid cell with a particular length scale before applying a forecast metric (Ebert, 2008). A commonly used metric in this class is the Fractions Skill Score (FSS) (N. Roberts, 2008; N. M. Roberts & Lean, 2008).

To calculate the FSS, we rst choose a threshold rainfall value d for each grid cell of the forecast and observations we calculate the fracting icansd  $O_{tij}$  of neighbour-

{15{

ing cells for which the rainfall exceeds the threshold, white j index the time, latitude and longitude axes respectively. The FSS is then de ned as:

$$FSS(n,r) := \frac{\sum_{t=1}^{T} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} 2F_{tij} O_{tij}}{\sum_{t=1}^{T} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} F_{tij}^2 + O_{tij}^2}$$
(4)

We use the pySTEPS implementation of the FSS (Pulkkinen et al., 2019), which performs the averaging using square convolutions with zero-padding.

In the limit of large neighbourhood size, the FSS approaches the asymptotic limit FSS<sub> $\infty$ </sub> where (N. M. Roberts & Lean, 2008):

FSS<sub>$$\infty$$</sub> =  $\frac{2f_0 f_m}{f_0^2 + f_m^2}$  (5)

where *f*<sub>o</sub>, *f*<sub>m</sub> are the observed and modelled frequency of exceeding the threshold, respectively. Thus the value that the FSS reaches at large neighbourhood sizes indicates
the level of bias in the average number of grid boxes exceeding the threshold, with FSS
1 for no bias.

## 434 2.6.4 Spread error

422 423

In order to assess how well calibrated the probabilities of the generated forecast are, we use a spread-error plot, commonly used to assess ensemble calibration (Leutbecher & Palmer, 2008). For an ensemble of forec $as_{its}g_{i=1}^{M}$  with ensemble meaple, and an observationy<sub>t</sub>, the spreads<sub>t</sub> and errore<sub>t</sub> are de ned as:

439 
$$s_{\rm t}^2 = \frac{1}{M} \sum_{\rm i=1}^{\rm M} (f_{\rm i;t} - \mu_{\rm t})^2$$
(6)

$$e_t^{440} = (y_t \quad \mu_t)^2 \tag{7}$$

<sup>442</sup> Note that for a nite number of ensemble memble rsve also include a correction to<sup>443</sup> the spread:

$$s_{t} = \frac{M+1}{M-1}s_{t}$$
(8)

To produce a spread-error plot, we rst calculate the spread values for each grid cell and each time value. Then we split the paired observations and forecasts into bins (of size 100 in our case) according to this spread value, and for each bin we calculate the root mean squared spread and error values. For a perfect forecast ensemble with in nite members, the spread of the ensemble will equal the average error between the ensemble mean and observations, so that an ideal spread-error plot is a straight line with gradient 1.

#### 453 2.6

## 2.6.5 Rank histogram

Another method to assess the statistical calibration of an ensemble forecast is the 454 use of a rank histogram (also known as a Talagrand diagram) (Wilks, 2019a). To con-455 struct a rank histogram, for each sample we rank the observed value relative to the en-456 semble members, and then average this over all the samples. This gives a frequency of 457 how many times the observations were seen to be in each rank, which can be plotted as 458 a histogram. A perfectly calibrated ensemble produces a at histogram. For an imper-459 fect ensemble, the spread of the ensemble members may be too wide, such that the ob-460 servations will rank in the middle most of the time, producing a peak in the histogram. 461 For the reverse scenario, the spread is too narrow leading to a U-shaped histogram in-462 dicating the ensemble members are too narrowly spread. A histogram sloping to the left 463 or right is also indicative of conditional under-forecasting or over-forecasting respectively 464 (Hamill, 2001; Wilks, 2019a). In this work we use the pySTEPS implementation of the 465 rank histogram (Pulkkinen et al., 2019). 466

467

## 2.6.6 Continuous Ranked Probability Score (CRPS)

The Continuous Ranked Probability Score, or CRPS, is a particularly important score in assessing probabilistic forecasts, as *idtriteday proper* score (Wilks, 2019b), which means that the score is only maximised when the forecast distribution equals the target distribution. For a cumulative forecast distribution, the CRPS for an ob-

served occurrence of(e.g. observed rainfall value) is de ned as:

473  
474 
$$CRPS(F,x) = \int_{-\infty}^{\infty} [F(y) \quad 1_{y \ge x}]^2 dy$$
(9)



Figure 3: (a) Quantile-quantile plot, up to the **9999** percentile. Red circles (diamonds) indicate quantiles for the IFS (IFS-qm) model. Blue circles (diamonds) indicate quantiles for the cGAN (cGAN-qm) model. The black dashed line is the line along which a perfectly calibrated forecast would sit. The error bars indicate an estimate of 2 standard errors from 1000 bootstrap samples (only shown for **90 9 9 9 9 9 9 9 9 9** centiles). (b) A histogram showing the distribution of rainfall values; the vertical dashed blue line indicates the **999**<sup>h</sup> percentile of observed rainfall

where  $1_{y \ge x}$  is the Heaviside step function. The CRPS is a univariate measure, so does

not properly account for spatial correlations. Whilst there is a multivariate generalisa-

tion of the CRPS, the energy score (Gneiting & Raftery, 2007), the CRPS is more com-

478 monly used and has been used in previous related works, so we use it in this work as one

of many validation metrics, in order to choose the best model.

## 480 3 Evaluation

In this section we presents results of evaluating the model on unseen data. Eval uations are performed on the primary test dataset except for Sec. 3.4 which is evaluated
 on the extreme Long Rains of March{May 2018.



Figure 4: The approximate return periods for di erent rainfall thresholds to be exceeded at any grid point in the spatial domain in the training set. The dashed lines indicate the values of particular high percentiles.

## 3.1 Climatological properties of the forecasts

We rst assess how well the forecasts capture the distribution of rainfall, shown by a quantile-quantile plot and a histogram of rainfall distribution for 4000 samples, shown in Fig. 3 (a) and (b) respectively. The unpostprocessed model outputs are shown together with quantile-mapped outputs. From these we can see that, without postprocessing, the distribution of cGAN output is an improvement upon that of the IFS output up to extremely high levels of rainfall (around 50mm/hr) beyond which point the IFS is closer to the distribution.

After both forecasts have been quantile mapped, they are much closer to the ideal line, with deviations at high quantiles. The scale of sampling variability due to variability of samples within the test year was quanti ed by performing 1000 iterations of bootstrapping (see Sec. 2.5) to estimate the standard error of the quantiles. These are shown in Fig. 3 (a), where each error bar shows 2 standard errors. The quantile-mapped forecasts' extreme values are slightly larger than in the observations, which we attribute to sampling variability between the training and test periods.

In order to also get a sense of how extreme these quantile values are, we plot the
 approximate return period in days for a range of thresholds in Fig. 4, calculated over all
 hours in the test period. These are calculated as the average time gap between instances

{19{

Figure 5: Example precipitation forecasts following postprocessing by the cGAN-qm model, for a selection of hours throughout the primary test dataset (rst column). The examples are from randomly chosen dates, but Itered to ensure that a diverse range of months in the year are represented, and so that the examples show di erent behaviours for periods with high and low rainfall. The columns from left to right show: a single member of the cGAN-qm ensemble, the average of 20 cGAN-qm ensemble members, the IMERG observations, the IFS-qm forecast. Each row corresponds to one time value, shown in the title of the IMERG sample.



Figure 18: Fractions Skill Score in the extreme March{May 2018 season, over the Kenya subregion, for di erent quantile thresholds, with quantiles calculated for this season and subregion (a) 90 percentile (b) 99 percentile (c) 99 percentile (d) 99 percentile (e) 9999 percentile. Shading indicates 2 standard errors estimated from bootstrapping with 50 samples.

els to capture. The cGAN-qm model demonstrated a substantial correction to the timing of peak mean rainfall over the whole domain, which persisted when looking at the
diurnal cycle of high quantiles. However, there is substantial spatial noise in the cGANqm peak rainfall hour (Fig. 9). Using time as an input variable may be one way to improve the learnt relationships. The frequency distributions of rainfall for cGAN-qm and
IFS-qm (Fig. 3) were both comparable, with small biases which we attribute to sampling
variability.

The cGAN-qm improved forecast skill scores in some respects. The cGAN-qm shows generally higher Fractions Skill Scores (Fig. 12) up to a high percentil@<sup>th</sup>()? Particularly at larger neighbourhood sizes (abov OKm). For higher percentiles the IFSqm forecast demonstrated a higher score. The IFS-qm model also achieved higher ETS at the grid scale at all thresholds, although the scores were nevertheless quite low (Fig. 13).

Both models were also evaluated on the 2018 Long Rains, which were signi cantly 686 wetter than any Long Rains season seen in training and across the whole IMERG dataset. 687 It may be expected that machine learning-based methods would show degraded perfor-688 mance on situations outside their training data, and this is highly important to evalu-689 ate for forecasting applications (Watson, 2022). In fact, we found that unmodi ed cGAN 690 forecasts actually had a more realistic frequency distribution of rainfall than that of the 691 IFS (Fig. 16), and cGAN-qm forecasts had higher skill than IFS-qm when evaluated us-692 ing scatter plots and the FSS (Figs. 17 and 18). However, more evaluation work would 693 be required to have high con dence that the method will generally perform well in ex-694 treme situations. 695

An important advantage that generative machine learning models provide over other 696 non-generative models is the ability to create an ensemble of predictions from a single 697 forecast. It is therefore an interesting question as to whether this machine learning model 698 can provide well-calibrated probability distributions. Our assessment indicates that the 699 spread of the model correlates well with the observed error, although it also demonstrates 700 a mixture of under- and over-dispersive behaviour (Sec. 3.3). Note that the stochastic 701 component of the cGAN predictions is not temporally coherent, so that combining pre-702 dictions from di erent times would produce a time series with hour-to-hour variability 703 that is likely too large. This could be addressed in future work by postprocessing using 704 models like those applied in conditional video generation (e.g. Xing et al. (2023)). 705

{34{

Many studies on the performance of machine learning models compare the output 706 of a model to an unpostprocessed IFS forecast, or similar (e.g. Bi et al. (2022); Lam et 707 al. (2023)). Whilst the cGAN without quantile mapping improves on the unpostprocessed 708 IFS forecast rainfall distribution (Fig. 2), diurnal cycle, and skill scores (Fig. 12), by eval-709 uating our model against a strong baseline of quantile mapping our comparison reveals 710 what machine learning can do that is not achieved by conventional postprocessing meth-711 ods like quantile mapping, which is a sterner test than comparing to unprocessed fore-712 cast. 713

A strength of using machine learning to postprocess existing forecasts is that we 714 incorporate the skill and physical understanding of physics-derived models (Watson, 2019). 715 However, we implicitly assume that the IFS forecast captures all the useful information 716 about phenomena such as the Indian Ocean Dipole and Madden-Julian oscillations, and 717 our model may be improved by including indexes of these drivers and/or sea surface tem-718 peratures as additional inputs. It would be interesting as well to include the initial state 719 of the forecast, if available, to see whether the machine learning model is able to correct 720 for errors in how the IFS evolves this initial state. 721

There are also other state-of-the-art machine learning models that would be interesting to compare with, to see if performance improvements can be made. Particularly promising approaches would be di usion models, which have recently started to be applied in weather and climate prediction (Addison et al., 2022; Leinonen et al., 2023, e.g.), since they appear to perform well and are easier to train, and models trained directly on a suitable loss function such as the energy score (Pacchiardi et al., 2022).

## 728 Data Availability Statement

The code for the GAN and quantile mapping used in this paper is availabletaps:// github.com/bobbyantonio/downscaling-cgan; this code was forked from the code in L. Harris et al. (2022). All experiments in this paper were performed within TensorFlow 2.7.0, and some of the analysis in this work utilised the PySteps package (Pulkkinen et al., 2019). The ECMWF IFS forecasts can be obtained through MARS, for which academic accounts are freely available subject to conditionhtspee://www.ecmwf.int/ en/forecasts/accessing-forecasts/licences-available. The IMERG satellite pre-

{35{

cipitation data (Hu man et al., 2022) is freely available after registration, type //
 gpm.nasa.gov/data/directory.

## 738 Acknowledgments

This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (Grant No
741112, ITHACA). BA was supported by a NERC Cross-disciplinary Research for Environmental Science award. PW was supported by a NERC Independent Research Fellowship (Grant No. NE/S014713/1).

This work was carried out using the computational facilities of the Advanced Com puting Research Centre, University of Bristol - http://www.bristol.ac.uk/acrc/.

## 746 References

- Addison, H., Kendon, E., Ravuri, S., Aitchison, L., & Watson, P. A. (2022, November). Machine learning emulation of a local-scale UK climate model. Retrieved fromhttp://arxiv.org/abs/2211.16116 (arXiv: 2211.16116)
  Ageet, S., Fink, A. H., Maranan, M., Diem, J. E., Hartter, J., Ssali, A. L., & Ayabagabo, P. (2022, February). Validation of Satellite Rainfall Estimates
- <sup>752</sup> over Equatorial East Africa. *Journal of Hydrometeorology*, 23(2), 129(151.
- Retrieved 2023-09-21, filorops://journals.ametsoc.org/view/journals/
- hydr/23/2/JHM-D-21-0145.1.xml (Publisher: American Meteorological Soci-
- ety Section: Journal of Hydrometeorology) doi: 10.1175/JHM-D-21-0145.1
- Arjovsky, M., & Bottou, L. (2016, November). Towards Principled Methods for
   Training Generative Adversarial Networks. IProceedings of the 5th Interna tional Conference on Learning Representations. Retrieved 2023-07-05, from
   https://openreview.net/forum?id=Hk4\_qw5xe
- Arjovsky, M., Chintala, S., & Bottou, L. (2017, July). Wasserstein Generative Ad versarial Networks. In Proceedings of the 34th International Conference on Ma-
- <sup>762</sup> chine Learning (pp. 214{223). PMLR. Retrieved 2023-07-05, fitumps://
- 763 proceedings.mlr.press/v70/arjovsky17a.html (ISSN: 2640-3498)
- Bechtold, P., Chaboureau, J. P., Beljaars, A., Betts, A. K., Kohler, M., Miller, M.,
- <sup>765</sup> & Redelsperger, J. L. (2004, October). The simulation of the diurnal cycle of
   <sup>766</sup> convective precipitation over land in a global modeQuarterly Journal of the

767	Royal Meteorological Society, 130 C(604), 3119{3137. doi: 10.1256/qj.03.103			
768	Bechtold, P., Semane, N., Lopez, P., Chaboureau, J. P., Beljaars, A., & Bormann, N.			
769	(2014, February). Representing equilibrium and nonequilibrium convection in			
770	large-scale models Journal of the Atmospheric Sciences, 71(2), 734{753. doi:			
771	10.1175/JAS-D-13-0163.1			
772	Ben-Bouallegue, Z., Weyn, J. A., Clare, M. C. A., Dramsch, J., Dueben, P., &			
773	Chantry, M. (2023, August). Improving medium-range ensemble weather			
774	forecasts with hierarchical ensemble transformers. arXiv. Retrieved 2023-			
775	08-30, frommttp://arxiv.org/abs/2303.17195 (arXiv:2303.17195) doi:			
776	10.48550/arXiv.2303.17195			
777	Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2022, November). Pangu-			
778	Weather: A 3D High-Resolution Model for Fast and Accurate Global Weather			
779	Forecast. Retrieved from http://arxiv.org/abs/2211.02556 (arXiv:			
780	2211.02556)			
781	Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2023, July). Accurate			
782	medium-range global weather forecasting with 3D neural networks ature,			
783	619(7970), 533{538. Retrieved 2024-01-11hfitpen://www.nature.com/			
784	articles/s41586-023-06185-3 (Number: 7970 Publisher: Nature Publishing			
785	Group) doi: 10.1038/s41586-023-06185-3			
786	Bœ, J., Terray, L., Habets, F., & Martin, E. (2007, October). Statistical and			
787	dynamical downscaling of the Seine basin climate for hydro-meteorological			
788	studies. International Journal of Climatology, 27(12), 1643{1655. doi:			
789	10.1002/joc.1602			
790	Cafaro, C., Woodhams, B. J., Stein, T. H., Birch, C. E., Webster, S., Bain, C. L.,			
791	Hill, P. (2021, April). Do convection-permitting ensembles lead to more skillful			
792	short-range probabilistic rainfall forecasts over tropical east a friction $deta$			
793	and Forecasting, 36(2), 697{716. (Publisher: American Meteorological Society)			
794	doi: 10.1175/WAF-D-20-0172.1			
795	Camberlin, P., Gitau, W., Planchon, O., Dubreuil, V., Funatsu, B. M., & Philippon,			
796	N. (2018). Major role of water bodies on diurnal precipitation regimes in East-			
797	ern Africa. International Journal of Climatology, 38(2), 613(629. Retrieved			
798	2023-06-15, fromtps://onlinelibrary.wiley.com/doi/abs/10.1002/			
799	joc.5197 (_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/joc.5197)			

800	doi: 10.1002/joc.5197					
801	Chamberlain, J. M., Bain, C. L., Boyd, D. F., Mccourt, K., Butcher, T., & Palmer,					
802	S. (2014). Forecasting storms over Lake Victoria using a high resolution model.					
803	Meteorological Applications, 21(2), 419{430. (Publisher: John Wiley and Sons					
804	Ltd) doi: 10.1002/met.1403					
805	Clark, P., Roberts, N., Lean, H., Ballard, S. P., & Charlton-Perez, C. (2016).					
806	Convection-permitting models: a step-change in rainfall forecasting. $Me$ -					
807	teorological Applications, 23(2), 165{181. Retrieved 2024-01-23, from					
808	https://onlinelibrary.wiley.com/doi/abs/10.1002/met.1538 (_eprint:					
809	https://onlinelibrary.wiley.com/doi/pdf/10.1002/met.1538) doi: 10.1002/					
810	met.1538					
811	Creswell, A., White, T., Dumoulin, V., Arulkumaran, K., Sengupta, B., & Bharath,					
812	A. A. (2018, January). Generative Adversarial Networks: An Overview $EEE$					
813	Signal Processing Magazine, 35(1), 53(65. (Conference Name: IEEE Signal					
814	Processing Magazine) doi: 10.1109/MSP.2017.2765202					
815	Dezfuli, A. K., Ichoku, C. M., Hu man, G. J., Mohr, K. I., Selker, J. S., De Giesen,					
816	N. V., Annor, F. O. (2017). Validation of IMERG Precipitation in					
817	Africa. Journal of Hydrometeorology, 18(10), 2817{2825. Retrieved from					
818	https://doi.org/10.1175/ doi: 10.1175/JHM-D-17-0139.s1					
819	Dinku, T., Connor, S. J., & Ceccato, P. (2010). Comparison of CMORPH and					
820	TRMM-3B42 over Mountainous Regions of Africa and South America. In					
821	M. Gebremichael & F. Hossain (Eds.) Satellite Rainfall Applications for Sur-					
822	face Hydrology (pp. 193(204). Dordrecht: Springer Netherlands. Retrieved					
823	2023-09-21, fromtps://doi.org/10.1007/978-90-481-2915-7_11 doi:					
824	10.1007/978-90-481-29115-7					
825	Duncan, J., Subramanian, S., & Harrington, P. (2022) Generative Modeling of High-					
826	resolution Global Precipitation Forecasts. (arXiv:2210.12504)					
827	Deque, M. (2007, May). Frequency of precipitation and temperature extremes over					
828	France in an anthropogenic scenario: Model results and statistical correction					
829	according to observed values. Global and Planetary Change, 57(1), 16{26.					
830	Retrieved 2023-06-16, fitureps://www.sciencedirect.com/science/					
831	article/pii/S0921818106002748 doi: 10.1016/j.gloplacha.2006.11.030					
832	Ebert, E. E. (2008). Fuzzy veri cation of high-resolution gridded forecasts: A review					

{38{

833	and proposed framework. Meteorological Applications (Vol. 15, pp. 51(64).			
834	John Wiley and Sons Ltd. (Issue: 1 ISSN: 14698080) doi: 10.1002/met.25			
835	ECMWF. (2023, February). Changes to the forecasting system - Forecast User -			
836	ECMWF Confluence Wiki. Retrieved 2023-11-02, from ps://confluence			
837	$. \verb+ecmwf.int/display/FCST/Changes+to+the+forecasting+system+$			
838	(https://con uence.ecmwf.int/display/FCST/Changes+to+the+forecasting+system,			
839	Accessed 2nd November 2023)			
840	ECMWF. (2023, February). Operational configurations of the ECMWF Inte-			
841	grated Forecasting System (IFS). Retrieved 2023-01-09, filorops://			
842	<pre>confluence.ecmwf.int/pages/viewpage.action?pageId=324860211</pre>			
843	(https://con uence.ecmwf.int/pages/viewpage.action?pageId=324860211.			
844	Accessed 1st September)			
845	ECMWF. (2023, February). Parameter database. Retrieved			
846	2023-09-20, from ps://codes.ecmwf.int/grib/param-db/			
847	(https://codes.ecmwf.int/grib/param-db/, Accessed February 2023)			
848	Efron, B., & Tibshirani, R. (1986). Bootstrap Methods for Standard Errors, Con-			
849	dence Intervals, and Other Measures of Statistical Accuracy. Statistical			
850	Science, 1(1), 54{75. Retrieved 2023-11-02, fromps://www.jstor.org/			
851	stable/2245500 (Publisher: Institute of Mathematical Statistics)			
852	Finney, D. L., Marsham, J. H., Jackson, L. S., Kendon, E. J., Rowell, D. P., Boor-			
853	man, P. M., Senior, C. A. (2019, April). Implications of Improved Repre-			
854	sentation of Convection for the East Africa Water Budget Using a Convection-			
855	Permitting Model. Journal of Climate, 32(7), 2109{2129. Retrieved 2023-			
856	O7-27, from ttps://journals.ametsoc.org/view/journals/clim/32/7/			
857	jcli-d-18-0387.1.xml (Publisher: American Meteorological Society Section:			
858	Journal of Climate) doi: 10.1175/JCLI-D-18-0387.1			
859	Floodlist. (2023, May). Somalia – Flooding Displaces Thousands, Prompts Urgent			
860	Humanitarian Response. Retrieved 2023-08-04, from ps://floodlist.com/			
861	africa/somalia-floods-may-2023 (https:// oodlist.com/africa/somalia-			
862	oods-may-2023. Accessed May 2023.)			
863	Gebremeskel Haile, G., Tang, Q., Sun, S., Huang, Z., Zhang, X., & Liu, X. (2019,			
864	June). Droughts in East Africa: Causes, impacts and resilience. Earth-			
865	Science Reviews, 193, 146{161. Retrieved 2023-08-03, fittages://			

866	www.sciencedirect.com/science/article/pii/S0012825218303519 doi:
867	10.1016/j.earscirev.2019.04.015
868	Gilleland, E., Ahijevych, D., Brown, B. G., Casati, B., & Ebert, E. E. (2009, Octo-
869	ber). Intercomparison of spatial forecast veri cation methods ather and
870	Forecasting, 24(5), 1416{1430. doi: 10.1175/2009WAF2222269.1
871	Gneiting, T., & Raftery, A. E. (2007, March). Strictly proper scoring rules, pre-
872	diction, and estimation. Journal of the American Statistical Association,
873	102(477), 359{378. doi: 10.1198/016214506000001437
874	Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.
875	Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S.,
876	Bengio, Y. (2014). Generative Adversarial Nets. Indvances in Neural
877	Information Processing Systems (Vol. 27). Curran Associates, Inc. Retrieved
878	2023-07-05, fr <b>om</b> ps://proceedings.neurips.cc/paper_files/paper/
879	2014/hash/5ca3e9b122f61f8f06494c97b1afccf3-Abstract.html
880	Gudmundsson, L., Bremnes, J. B., Haugen, J. E., & Skaugen, T. E. (2012). Quantile
881	mapping Hydrology and Earth System Sciences Discussions Technical Note:
882	Downscaling RCM precipitation to the station scale using quantile mapping-a
883	comparison of methods Quantile mapping. Hydrol. Earth Syst. Sci. Discuss,
884	9,6185{6201. Retrieved fromm.hydrol-earth-syst-sci-discuss.net/9/
885	6185/2012/ doi: 10.5194/hessd-9-6185-2012
886	Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., & Courville, A. C. (2017).
887	Improved Training of Wasserstein GANs. Indvances in Neural Information
888	Processing Systems (Vol. 30). Curran Associates, Inc. Retrieved 2023-07-05,
889	<pre>fromhttps://proceedings.neurips.cc/paper_files/paper/2017/hash/</pre>
890	892c3b1c6dccd52936e27cbd0ff683d6-Abstract.html
891	Haiden, T., Rodwell, M. J., Richardson, D. S., Okagaki, A., Robinson, T., & Hew-
892	son, T. (2012, August). Intercomparison of global model precipitation forecast
893	skill in 2010/11 using the SEEPS score. Monthly Weather Review, 140(8),
894	2720{2733. doi: 10.1175/MWR-D-11-00301.1
895	Hamill, T. M. (2001, March). Interpretation of Rank Histograms for Verify-
896	ing Ensemble Forecasts. Monthly Weather Review, 129(3), 550{560. Re-
897	trieved 2023-09-01, fitureps://journals.ametsoc.org/view/journals/
898	mwre/129/3/1520-0493_2001_129_0550_iorhfv_2.0.co_2.xml (Publisher:

{40{

899	American Meteorological Society Section: Monthly Weather Review) doi:			
900	10.1175/1520-0493(200 <b>11)529</b> :IORHFVi2.0.CO;2			
901	Harris, D., Foufoula-Georgiou, E., Droegemeier, K. K., & Levit, J. J. (2001, Au-			
902	gust). Multiscale Statistical Properties of a High-Resolution Precipitation			
903	Forecast. Journal of Hydrometeorology, 2(4), 406{418. Retrieved 2023-			
904	O8-O3, frommttps://journals.ametsoc.org/view/journals/hydr/2/			
905	4/1525-7541_2001_002_0406_mspoah_2_0_co_2.xml (Publisher: Amer-			
906	ican Meteorological Society Section: Journal of Hydrometeorology) doi:			
907	10.1175/1525-7541(200h0@2:MSPOAH2.0.CO;2			
908	Harris, L., McRae, A. T. T., Chantry, M., Dueben, P. D., & Palmer, T. N. (2022,			
909	April). A Generative Deep Learning Approach to Stochastic Downscaling			
910	of Precipitation Forecasts. Journal of Advances in Modeling Earth Systems,			
911	<i>14</i> (10).			
912	Hu man, G., Bolvin, D., Braithwaite, D., Hsu, K., Joyce, R., & Xie, P. (2022,			
913	October). Integrated multi-satellite retrievals for GPM (IMERG), VO6B.			
914	[dataset]. NASA's Precipitation Processing Center, accessed 1st October 2022,			
915	https://arthurhouhttps.pps.eos dis.nasa.gov/text/gpmallversions/V06/YYYY/MM/DD/imerg/Without the state of t			
916				
917	IFRC. (2014). World Disasters Report, Focus on culture and risk (Tech.			
918	Rep.). International Federation of Red Cross and Red Crescent Soci-			
919	eties. Retrieved 2023-07-27, from s://www.ifrc.org/document/			
920	world-disasters-report-2014			
921	Jeong, C. H., & Yi, M. Y. (2023, February). Correcting rainfall forecasts of a numer-			
922	ical weather prediction model using generative adversarial networksal of			
923	Supercomputing, 79(2), 1289{1317. (Publisher: Springer) doi: 10.1007/s11227			
924	-022-04686-y			
925	Karras, T., Laine, S., & Aila, T. (2019, June). A Style-Based Generator Archi-			
926	tecture for Generative Adversarial Networks. Phaceedings of the IEEE/CVF			
927	Conference on Computer Vision and Pattern Recognition (CVPR).			
928	Kilavi, M., MacLeod, D., Ambani, M., Robbins, J., Dankers, R., Graham, R.,			
929	Todd, M. C. (2018, November). Extreme rainfall and ooding over Central			
930	Kenya Including Nairobi City during the long-rains season 2018: Causes, pre-			
931	dictability, and potential for early warning and actions. Atmosphere, 9(12).			

932	(Publisher: MDPI AG) doi: 10.3390/atmos9120472			
933	Kim, J. E., & Joan Alexander, M. (2013, May). Tropical precipitation variabil-			
934	ity and convectively coupled equatorial waves on submonthly time scales			
935	in reanalyses and TRMM. Journal of Climate, 26(10), 3013{3030. doi:			
936	10.1175/JCLI-D-12-00353.1			
937	Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M.,			
938	Alet, F., Battaglia, P. (2023, November). Learning skillful medium-			
939	range global weather forecasting Science, $\theta$ (0), eadi2336. Retrieved 2023-			
940	11-22, from https://www.science.org/doi/10.1126/science.adi2336			
941	(Publisher: American Association for the Advancement of Science) doi:			
942	10.1126/science.adi2336			
943	Leinonen, J., Hamann, U., Nerini, D., Germann, U., & Franch, G. (2023, April). La-			
944	tent diffusion models for generative precipitation nowcasting with accurate un-			
945	certainty quantification. Retrieved from http://arxiv.org/abs/2304.12891			
946	(arXiv: 2304.12891)			
947	Leinonen, J., Nerini, D., & Berne, A. (2020, November). Stochastic Super-Resolution			
948	for Downscaling Time-Evolving Atmospheric Fields With a Generative Adver-			
949	sarial Network. IEEE Transactions on Geoscience and Remote Sensing, 59(9),			
950	7211{7223. (arXiv: 2005.10374 Publisher: Institute of Electrical and Electron-			
951	ics Engineers (IEEE)) doi: 10.1109/tgrs.2020.3032790			
952	Leutbecher, M., & Palmer, T. N. (2008, March). Ensemble forecastidgurnal of			
953	Computational Physics, 227(7), 3515{3539. (Publisher: Academic Press Inc.)			
954	doi: 10.1016/j.jcp.2007.02.014			
955	MacLeod, D., Kilavi, M., Mwangi, E., Ambani, M., Osunga, M., Robbins, J.,			
956	Todd, M. C. (2021, January). Are Kenya Meteorological Department heavy			
957	rainfall advisories useful for forecast-based early action and early preparedness			
958	for ooding? Natural Hazards and Earth System Sciences, 21(1), 261{277.			
959	Retrieved 2023-06-15, filontps://nhess.copernicus.org/articles/21/			
960	261/2021/ (Publisher: Copernicus GmbH) doi: 10.5194/nhess-21-261-2021			
961	Macleod, D. A., Dankers, R., Graham, R., Guigma, K., Jenkins, L., Todd, M. C.,			
962	Mwangi, E. (2021). Drivers and Subseasonal Predictability of Heavy Rainfall			
963	in Equatorial East Africa and Relationship with Flood Risk. Journal of Hy-			
964	drometeorology, 22(4), 887{903. Retrieved from tps://doi.org/10.1175/			

965	JHM-D-20- doi: 10.1175/JHM-D-20				
966	Manzato, A., & Jolli e, I. (2017, April). Behaviour of veri cation measures for de-				
967	terministic binary forecasts with respect to random changes and thresholding.				
968	Quarterly Journal of the Royal Meteorological Society, 143(705), 1903(1915.				
969	(Publisher: John Wiley and Sons Ltd) doi: 10.1002/qj.3050				
970	Maraun, D., & Widmann, M. (2017, December). Model Output Statistics. Sta				
971	tistical Downscaling and Bias Correction for Climate Research (pp. 170{200).				
972	Cambridge University Press. doi: 10.1017/9781107588783.013				
973	Marsham, J. H., Dixon, N. S., Garcia-Carreras, L., Lister, G. M., Parker, D. J.,				
974	Knippertz, P., & Birch, C. E. (2013, May). The role of moist convection in				
975	the West African monsoon system: Insights from continental-scale convection-				
976	permitting simulations. Geophysical Research Letters, 40(9), 1843(1849. doi:				
977	10.1002/grl.50347				
978	Mittermaier, M., Roberts, N., & Thompson, S. A. (2013). A long-term assessment of				
979	precipitation forecast skill using the Fractions Skill Schlieteorological Appli-				
980	cations, 20(2), 176{186. (Publisher: John Wiley and Sons Ltd) doi: 10.1002/				
981	met.296				
982	Nguyen, T., Brandstetter, J., Kapoor, A., Gupta, J. K., & Grover, A. (2023, Jan-				
983	uary). ClimaX: A foundation model for weather and climate. Retrieved from				
984	http://arxiv.org/abs/2301.10343 (arXiv: 2301.10343)				
985	Nicholson, S. (2016). The Turkana low-level jet: mean climatology and				
986	association with regional aridity. International Journal of Clima-				
987	tology, 36(6), 2598{2614. Retrieved 2023-08-08, frpsn//				
988	onlinelibrary.wiley.com/doi/abs/10.1002/joc.4515 (_eprint:				
989	https://onlinelibrary.wiley.com/doi/pdf/10.1002/joc.4515) doi: 10.1002/				
990	joc.4515				
991	Nicholson, S. E. (2017). Climate and climatic variability of rainfall over eastern				
992	Africa. Reviews of Geophysics, 55(3), 590(635.				
993	Pacchiardi, L., Adewoyin, R., Dueben, P., & Dutta, R. (2022, May). Probabilis-				
994	tic Forecasting with Generative Networks via Scoring Rule Minimization.				
995	arXiv. Retrieved 2023-07-20, fromp://arxiv.org/abs/2112.08217				
996	(arXiv:2112.08217 [cs, stat])				
997	Palmer, P. I., Wainwright, C. M., Dong, B., Maidment, R. I., Wheeler, K. G., Ged-				

{43{

998	ney, N., Turner, A. G. (2023, April). Drivers and impacts of Eastern					
999	African rainfall variability. Nature Reviews Earth & Environment, 4(4),					
1000	254{270. Retrieved 2023-09-01, httpps://www.nature.com/articles/					
1001	s43017-023-00397-x (Number: 4 Publisher: Nature Publishing Group) doi:					
1002	10.1038/s43017-023-00397-x					
1003	Price, I., & Rasp, S. (2022). Increasing the accuracy and resolution of precipitation					
1004	forecasts using deep generative model Proceedings of the 25th International					
1005	Conference on Artificial Intelligence and Statistics (AISTATS), 151. doi:					
1006	https://doi.org/10.48550/arXiv.2203.12297					
1007	Pulkkinen, S., Nerini, D., Perez Hortal, A. A., Velasco-Forero, C., Seed, A., Ger-					
1008	mann, U., & Foresti, L. (2019, October). Pysteps: an open-source Python					
1009	library for probabilistic precipitation nowcasting (v1.0) [softv@ave3_ientific					
1010	Model Development, 12(10), 4185{4219. doi: 10.5194/gmd-12-4185-2019					
1011	Rasp, S., & Lerch, S. (2018). Neural Networks for Postprocessing Ensemble Weather					
1012	Forecasts. Monthly Weather Review, 146(11), 3885{3900. Retrieved from					
1013	https://doi.org/10.1175/MWR-D-18- doi: 10.1175/MWR-D-18					
1014	Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., Mo-					
1015	hamed, S. (2021, September). Skilful precipitation nowcasting using deep					
1016	generative models of radar: Supplementary InformationNature, 597(7878),					
1017	672{677. Retrieved 2023-07-06, httops://www.nature.com/articles/					
1018	s41586-021-03854-z doi: 10.1038/s41586-021-03854-z					
1019	Roberts, N. (2008). Assessing the spatial and temporal variation in the skill of					
1020	precipitation forecasts from an NWP model. Mateorological Applications					
1021	(Vol. 15, pp. 163{169). John Wiley and Sons Ltd. (Issue: 1 ISSN: 14698080)					
1022	doi: 10.1002/met.57					
1023	Roberts, N. M., & Lean, H. W. (2008, January). Scale-selective veri cation of rain-					
1024	fall accumulations from high-resolution forecasts of convective eMem $tbly$					
1025	Weather Review, 136(1), 78{97. doi: 10.1175/2007MWR2123.1					
1026	Roca, R., Chambon, P., Jobard, I., Kirstetter, P. E., Gosset, M., & Berges, J. C.					
1027	(2010). Comparing satellite and surface rainfall products over West Africa					
1028	at meteorologically relevant scales during the AMMA campaign using error					
1029	estimates. Journal of Applied Meteorology and Climatology, 49(4), 715(731.					
1030	(Publisher: American Meteorological Society) doi: 10.1175/2009JAMC2318.1					

1031	Schaefer, J. T. (1990, December). The Critical Success Index as an Indi-
1032	cator of Warning Skill. Weather and Forecasting, 5(4), 570{575. Re-
1033	trieved 2023-07-03, fi <b>tom</b> ps://journals.ametsoc.org/view/journals/
1034	wefo/5/4/1520-0434_1990_005_0570_tcsiaa_2_0_co_2.xml (Publisher:
1035	American Meteorological Society Section: Weather and Forecasting) doi:
1036	10.1175/1520-0434(199 <b>b)605</b> :TCSIAAi 2.0.CO;2
1037	Sinclair, S., & Pegram, G. G. S. (2005, July). Empirical Mode Decomposition in 2-
1038	D space and time: a tool for space-time rainfall analysis and nowcasting.
1039	drology and Earth System Sciences, 9(3), 127{137. Retrieved 2023-08-03, from
1040	https://hess.copernicus.org/articles/9/127/2005/ (Publisher: Coper-
1041	nicus GmbH) doi: 10.5194/hess-9-127-2005
1042	Vogel, P., Knippertz, P., Fink, A. H., & Schlueter, A. (2018). Skill of Global
1043	Raw and Postprocessed Ensemble Predictions of Rainfall over Northern
1044	Tropical Africa. Weather and Forecasting, 33(2), 369(388. Retrieved from
1045	https://doi.org/10.1175/WAF-D-17- doi: 10.1175/WAF-D-17
1046	Vogel, P., Knippertz, P., Fink, A. H., Schlueter, A., & Gneiting, T. (2020, Decem-
1047	ber). Skill of Global Raw and Postprocessed Ensemble Predictions of Rainfall
1048	in the Tropics. Weather and Forecasting, 35(6), 2367(2385. Retrieved 2024-
1049	O1-19, from ttps://journals.ametsoc.org/view/journals/wefo/35/6/
1050	WAF-D-20-0082.1.xml (Publisher: American Meteorological Society Section:
1051	Weather and Forecasting) doi: 10.1175/WAF-D-20-0082.1
1052	Vosper, E., Watson, P., Harris, L., McRae, A., Santos-Rodriguez, R., Aitchison,
1053	L., & Mitchell, D. (2023). Deep Learning for Downscaling Tropical Cy-
1054	clone Rainfall to Hazard-Relevant Spatial Scales. Journal of Geophysical
1055	Research: Atmospheres, 128(10), e2022JD038163. Retrieved 2023-10-27,
1056	<pre>fromhttps://onlinelibrary.wiley.com/doi/abs/10.1029/2022JD038163</pre>
1057	(_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1029/2022JD038163) doi
1058	10.1029/2022JD038163
1059	Wainwright, C. M., Finney, D. L., Kilavi, M., Black, E., & Marsham, J. H. (2021,
1060	January). Extreme rainfall in East Africa, October 2019{January 2020 and
1061	context under future climate change. Weather, 76(1), 26{31. Retrieved
1062	<pre>fromhttps://onlinelibrary.wiley.com/doi/10.1002/wea.3824 doi:</pre>
1063	10.1002/wea.3824

{45{

1064	Walker, D. P., Birch, C. E., Marsham, J. H., Scaife, A. A., Graham, R. J., & Segele,				
1065	Z. T. (2019, October). Skill of dynamical and GHACOF consensus seasonal				
1066	forecasts of East African rainfall. Climate Dynamics, 53(7-8), 4911{4935.				
1067	(Publisher: Springer Verlag) doi: 10.1007/s00382-019-04835-9				
1068	Warner, J. L., Petch, J., Short, C. J., & Bain, C. (2023). Assessing the impact of				
1069	a NWP warm-start system on model spin-up over tropical AfricaQuarterly				
1070	Journal of the Royal Meteorological Society, 149(751), 621(636. Retrieved				
1071	2023-11-02, from tps://onlinelibrary.wiley.com/doi/abs/10.1002/				
1072	qj.4429 (_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/qj.4429)				
1073	doi: 10.1002/qj.4429				
1074	Watkiss, P., & Cimato, F. (2021, September). Socio-Economic Benefits of				
1075	the WISER Programme. Synthesis of Results. (Tech. Rep.). Met Of-				
1076	ce UK. Retrieved from https://www.metoffice.gov.uk/binaries/				
1077	content/assets/metofficegovuk/pdf/business/international/wiser/				
1078	wiser-seb-results_final-web.pdf				
1079	Watkiss, P., Powell, R., Hunt, A., & Cimato, F. (2020, September). The Socio-				
1080	Economic Benefits of the HIGHWAY project (Tech. Rep.). Weather and Cli-				
1081	mate Information Services for Africa (WISER).				
1082	Watson, P. A. G. (2019). Applying Machine Learning to Improve Simulations of				
1083	a Chaotic Dynamical System Using Empirical Error Correction. Journal of				
1084	Advances in Modeling Earth Systems, 11(5), 1402{1417. Retrieved 2023-09-14,				
1085	<pre>fromhttps://onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001597</pre>				
1086	(_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1029/2018MS001597) doi:				
1087	10.1029/2018MS001597				
1088	Watson, P. A. G. (2022, November). Machine learning applications for weather				
1089	and climate need greater focus on extremeds avironmental Research Letters,				
1090	17(11), 111004. Retrieved 2023-08-04, https://dx.doi.org/10.1088/				
1091	1748-9326/ac9d4e (Publisher: IOP Publishing) doi: 10.1088/1748-9326/				
1092	ac9d4e				
1093	Wilks, D. S. (2019a). Forecast Veri cation. In the Atmo-				
1094	spheric Sciences (pp. 369{483). Elsevier. doi: 10.1016/b978-0-12-815823-4				
1095	.00009-2				
1096	Wilks, D. S. (2019b). Statistical Forecasting. Statistical Methods in the At-				

1097	mospheric Sciences (pp. 235{312). Elsevier. doi: 10.1016/b978-0-12-815823-4			
1098	.00007-9			
1099	Woodhams, B. J., Birch, C. E., Marsham, J. H., Bain, C. L., Roberts, N. M., &			
1100	Boyd, D. F. (2018, September). What is the added value of a convection-			
1101	permitting model for forecasting extreme rainfall over tropical East Africa?			
1102	Monthly Weather Review, 146(9), 2757{2780. (Publisher: American Meteoro-			
1103	logical Society) doi: 10.1175/MWR-D-17-0396.1			
1104	Woodhams, B. J., Birch, C. E., Marsham, J. H., Lane, T. P., Bain, C. L., & Web-			
1105	ster, S. (2019). Identifying key controls on storm formation over the lake			
1106	Victoria basin. Monthly Weather Review, 147(9), 3365{3390. (Publisher:			
1107	American Meteorological Society) doi: 10.1175/MWR-D-19-0069.1			
1108	Xing, Z., Feng, Q., Chen, H., Dai, Q., Hu, H., Xu, H., Jiang, YG. (2023). A			
1109	survey on video di usion modelsarXiv preprint arXiv:2310.10647.			
1110	Youds, L. H., Parker, D. J., Ade san, E. A., Antwi-Agyei, P., Bain, C. L., Black,			
1111	E. C. L., Woolnough, S. J. (2021, November). GCRF African SWIFT			
1112	and ForPAc SHEAR White Paper on the Potential of Operational Weather			
1113	Prediction to Save Lives and Improve Livelihoods and Economies in Sub-			
1114	Saharan Africa [Monograph]. Retrieved 2023-07-04, fitops://			
1115	eprints.whiterose.ac.uk/181045/ (Publisher: University of Leeds)			
1116	Zhang, Y., Long, M., Chen, K., Xing, L., Jin, R., Jordan, M. I., & Wang, J. (2023,			
1117	July). Skilful nowcasting of extreme precipitation with NowcastMature,			
1118	1{7. Retrieved 2023-07-14, fnormps://www.nature.com/articles/s41586			
1119	-023-06184-4 (Publisher: Nature Publishing Group) doi: 10.1038/s41586-023			
1120	-06184-4			
1121	Appendix A Forecast variables used			
1122	The IFS variables and constant elds used to train the model are shown in Table A1,			
1123	de nitions taken from (ECMWF, 2023).			
1124	The preprocessing methods mentioned in the table are as follows, using the year			
1125	2017 as the reference period:			

Minmax: calculate the minimum  $d_{min}$  and maximum  $d_{max}$  over the reference period, and then transform each value according to  $(d_{min})/(d_{max} - d_{min})$ .

- Max: calculate the and maximum over the reference period, and then trans-
- form each value according to  $d_{max}$ .
- Log: Transform each value according to  $\log(1 + v)$ .

Variable name	Symbol	Pre-processing applied
2 metre temperature	2t	Minmax
Convective available potential energy	cape	Log
Convective inhibition	cin	Max
Convective precipitation	ср	Log
Surface pressure	sp	Minmax
Total column cloud liquid water	tclw	Max
Total column vertically-integrated water vapo	ourtcwv	Log
Top of atmosphere incident solar radiation	tisr	Max
Total precipitation	tp	Log
Relative humidity at 200hPa	r200	Max
Relative humidity at 700hPa	r700	Max
Relative humidity at 950hPa	r950	Max
Temperature at 200hPa	t200	Minmax
Temperature at 700hPa	t700	Minmax
Eastward component of wind at 200hPa	u200	Max
Eastward component of wind at 700hPa	u700	Max
Northward component of wind at 200hPa	v200	Max
Northward component of wind at 700hPa	v700	Max
Vertical velocity at 200hPa	w200	Max
Vertical velocity at 500hPa	w500	Max
Vertical velocity at 700hPa	w700	Max
Orography	h	Max
Land-sea mask	Ism	N/A

Table A1: IFS variables used to train the model, as well as the normalisation applied to each variable. See text for description of the di erent preprocessing types.

