

Spatial scale evaluation of forecast flood inundation maps

Article

Published Version

Creative Commons: Attribution 4.0 (CC-BY)

open access

Hooker, H. ORCID: <https://orcid.org/0000-0002-5135-3952>,
Dance, S. L. ORCID: <https://orcid.org/0000-0003-1690-3338>,
Mason, D. C. ORCID: <https://orcid.org/0000-0001-6092-6081>,
Bevington, J. and Shelton, K. (2022) Spatial scale evaluation
of forecast flood inundation maps. *Journal of Hydrology*, 612
(Part B). 128170. ISSN 0022-1694 doi:
<https://doi.org/10.1016/j.jhydrol.2022.128170> Available at
<https://centaur.reading.ac.uk/106342/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1016/j.jhydrol.2022.128170>

Publisher: Elsevier

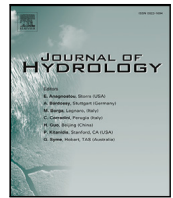
All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

www.reading.ac.uk/centaur

CentAUR

Central Archive at the University of Reading

Reading's research outputs online



Research papers

Spatial scale evaluation of forecast flood inundation maps

Helen Hooker^{a,*}, Sarah L. Dance^{a,b,d}, David C. Mason^c, John Bevington^e, Kay Shelton^e^a Department of Meteorology, University of Reading, UK^b Department of Mathematics and Statistics, University of Reading, UK^c Department of Geography and Environmental Science, University of Reading, UK^d National Centre for Earth Observation (NCEO), University of Reading, UK^e Jeremy Benn Associates Limited (JBA Consulting), Skipton, UK

ARTICLE INFO

This manuscript was handled by Emmanouil Anagnostou, Editor-in-Chief, with the assistance of Guy Schumann, Associate Editor.

Keywords:

Flood maps
Spatial verification
Scale selective
SAR

ABSTRACT

Flood inundation forecast maps provide an essential tool to disaster management teams for planning and preparation ahead of a flood event in order to mitigate the impacts of flooding on the community. Evaluating the accuracy of forecast flood maps is essential for model development and improving future flood predictions. Conventional, quantitative binary verification measures typically provide a domain-averaged score, at grid level, of forecast skill. This score is dependent on the magnitude of the flood and the spatial scale of the flood map. Binary scores have limited physical meaning and do not indicate location-specific variations in forecast skill that enable targeted model improvements to be made. A new, scale-selective approach is presented here to evaluate forecast flood inundation maps against remotely observed flood extents. A neighbourhood approach based on the Fraction Skill Score is applied to assess the spatial scale at which the forecast becomes skilful at capturing the observed flood. This skilful scale varies with location and when combined with a contingency map creates a novel categorical scale map, a valuable visual tool for model evaluation and development. The impact of model improvements on forecast flood map accuracy skill scores are often masked by large areas of correctly predicted flooded/unflooded cells. To address this, the accuracy of the flood-edge location is evaluated. The flood-edge location accuracy proves to be more sensitive to variations in forecast skill and spatial scale compared to the accuracy of the entire flood extent. Additionally, the resulting skilful scale of the flood-edge provides a physically meaningful verification measure of the forecast flood-edge discrepancy. The methods are illustrated by application to a case study flood event (with an estimated return period of 120 to 550 years) of the River Wye and River Lugg (UK) in February 2020.

Representation errors are introduced where remote sensing observations capture flood extent at different spatial resolutions in comparison with the model. The sensitivity of the verified skilful scale to the resolution of the observations is investigated. Re-scaling and interpolating observations leads to a small reduction in skill score compared with the observation flood map derived at the model resolution. The domain-averaged skilful scale remains the same with slight location-specific variations in skilful scale evident on the categorical scale map. Overall, our novel emphasis on scale, rather than domain-average score, means that comparisons can be made across different flooding scenarios and forecast systems and between forecasts at different spatial scales.

1. Introduction

Timely predictions of flood extent and depth from flood forecasting systems provide essential information to flood risk managers that enable anticipatory action prior to the occurrence of a potential flooding event. Evaluating the accuracy of flood extent forecasts against observations forms an essential part of model development (Schumann, 2019). Forecast flood inundation footprints are typically validated against remote sensing images using binary performance measures (Stephens et al., 2014) calculated at grid level.

In order to produce a forecast flood map, hydrodynamic or hydraulic flood models in two-dimensions simulate the flow of water using a local digital terrain model (DTM). The spatial resolution of DTMs has increased over recent years and is important for accurate flood mapping. For example, in the UK, the Environment Agency National LIDAR Programme offers open source 1 m surface elevation data for the whole of England (Environment Agency, 2021). Additional surface detail to 0.3 m spatial resolution from unmanned aerial vehicle UAV-LIDAR data acquired in urban areas is now possible (Trepecki

* Corresponding author.

E-mail address: h.hooker@pgr.reading.ac.uk (H. Hooker).<https://doi.org/10.1016/j.jhydrol.2022.128170>

Received 18 February 2022; Received in revised form 17 May 2022; Accepted 2 July 2022

Available online 14 July 2022

0022-1694/© 2022 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

et al., 2021). This means forecast flood maps could be presented at this very high resolution. It is questionable how meaningful it is to present highly detailed flood maps as a deterministic forecast (Savage et al., 2016), particularly at longer lead times where the skill of the flood forecasting system becomes increasingly dependent on the accuracy of the meteorological forecast (ECMWF, 2022). Speight et al. (2021) note for surface water flooding that more detail is included in local scale flood maps than can be justified by the predictability of the forecast. A high resolution, fine scale forecast flood map will show greater detail of the flood extent and the flood-edge location compared to a low resolution, coarse scale flood map. At a high resolution the discrepancy between the forecast and observed flood maps may be closer in terms of distance, however a small mismatch will lead to a double penalty impact on forecast verification. The model is penalised twice for the over-prediction (false alarm) and the under-prediction (miss) (Stein and Stoop, 2019). When high resolution forecasts are verified against observations at grid level, the predictability can appear to worsen and the high resolution forecast would need to perform better than the low resolution forecast to achieve the same verification score. It is not meaningful to compare verification scores across different spatial scales. Spatial verification methods for flood inundation mapping have only received limited attention over the past decade (Schumann, 2019).

Verification approaches that account for uncertainties in observations and small discrepancies in gridded data using a fuzzy set approach (Hagen, 2003) have previously been applied to flood mapping (Pappenberger et al., 2007; Dasgupta et al., 2018). However, the fuzzy set method does not incorporate variations in spatial scale (Croke and Pappenberger, 2008). In atmospheric sciences, verification approaches that account for changes in spatial scale are well established. These approaches include the Fraction Skill Score (FSS), which applies a neighbourhood approach to assess a useful/skilful scale (Roberts and Lean, 2008) of a precipitation forecast. Dey et al. (2014, 2016) developed the FSS approach to produce location-specific agreement scales between the forecast and observed fields to understand the spatial predictability of an ensemble forecast. Other spatial scale approaches include the wavelet method of scale decomposition, where the forecast and observed fields are decomposed into maps at different scales by wavelet transformation and subsequently verified (Briggs and Levine, 1997; Casati and Wilson, 2007). Croke and Pappenberger (2008) note that this method is extremely sensitive to offsetting of maps.

In general, the performance of forecast flood maps are evaluated for the entire flood extent, regardless of flood magnitude, adding bias to binary performance measures (Stephens et al., 2014). Stephens et al. (2014) question whether it is important to validate all flooded cells, when only cells that are close to the flood margin are difficult to predict. Pappenberger et al. (2007) evaluated model performance only on cells that were subject to change between differing model runs to address the issue of large areas of correctly predicted flooded/unflooded cells masking variations in forecast skill scores.

Satellite based Synthetic Aperture Radar (SAR) sensors are well known for their flood detection capability. Unobstructed flood waters appear dark on SAR images due to the low backscatter return from the relatively smooth water surface. SAR sensors also have an advantage over optical instruments as they can scan at night and are not impacted by cloud and weather, usually associated with a flooding situation. Due to improvements in spatial resolution and more frequent revisit times, SAR data has been used successfully to calibrate and validate hydrodynamic and hydraulic forecast models (Schumann et al., 2009; Grimaldi et al., 2016). Further model improvements have been shown through the assimilation of SAR data (e.g. García-Pintado et al., 2015; Hostache et al., 2018; Cooper et al., 2019; Di Mauro et al., 2020; Dasgupta et al., 2018, 2021a,b). Recent techniques have improved the flood detection in urban areas using medium and high resolution SAR (Mason et al., 2018, 2021a,b). The Copernicus Emergency Management Service (CEMS) (Copernicus Programme, 2021) offers freely available, open access Sentinel-1 SAR data. Currently (due to

the malfunction of Sentinel-1B in December, 2021) one satellite is in orbit, at 10 m ground resolution and a six day revisit time (for the mid-latitudes). Nevertheless, Sentinel-1 data offers good coverage of a potential flood event. For a major flood event CEMS can be triggered to offer additional rapid flood mapping. From 2022, the new Global Flood Monitoring (GFM) product (GFM, 2021; Hostache et al., 2021) of the Copernicus Emergency Management Service (CEMS) (Copernicus Programme, 2021) produces Sentinel-1 SAR-derived flood inundation maps using three flood detection algorithms providing uncertainty and population affected estimates within 8 hours of the image acquisition.

Representation errors arise where observation spatial scales are different from the model spatial scale (Janjić et al., 2018). The spatial resolution of SAR imagery suitable for flood detection varies across satellite constellation both historically and presently and continues to improve. Very high resolution (less than 3 m) imaging capabilities are increasingly available including TerraSAR-X, ALOS-2/PALSAR-2, and the COSMO-SkyMed, RADARSAT-2, and ICEYE constellations (Mason et al., 2021a). It is common practice to re-scale SAR-derived flood maps to match the model grid size for validation or assimilation with model data.

The objective of this paper is to present a scale-selective approach to evaluate flood inundation forecast maps and to develop a physically meaningful measure of flood-edge location accuracy that can be automated and easily applied in practice. The method has been developed with operational forecast verification in mind, but it is applicable to all flood inundation maps. A new approach is described and applied here to evaluate the spatial scale at which the forecast becomes useful/skilful at capturing the remotely observed flood extent and specifically the flood-edge location. The spatial skill of a forecast flood map varies with location. We aim to improve the conventional contingency map by incorporating the skilful scale to create a new *categorical scale map*. Also, we address how representation errors arising from observation spatial scale variations and interpolation have an impact on model evaluation.

In the rest of this paper we explore the features of a novel scale-selective evaluation approach illustrated through application to a case study. In Section 2 we describe the case study, a recent flooding event in the UK following Storm Dennis, February 2020, along with catchment descriptions for three chosen domains. The flood inundation forecasting system developed by JBA Consulting, Flood Foresight, (Revilla-Romero et al., 2017) is used to produce forecast flood maps for the event and is detailed in Section 3.1. Section 3.2 explains two methods that are used to derive remotely observed flood maps from SAR imagery. Our new approach to the spatial evaluation of flood maps is detailed in Section 4 along with descriptions of other binary performance measures. The novel categorical scale map is applied to the case studies in Section 5, and the evaluation results are discussed. We conclude in Section 6 and discuss the wider applications of a spatial scale approach to flood map skill evaluation.

2. Flood event

This extreme flooding event is chosen here as a case study to demonstrate the features of a spatial scale approach to forecast flood map evaluation. During February 2020, three named Storms, Ciara, Dennis and Jorge, arrived in quick succession delivered by a powerful and ideally positioned jet-stream that enabled rapid cyclogenesis (Davies et al., 2021). Each storm rapidly intensified and deepened bringing damaging winds and exceptionally heavy rainfall across the UK (Met Office, 2020). This led to the River Wye reaching its highest ever recorded water level at the Old Bridge in Hereford (riverlevels.uk, 2020). The annual exceedance probability (AEP) for the recorded peak flow of the Lugg and Wye rivers was 0.2–0.8% (return period 120–550 years) and 0.6–2.0% (160–550 years) respectively (Sefton et al., 2021).

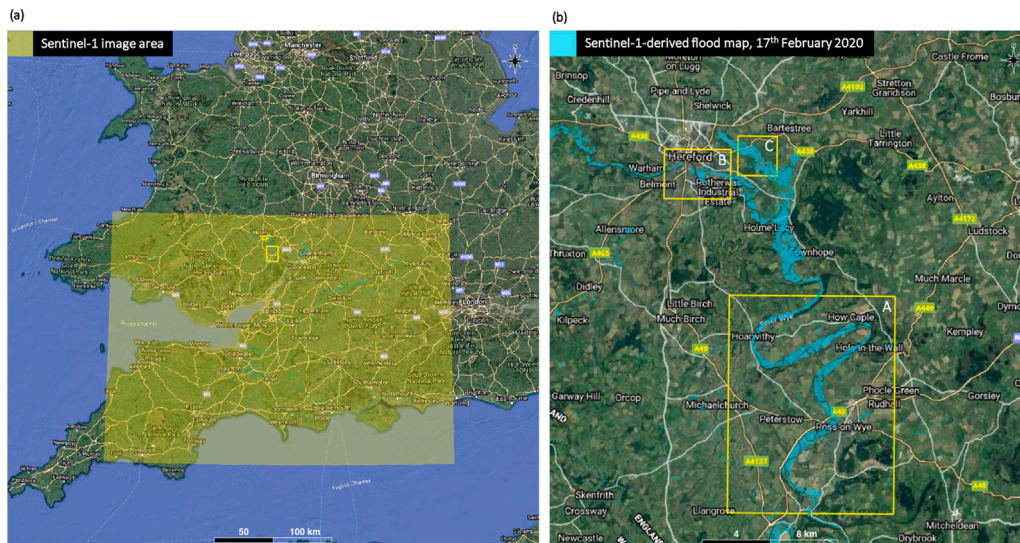


Fig. 1. Location of Sentinel-1 image acquisition over southeast UK (a) and flood map evaluation domains (b). Domain A: 28.4 km length of the River Wye centred at Ross-on-Wye, domain size 9.8×12.8 km. Domain B: 5.8 km of the River Wye at Hereford, domain size 3.0×4.0 km. Domain C: 4 km of the River Lugg at Lugwardine, domain size 2.3×2.3 km. Base map from Google Maps.

2.1. February 2020

February 2020 was the UK's wettest February on record and the fifth wettest month ever recorded. The UK average rainfall total exceeded the 1981–2010 average by 237% (Kendon, 2020). Locally, in northwest England and north Wales the rainfall exceedance was three to four times the typical monthly average rainfall. During this period around 4000 to 5000 properties were flooded in the UK, with significant river water levels recorded in Wales, west and northwest England (Sefton et al., 2021). With six days between Ciara and Dennis, groundwater and river levels were high and soils saturated. The Environment Agency issued a record number of over 600 flood alerts and warnings for England (JBA, 2021).

2.2. Catchment location and description

Three domains, each differing in hydrological characteristics, have been selected for forecast flood map evaluation during the storm Dennis flooding event. Two domains (A and B) have been chosen from the Wye catchment (Fig. 1), a 28.4 km length centred upon Ross-on-Wye (A) and the Wye at Hereford (B), a 5.8 km section. A third domain (C) includes 4 km of the River Lugg.

2.2.1. The River Wye (domains A and B)

The River Wye flows for approximately 215 km from Plynlimon at 750 metres above ordnance datum (mAOD) in the Cambrian Mountains, mid Wales. It initially travels southeastwards into England where it meanders southwards to ultimately join the Severn Estuary. The upper catchment land cover is predominantly grassland with some forest cover with highly impermeable bedrock and superficial deposits of sand and gravel in the Hereford area (National River Flow Archive, 2021). The upstream catchment area of Hereford is 1896 km². At Hereford, the only city situated on the Wye, the river is embanked on the north side by a deep flood wall with further embankments on the opposite side. Hereford is characterised by the Old Bridge, a 15th century stone bridge that creates a damming effect during high river flows. As the Wye flows south of Hereford, the topography flattens and the floodplain widens, with large river meanders and a distinctive U-shaped valley.

2.2.2. River Lugg at Lugwardine (domain C)

The River Lugg has an upstream catchment area of 886 km² and a maximum altitude of 660 mAOD and flows across the grasslands and agricultural fields of the Herefordshire plain. It has similar bedrock to the Wye catchment and a higher proportion of more permeable superficial fluvial deposits of sand and gravel. This is particularly evident in the Lugwardine region where the topography is relatively flat with little to impede the flow of floodwaters across the plain. The Lugg flows into the River Wye, 2 km south of domain C.

2.2.3. Event hydrology

The observed catchment rainfall (which also includes a downstream section of the River Wye) shows that 50 mm fell on the 15th, 10 mm on the 16th and 1 mm on the 17th February 2020 (UK Water Resources Portal, 2022). There were further heavy showers forecast for the 16/17th and whilst these have not been captured by the rain gauges on the 17th, they cannot be ruled out as contributing to surface water flooding in Hereford. The nearest hourly rainfall-rate observation is a citizen science observation from the Met Office WOW database (Met Office, 2022) for a site at Sutton St Nicholas near the River Lugg and this shows the highest rainfall rate of 5.8 mm/hr at 0300 on the 16th and a total accumulation of 12.5 mm on the 16th and 0.3 mm on the 17th.

Daily maximum river levels recorded at Ross-on-Wye, the Old Bridge, Hereford and Lugwardine for January to March 2020 are plotted in Fig. 2 (riverlevels.uk, 2020). The impact of the three storms on the River Wye is indicated by a very sharp rise in water levels from the 8th to the 10th February following storm Ciara. Further heavy showers maintained high water levels before storm Dennis brought an exceptional rise in water levels, peaking on the morning of the 17th February with record levels recorded at Hereford (6.11 m at 9.30 am UCT) and Ross-on-Wye (4.77 m at 5.45 am UTC). Unfortunately there are two days of missing data at Ross-on-Wye following the flood event. By analysing the trend between the Hereford and Ross-on-Wye river levels, the peak level at Ross-on-Wye was likely higher and later than recorded. The response of the Wye at Hereford is faster than at Ross-on-Wye, most likely due to the upstream location of Hereford and a more constrained embankment with the city center located either side of the river. In comparison to the fast, rapid response of the Wye, the River Lugg displays a distinctively dampened response. Whilst the Lugg initially responded quickly to the heavy rainfall, once bankfull was reached and overtopping occurred the water levels remained

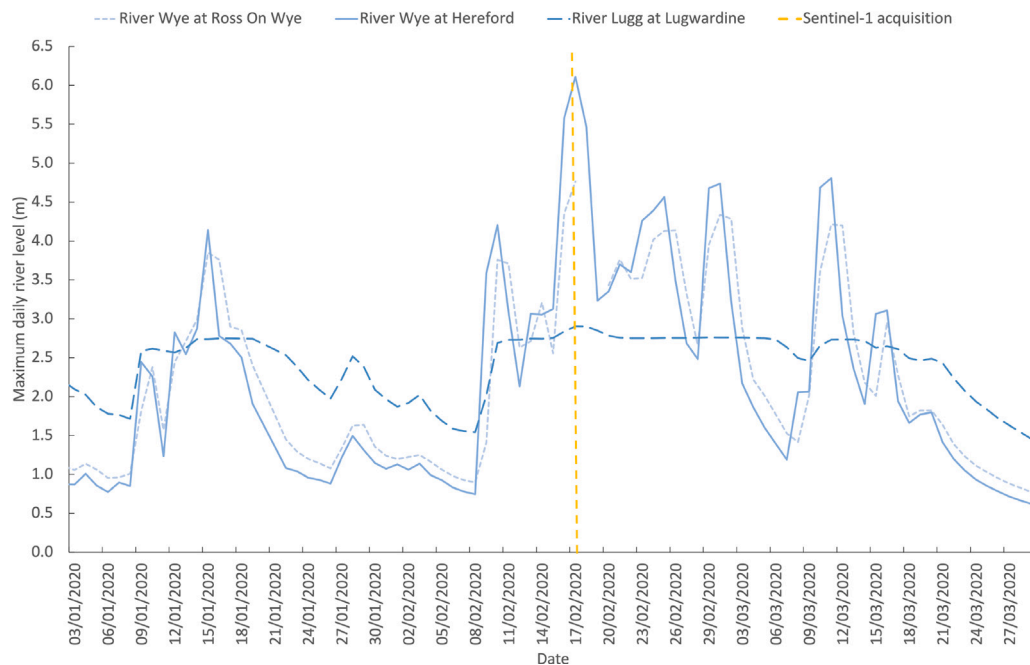


Fig. 2. Daily maximum river levels (m) at Ross-on-Wye, Hereford and Lugwardine. The dashed yellow line indicates Sentinel-1 SAR acquisition date. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

consistently high, with floodwaters extending across the relatively flat flood plain.

3. Data

In this section we describe the model and observation data that we will use to illustrate our novel scale selective verification approach.

3.1. Flood Foresight

Flood Foresight (Fig. 3), developed and run routinely by JBA Consulting, is a fluvial flood inundation mapping system that can be implemented in any catchment around the globe. Flood Foresight utilises a simulation library approach to generate maps of real time and forecast flood inundation and water depth. The simulation library approach saves valuable computing time and allows the application of Flood Foresight in near continuous real-time at national and international scales. A library of flood maps is pre-computed using JFlow®, a 2D hydrodynamic model (Bradbrook, 2006). Note that in this study the flood maps are undefended i.e. temporary flood defences are not included. JFlow uses a raster-based approach with a detailed underlying DTM and a simplified form of the full 2D hydrodynamic equations that capture the main controls of the flood routing for shallow, topographically driven flow. Five flood maps at 5 m resolution are created for 20, 75, 100, 200 and 1000 year return period flood events (corresponding to annual exceedance probabilities (AEPs) of 5%, 1.3%, 1%, 0.5% and 0.1% respectively). These are interpolated to derive five intermediate maps between each adjacent pair of the JFlow maps, equally spaced in return period creating a full library of thirty flood maps. Flood Foresight takes inputs of rainfall from numerical weather prediction (NWP) models, river gauge data (both historical and real-time) and forecast streamflow and uses these to select the most appropriate flood map for the location and forecast time period. The UK and Ireland configurations of the Flood Forecasting Module use deterministic streamflow forecast data from the Swedish Meteorological and Hydrological Institute (SMHI) European Hydrological Predictions for the Environment (E-HYPE). The meteorological input data for the E-HYPE model is the European Centre for Medium-range Weather Forecasts (ECMWF) Atmospheric Model high resolution (HRES) numerical weather prediction

(NWP) model on a $0.1^\circ \times 0.1^\circ$ grid with forecasts issued daily out to 10 days lead time. Forecast flood maps for the UK are produced on a 25 m grid length out to 10 days ahead (see Mason et al. (2021b) Section 2.1 for additional details).

3.2. SAR-derived flood maps

Two methods are applied to derive a flood map from SAR backscatter values captured close to the flood peak. The second method was included as it provides derivation of flood maps at different spatial resolutions. A Sentinel-1 (S1B) image was acquired in interferometric wide swath mode (swath width 250 km) just prior to the flood peak at 0622 on the 17th February. A pre-flood image (September 2019) from the same satellite sensor and track was used to derive the flood map in both methods.

In the first method, the ESA Grid Processing on Demand (GPOD) HASARD service (<http://gpod.eo.esa.int/>) has been utilised. The automated flood mapping algorithm (Chini et al., 2017) uses a statistical, hierarchical split-based approach to distinguish the two classes (flood and background) using a pre-flood and flood image. Level-1 GRD product SAR images (VV) are preprocessed, which involves; precise orbit correction, radiometric calibration, thermal noise removal, speckle reduction, terrain correction, and reprojection to the WGS84 coordinate system. The HASARD mapping algorithm removes permanent water bodies, including the river water. Flooded areas beneath vegetation, bridges and near to buildings are not detected using this method. The HASARD flood map at 20 m spatial scale is used to evaluate the performance of Flood Foresight for each of the three domains out to 10 days lead time.

In the second method, the same Sentinel-1 SAR image (in this case using both VV and VH) was processed using Google Earth Engine (GEE) to derive flood maps at a range of spatial resolutions (5 m to 25 m). GEE holds a catalogue of level-1 preprocessed Sentinel-1 SAR images (Google Earth Engine Catalog, 2021). A smoothing filter is applied to reduce speckle and a pre and post flood image are used to train a Classification And Regression Tree (CART) classifier (Breiman et al., 1984; Google Earth Engine CART, 2021). The classifier is applied to the whole image to produce a flood map at a specified scale. GEE uses

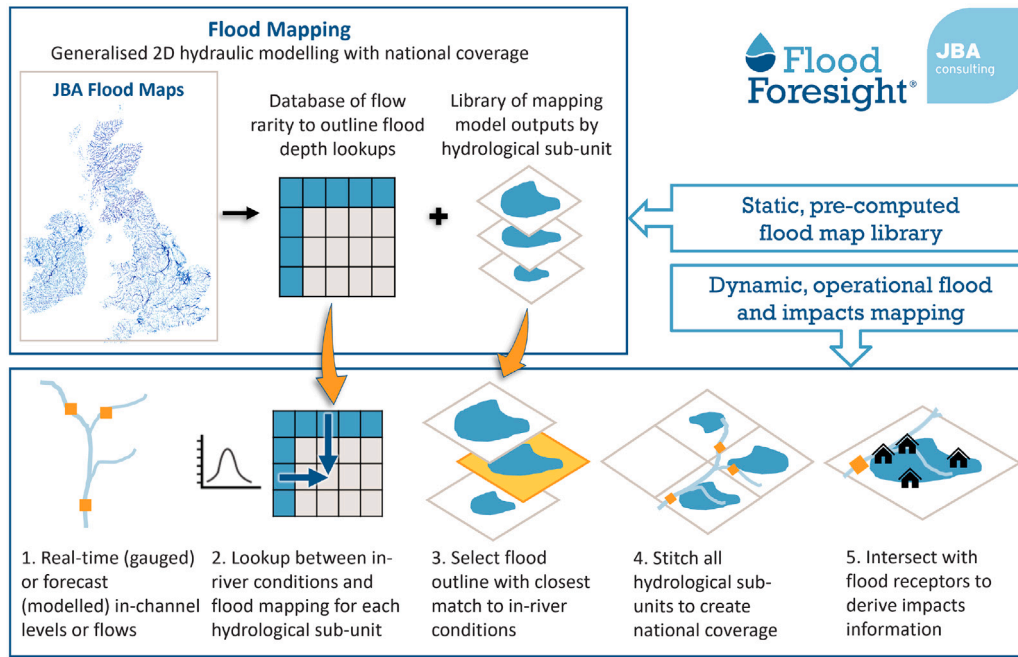


Fig. 3. Flood Foresight flood map simulation library selection process. Source JBA Consulting.

an image pyramid approach to scale, or pixel resolution, analysis. This means variations in the scale selected are determined from the scale of the input image (Google Earth Engine Scale, 2021). The variation of the flood extent detected at a range of spatial resolutions and the impact of re-scaling and interpolation errors on performance measures are investigated.

Flood Foresight forecast flood maps include the river channel and exclude surface features such as vegetation and buildings. To smooth the HASARD and GEE flood maps and allow a fairer comparison we apply a morphological closing operation (without impacting the location of the flood extent) to flood fill vegetation and buildings.

4. Flood map evaluation methods

The following subsections detail a new spatial scale-selective approach to forecast flood map evaluation. The Fraction Skill Score (FSS) developed by Roberts and Lean (2008) for validation of convective precipitation forecasts in atmospheric science uses a neighbourhood approach to determine the scale at which the forecast becomes skilful. Dey et al. (2016) developed this approach to determine an agreement scale between an ensemble forecast and observations at each grid cell to add location-specific information. Here we extend the technique to apply it to the new application of flood inundation mapping, and further develop a novel categorical scale map that combines an agreement scale map with a conventional contingency map.

4.1. Spatial scale-selective approach

Initially, the observed flood extent derived from SAR data is re-scaled to match the forecast flood map grid size using spline interpolation and both are converted into binary fields. A threshold approach is determined for the situation. For a flood map verification of spatial skill, the simplest example applied here is to assign each grid cell as flooded (1) or unflooded (0) for the whole domain. Alternative future threshold approaches for flood inundation maps could include applying thresholds to water depth percentiles. The location of the flood-edge cells can be extracted from the observed and modelled binary flood maps.

Given a domain of interest, we number all of the grid cells according to their spatial coordinates (i, j) , $i = 1 \dots N_x$ and $j = 1 \dots N_y$ where N_x is the number of columns in the domain and N_y is the number of rows. For each grid cell a square of length n forms an $n \times n$ neighbourhood surrounding the grid cell. The fraction of 1s in the square neighbourhood is calculated for each grid cell. This creates two fields of fractions over the domain for both the forecast M_{nij} and observed O_{nij} data. The fraction fields are compared against one another to calculate the mean squared error (MSE) for the neighbourhood

$$MSE_n = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [O_{nij} - M_{nij}]^2. \quad (1)$$

Based on the fractions calculated for the model and observed fields a worst possible MSE is calculated

$$MSE_{n(ref)} = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [O_{nij}^2 + M_{nij}^2]. \quad (2)$$

The FSS is given by

$$FSS_n = 1 - \frac{MSE_n}{MSE_{n(ref)}}. \quad (3)$$

Fig. 4 illustrates an example of the FSS application at grid level ($n = 1$) and at the next neighbourhood size $n = 3$. In this simple example, there is no agreement between the model and observation at grid level but at $n = 3$, the skill score improves to 0.92.

In general, the FSS is calculated for each length of neighbourhood n . For a given neighbourhood size an FSS of 1 is said to have perfect skill and 0 means no skill. The FSS will increase as n increases up to an asymptote (see Fig. 3 from Roberts and Lean (2008)). If there is no model bias across the whole domain of interest (observed and forecast flooded areas are the same) then the asymptotic fraction skill score (AFSS) at $n = 2N - 1$, where N is the number of grid cells along the longest side of the domain, will equal 1. Plotting FSS against spatial scale can indicate a range of scales where the model is deemed to be the most useful. This usefulness is a trade-off between being too smooth (larger n) or too fine, where the forecast skill is lost and the computation time lengthy. The gradient of the FSS curve versus neighbourhood size is another indicator of forecast skill with respect

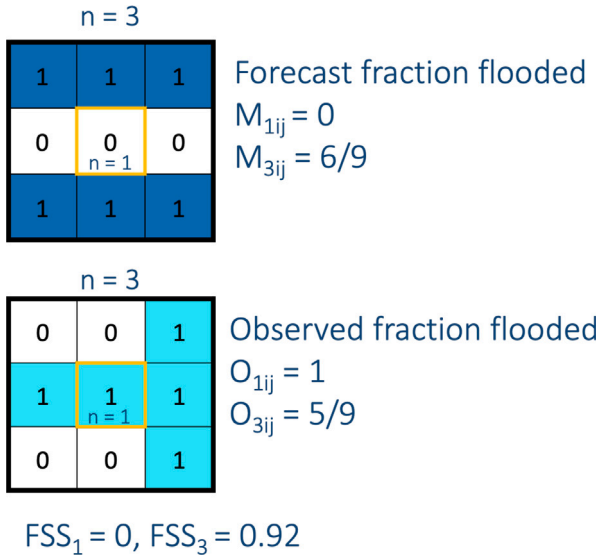


Fig. 4. FSS (see Section 4.1 for calculation details) example applied to a binary flooded (1) / unflooded (0) field at grid scale (yellow box, $n = 1$) and a 3×3 neighbourhood (black box, $n = 3$). The observed SAR-derived forecast is in turquoise and the forecast is shown in blue. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

to spatial scale. A steeper gradient indicates more rapidly improving skill over smaller grid sizes compared with a flatter curve, indicating a much wider neighbourhood is required to reach the same skill score. A target FSS score (FSS_T) is defined as

$$FSS_T \geq 0.5 + \frac{f_0}{2}, \quad (4)$$

where f_0 is the fraction of flood observed across the whole domain of interest and can be thought of as being equidistant between the skill of a random forecast and perfect skill. FSS_T will vary depending on the magnitude of the observed flood, relative to the domain area. This allows the comparison of the FSS_T scale across different domain sizes and floods of different magnitudes.

When the FSS is plotted against spatial scale (neighbourhood size), we can identify a spatial scale when the FSS first equals or exceeds FSS_T (Fig. 6 shows an example of this plot). The spatial scale (neighbourhood size) reached at FSS_T can tell us the displacement distance (D_T) between the observed and forecast flood, or more meaningfully the flood-edge locations. As the flood-edge represents a very small fraction of the domain, the scale at FSS_T will tend to $2D_T$, meaning the displacement distance is half of this scale (see Figure 4 in Roberts and Lean (2008)).

It has been shown by Skok and Roberts (2016) that care must be taken when calculating the FSS near to the domain boundary since increasingly larger neighbourhood sizes would extend further beyond the boundary edge. Skok and Roberts (2016) concluded that as long as the domain was sufficiently large, relative to the spatial errors, then the boundary effect could be considered to be insignificant. For flood mapping verification purposes the domain area should be selected to include the area of interest (e.g. the floodplain) with the neighbourhoods considered extending beyond the domain at the boundary. This assumes that the observations available allow this. If this is not that case then another boundary method could be applied, such as cropping at the domain edge.

4.2. Location dependent agreement scales

The FSS gives an overall domain-averaged measure of forecast performance and an average minimum scale at which the forecast is

Table 1

Contingency table (based on Stephens et al. (2014)).

	Forecast flooded	Forecast unflooded
Observed flooded	A (correct wet)	C (under-prediction/miss)
Observed unflooded	B (over-prediction/false alarm)	D (correct dry)

deemed skilful. Dey et al. (2016) describe a method for calculating an agreement scale at each grid cell located at coordinate position (i, j). A brief summary of the method is presented here. Two fields are considered f_{1ij} and f_{2ij} . In this application these are the forecast and observed fields. In alternative applications the method could be applied to measure similarity between members of an ensemble. The fields in this instance are not required to be thresholded and can be applied to flood depths. The aim is to find a minimum neighbourhood size (or scale) for every grid point such that there is an agreement between f_{1ij} and f_{2ij} . This is known as the agreement scale S_{ij} . The relationship between the agreement scale and the neighbourhood size described in Section 4.1 is given by $S_{ij} = (n - 1)/2$.

Firstly, all grid points are compared by calculating the relative MSE D_{ij}^S at the grid scale, $S = 0$ ($n = 1$),

$$D_{ij}^S = \frac{(f_{1ij}^S - f_{2ij}^S)^2}{(f_{1ij}^S)^2 + (f_{2ij}^S)^2}. \quad (5)$$

If $f_{1ij} = 0$ and $f_{2ij} = 0$ (both dry) then $D_{ij}^S = 0$ (correct at grid level). Note that D_{ij}^S varies from zero to 1. The fields are considered to be in agreement at the scale being tested if:

$$D_{ij}^S \leq D_{crit,ij}^S \quad \text{where} \quad D_{crit,ij}^S = \alpha + (1 - \alpha) \frac{S}{S_{lim}} \quad (6)$$

and S_{lim} is a predetermined, fixed maximum scale. The parameter value α is chosen to indicate the acceptable bias at grid level such that $0 \leq \alpha \leq 1$. Here we set $\alpha = 0$ (no background bias). If $D_{ij}^S \geq D_{crit,ij}^S$ then the next neighbourhood size up is considered ($S = 1$, a 3×3 square). The process continues with increasingly larger neighbourhoods until the agreement scale, or S_{lim} is reached for every cell in the domain of interest. The agreement scale at each grid cell is then mapped onto the domain of interest.

4.3. Categorical scale map

Currently, the agreement scale map proposed by Dey et al. (2016) provides a location-specific scale of agreement between the forecast and observed flood map. However, it does not show whether the model is over- or under-predicting the flood extent. In our work, we develop the agreement scale map further by combining with a contingency map for the forecast to create a new *categorical scale map*. This highlights the agreement scale for areas of over- or under-prediction. In a contingency map, each cell in the forecast and observed flood map are compared and classified using a contingency table (Table 1). The categories are re-classified numerically in the array for automated updating of the agreement scale map. Over-predicted cells (B) are set to -1 , under-predicted cells (C) are set to $+1$, correctly predicted flooded cells (A) are assigned NaN and correctly predicted unflooded cells are set to 0. The array element-wise product of the agreement scale map and the numerical contingency map produces the new categorical scale map.

4.4. Binary performance measures

It has been suggested by Cloke and Pappenberger (2008) that a range of performance measures should be applied so that a forecast can be assessed as rigorously as possible. A selection of commonly applied binary performance measures, each focusing on a different aspect of performance have been included here for comparison with the Fraction Skill Score results. Following the application of a contingency table

Table 2
Binary performance measures and formula based on contingency Table 1.

Performance measure	Formula	Description [range min, range max, perfect score]
Bias	$\frac{A+B}{A+C}$	[0, ∞ , 1] 1 implies forecast and observed flooded areas are equal > 1 indicates over-prediction, < 1 indicates under-prediction
Critical Success Index/Threat score $F_{<2>}$ (CSI)	$\frac{A}{A+B+C}$	[0, 1, 1] Fraction correct of observed and forecast flooded cells
$F_{<1>}$ Proportion correct	$\frac{A+D}{A+B+C+D}$	[0, 1, 1] Proportion correct (wet and dry) of total domain area
$F_{<3>}$	$\frac{A-C}{A+B+C}$	[-1, 1, 1] Score reduced by over-prediction
$F_{<4>}$	$\frac{A-B}{A+B+C}$	[-1, 1, 1] Score reduced by under-prediction
False Alarm Rate (FAR)	$\frac{B}{B+D}$	[0, 1, 0] Proportion of over-prediction of dry areas
Hit Rate (HR)	$\frac{A}{A+C}$	[0, 1, 1] Fraction correct of observed flooded area
Pierce Skill Score (PSS)	$HR - FAR$	[-1, 1, 1] Incorporates both under and over-prediction

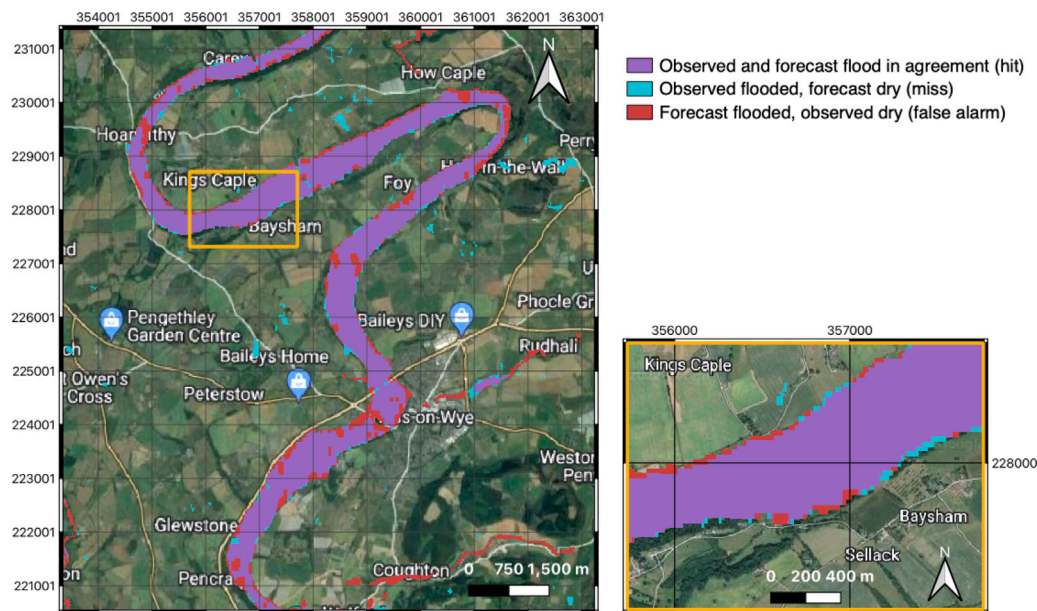


Fig. 5. Left panel: contingency map of a 0-day lead time forecast versus the HASARD SAR-derived flood map for the Wye valley indicates the model is predicting the flood extent accurately, including the position of the flood-edge. Right panel: Zoom of yellow box on the left panel. On closer inspection, at grid level, the flood-edge in many places is over- or under-predicted by around one grid length. Base map from Google Maps. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

(Table 1) to the forecast flood map, a number of binary performance measures can be calculated (Table 2). Table 2 describes the range of performance value, the ideal score and a description of which aspects of the forecast flood map performance each binary measure assesses.

5. Results

We illustrate and discuss our new method applied to the flood event in Sections 5.1 and 5.2. The scale-selective approach is applied to an extreme flooding event in the UK to determine a useful/skilful spatial scale for both the entire flood extent and the flood-edge location for three domains out to 10-days lead time. An example forecast flood map for 0-day lead time compared with the SAR-derived flood map is presented as a contingency map in Fig. 5. The zoomed in perspective shows the double penalty impact described in Section 1. The discrepancy at the flood-edge depends on the spatial scale of the forecast flood maps along with the model performance. Next, in Section 5.3 location-specific agreement scales are presented on categorical scale maps. The final Section 5.4 addresses the question of the impact of representation error caused by variations in SAR-derived flood map spatial resolution on the evaluation results.

5.1. Spatial scale variability of forecast flood extent and flood-edge location

An evaluation of the spatial skill of the Flood Foresight forecast flood maps against the SAR-derived flood map for the flood peak on the 17th February 2020 has been calculated for each domain (Fig. 1) for both the entire flood extent and the flood-edge location. The Fraction Skill Score (FSS) is applied to increasing neighbourhood sizes (n) to determine the spatial scale at which the forecast becomes skilful at capturing the observed flood. Fig. 6 shows FSS against n for one example, the River Lugg (domain C) for the entire flood (a) and the flood-edge (b). Each line represents a different model run date from the 10/02/2020 (7-day lead time) to the 17/02/2020 (0-day lead time). With the exception of the 7-day lead time, all forecasts for the whole flood (Fig. 6a) exceed the FSS_T at grid level ($n = 1$) with gradually improving skill as n increases. In contrast to this, the FSS applied to the flood-edge (Fig. 6b) shows all forecasts below FSS_T at grid level and $n = 3$ with the skill increasing more rapidly compared with the whole flood to reach FSS_T at $n = 5$ for all run dates within a 5-day lead time (except for 16/02/2020, which is just below FSS_T). This indicates that the flood-edge is forecast to be around 62.5 m

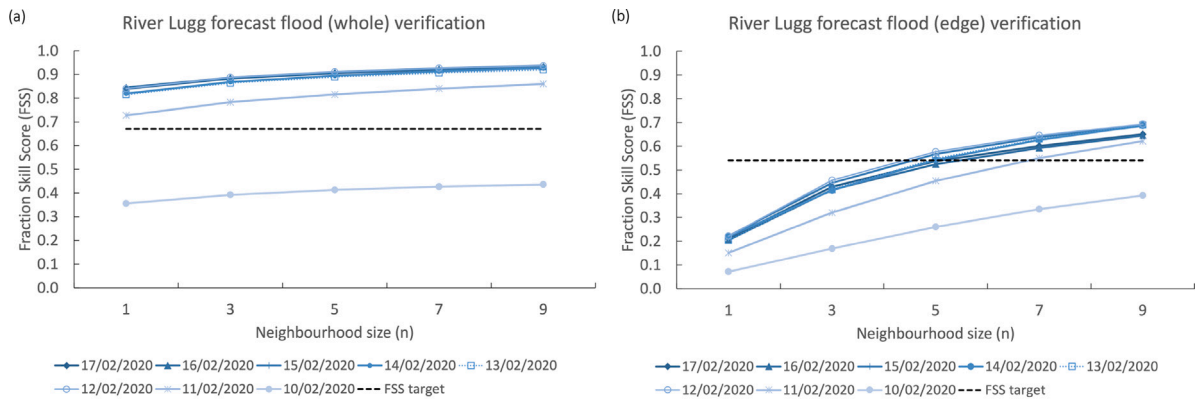


Fig. 6. FSS calculated for the River Lugg at Lugwardine for (a) entire flood extent and (b) the flood-edge for increasing neighbourhood sizes for daily forecast lead times up to 7 days.

from the observed flood-edge, on average, for a 5-day lead time. The difference between the gradients of the plots indicate the flood-edge is more sensitive to changes in spatial scale compared with evaluation of the whole flooded area. The whole flood verification here indicates a strong model performance. However, verifying the whole flood alone could mask the flood-edge location performance, which in this case has a coarser scale at FSS_T . Similar trends in FSS with neighbourhood size and comparisons between the entire flood and the flood-edge verification scales are found for all domains. The rate of FSS increase, or FSS gradient with n , tells us how quickly the forecast skill improves with increasing scale. A more spatially accurate forecast of the flood-edge will demonstrate a steeper gradient, reaching FSS_T at a smaller neighbourhood size.

5.2. Comparison of spatial scales at differing lead times and domain location

The performance measures for each domain for daily lead times out to 10 days are presented in Fig. 7. The FSS at $n = 1, 3$, and 5 are shown along with Critical Success Index (CSI), Hit Rate (HR), Pierce Skill Score (PSS) and the Bias (see Table 2 for definitions). The Bias score is an indicator of over- or under-prediction of the flood extent and is plotted on a separate axis to account for the larger range. For lead times within 5-days of the flood peak, $FSS > 0.8$ for the entire flooded area at grid level for the River Wye (domain A) indicates a strong model performance (Fig. 7a). There is a dip in the FSS on the 16/02/2020 where the forecast over-predicts the flood extent. This is also reflected in the CSI score. In contrast to this the HR and PSS increase, despite the over-prediction, as more observed flood cells are correctly predicted wet. We note that the PSS (HR - FAR) does account for over-prediction, however the FAR is the fraction of the dry area incorrectly predicted wet, which is very small relative to the HR (0.03 versus 0.90). Validation of the River Wye flood-edge (Fig. 7b) is more sensitive to changes in neighbourhood size compared with the whole flood validation. Here the flood-edge is very well forecast in terms of spatial location and exceeds FSS_T at $n = 3$ (on average, 37.5 m displacement) for a 5-day lead time (except for 1-day lead time where FSS_T is exceeded at $n = 5$). As shown previously in Section 5.1, the forecast of the River Lugg flood-edge is skilful at $n = 5$ (Fig. 7f) (on average, 62.5 m displacement) for a 5-day lead time. Differences in the hydrological characteristics might explain differences in model performance. The Wye valley flood plain is well defined with distinctive valley sides and this event proved to be valley filling in contrast to the Lugg flood plain which is relatively flat and extensive. This could explain the increased skill shown for the prediction of the Wye flood-edge. The average observed flood top width for the Lugg (domain C) is 740 m and for the Wye (domain A) 430 m. This gives a flood-edge

displacement as a fraction of the flood top width of 7.4% for the Lugg and 7.8% for the Wye.

The results for all three domains show that for this case study the forecasting system has limited skill beyond a five-day lead time. The forecast accuracy of the meteorological driving data diminishes with increasing lead time (ECMWF, 2022). Extratropical cyclones (ETCs) are the dominant meteorological driver of major winter flooding in the UK. This is particularly true when an Atmospheric River is associated with an ETC and when ETCs arrive in clusters (as was the case here) bringing multiple spells of heavy precipitation (Lavers et al., 2011; Griffith et al., 2020). The typical formation time of ETCs is 3–5 days, occasionally up to 10 days (Ulbrich et al., 2009) which limits the predictability of the meteorological system, particularly when the jet stream is very strong (as was the case here). The atmospheric (and precipitation) predictability will vary depending on the situation, for example a slow moving ETC close to the UK would potentially have a longer lead time of useful prediction. Conversely, flooding in the summer associated with convection would likely have a shorter skilful lead time. The scale selective approach presented here can be used to determine a meaningful scale to present flood inundation maps. This scale will vary with forecast lead time and will depend on the predictability of the meteorological situation.

There is more variation in skilful scale with lead time evident for the Wye at Hereford (domain B) in Fig. 7c and d compared with domain A and C. To achieve the same FSS for the whole flood as domain A and C up to a 5-day lead time, the neighbourhood size would need to exceed $n = 5$. The model is over-predicting the flood extent, in particular on the 16/02/2020 (1-day) lead time. This overprediction at 1-day lead time is evident for all domains as can be seen in the Bias scores but the impact of this is most noticeable at Hereford. Hereford has more complex topography compared to the other domains, particularly along the river bank with bridges, buildings, permanent and temporary flood defences deployed during the event affecting the flow of the flood wave through the city. The maps used in the simulation library of Flood Foresight are produced using a bare-earth DTM. Despite this, the model performs well, exceeding FSS_T at $n = 5$ at the 5-day and 2-day lead times for the flood-edge forecast.

Overall, the FSS indicates a similar trend in performance across all results as the commonly applied CSI. The value of FSS_T is determined by the magnitude of the observed flood, which means the skilful scale determined at FSS_T can be meaningfully compared across the domains. The skilful scale of the forecast flood-edge location gives an average discrepancy distance. A physically meaningful evaluation measure provides additional information compared to a conventional verification score.

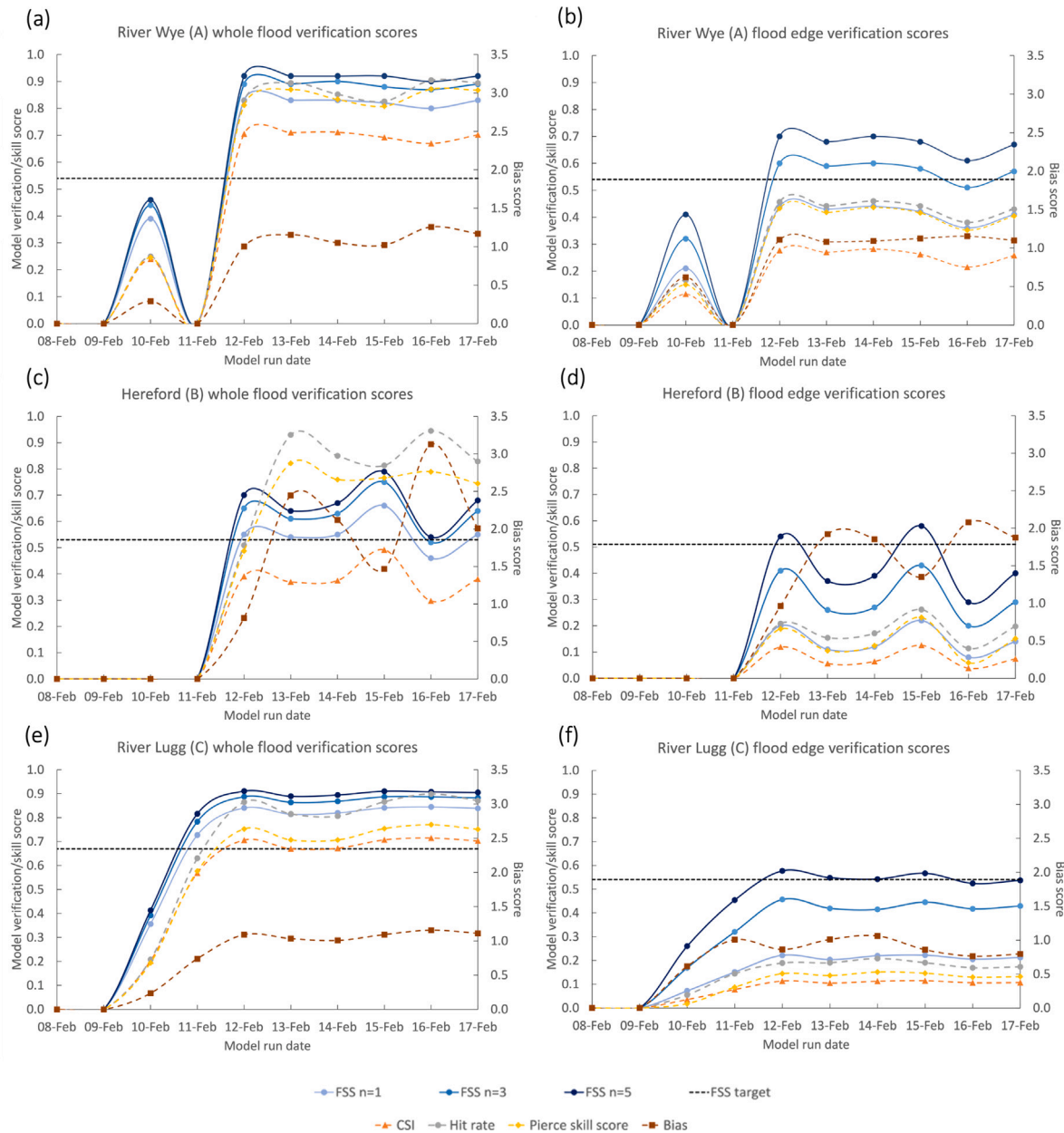


Fig. 7. Conventional binary performance measures (dashed lines) and FSS (solid lines) at $n = 1, 3$, and 5 for each domain for both the whole flooded area and the flood-edge for daily lead times out to 10 days for the River Wye (domain A, (a) and (b)), Hereford (domain B, (c) and (d)) and the River Lugg (domain C, (e) and (f)). Plots on the left show the verification scores applied to the entire flood extent and plots on the right show the flood-edge scores.

5.3. Categorical scale maps

Location dependent categorical scale maps (Section 4.3) have been calculated for all run dates for both the entire flooded area and the flood-edge. Fig. 8 shows categorical scale maps for the whole flood for three different lead times for each domain, longer lead times are on the left. The run dates vary with domain to present the most informative maps such that variation in forecast skill can be seen across the different lead times. The colours on the map indicate grid cell specific agreement scales (Section 4.2) between the forecast flood map and the SAR-derived flood map. Grey/white regions indicate correctly predicted flooded/unflooded cells, red shows the forecast flood extent is under-predicted (miss) and blue indicates over-prediction (false alarm). Increasingly darker shades of red/blue show that larger scales were needed for the agreement criteria to be met. The darkest blue at $S = 10$ indicates a total mismatch between forecast and observed flooding. The addition of the agreement scale information in comparison to a

conventional contingency map (for an example, see Fig. 5) quickly highlights regions of total mismatch through the darkest shading, with areas that are slightly misaligned in lighter shades. The agreement scale indicated gives a physical measure of distance at specific locations between the forecast and the observed flood map (where $S < S_{lim}$).

The location-specific skilful scale varies with location and lead time as indicated on the categorical scale maps. For a 7-day lead time forecast for the River Wye (Fig. 8a), the model is indicating some flooding could occur, although under-estimating the total extent as shown by the darkest red areas, which show the limits of the agreement scale have been reached. By 5-days lead time the forecast is in very close agreement with the observed flood at grid level (in grey) with larger agreement scales indicated by red/blue shading along some of the flood-edge locations (Fig. 8b) and a balance between under- and over-prediction. Over-prediction is more evident by 1-day lead time for the River Wye (Fig. 8c) and flooding is also over-predicted along smaller tributaries. There are several detached areas of flooding

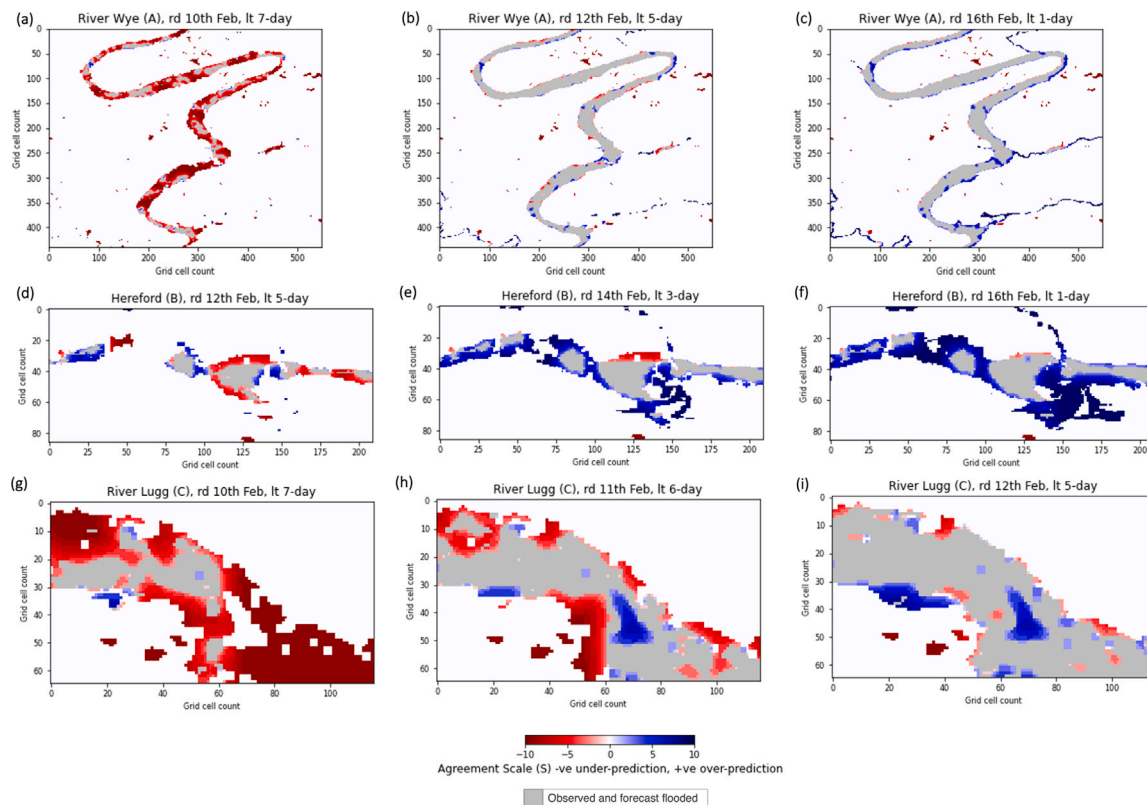


Fig. 8. Categorical scale maps for each domain at various lead times (Lt). Red indicates where the forecast flood extent is under-predicted, blue indicates over-prediction. The shading indicates the agreement scale, a measure of distance between the forecast and observed flood maps. Grey areas are correctly predicted flooded, white areas are correctly predicted unflooded. Each grid cell represents 25 m \times 25 m for all domains. (Note: rd (forecast run date) varies between location, all dates have been evaluated and the most illustrative maps selected.). (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

observed remotely that are most likely due to ponding of surface water flooding, which were not predicted by the fluvial flood forecasting system.

The Hereford forecast is most skilful on the 12th February (Fig. 8d) with over-prediction, particularly towards the southwest at 3-day and 1-day lead times (Fig. 8e and f). A small stream running southwards to the Wye, the Eign Brook, could be contributing to the over-prediction seen here. It is also worth mentioning that SAR will struggle to detect flood waters where buildings are closer together when the distance between them is less than the ground resolution of the SAR. Shadow and layover effects due to the side-looking nature of the SAR also mean flood detection is more difficult in urban areas (Mason et al., 2021a). This will likely only impact a small area of the Hereford domain but this observation uncertainty should be considered when interpreting these results. There is an area of under-prediction of the flood extent in the centre of the Hereford domain visible at all lead times. This could be due to surface water flooding, which most likely occurred due to the very high intensity rainfall observed. This combined with the urban area and steeply sloping gradient to the north of this area most likely contributed to rapid surface water runoff towards the river. Since Flood Foresight is a fluvial flood forecast system we would not expect surface water flooding such as this to be predicted.

Flood Foresight selects multiple flood maps and stitches them together when the return period threshold is exceeded for a given area. The Hereford section of the Wye does not trigger a flood map selection until a 5-day lead time, this area also influences part of the River Lugg flood map and can be seen as a mismatch on the lower left hand side of Fig. 8g and h. Once this is included the forecast flood map is in very good agreement from a 5-day lead time. There are areas that could be further improved, indicated by the lighter shading (Fig. 8i). An acceptable level of agreement scale could be determined

for a given situation, for example $n < 5$, and efforts made to understand/improve larger agreement scales at specific locations. These improvements might include changes to infrastructure included in the DTM used in the hydraulic modelling, for example.

5.4. SAR-derived flood map scale variation

In practice, particularly where a flood event is prolonged or the flooding extent covers a wide area, there may be multiple sources of SAR data available for model evaluation, usually at higher spatial resolutions compared to the model grid size (e.g. ICEYE in spot mode at 1 m and strip mode at 3 m ground resolution). It is important to consider the impact of using observations at different spatial scales on the scale-selective approach results. By conducting a simulation experiment we address the question of how re-scaling and interpolating three higher spatial resolution SAR-derived flood maps (relative to the forecast flood maps) affects the scale selective skill scores and location-specific forecast skill. In order to simulate a range of observation spatial scales, SAR-derived flood maps are produced using method two described in Section 3.2 at spatial resolutions from 5 m to 25 m. These are re-scaled by 0-order spline interpolation (ndimage.zoom, 2021; Briand and Monasse, 2018) to match the model resolution (25 m) and compared to the forecast flood map for the River Lugg (5-day lead time). A comparison of the GEE flood map against the HASARD flood map, both at 20 m spatial scale produce almost identical verification scores for all performance measures for the River Lugg ($\Delta F.S.S. < 0.01$).

The categorical scale maps for the comparison between the forecast flood map and the re-scaled simulated SAR-derived flood maps are shown in Fig. 9. The resulting domain-averaged skill scores for the same forecast flood map against the four SAR-derived flood maps are displayed in Fig. 10. The scores are calculated for the whole flood and the flood edge alone. In general, the categorical scale maps show similar

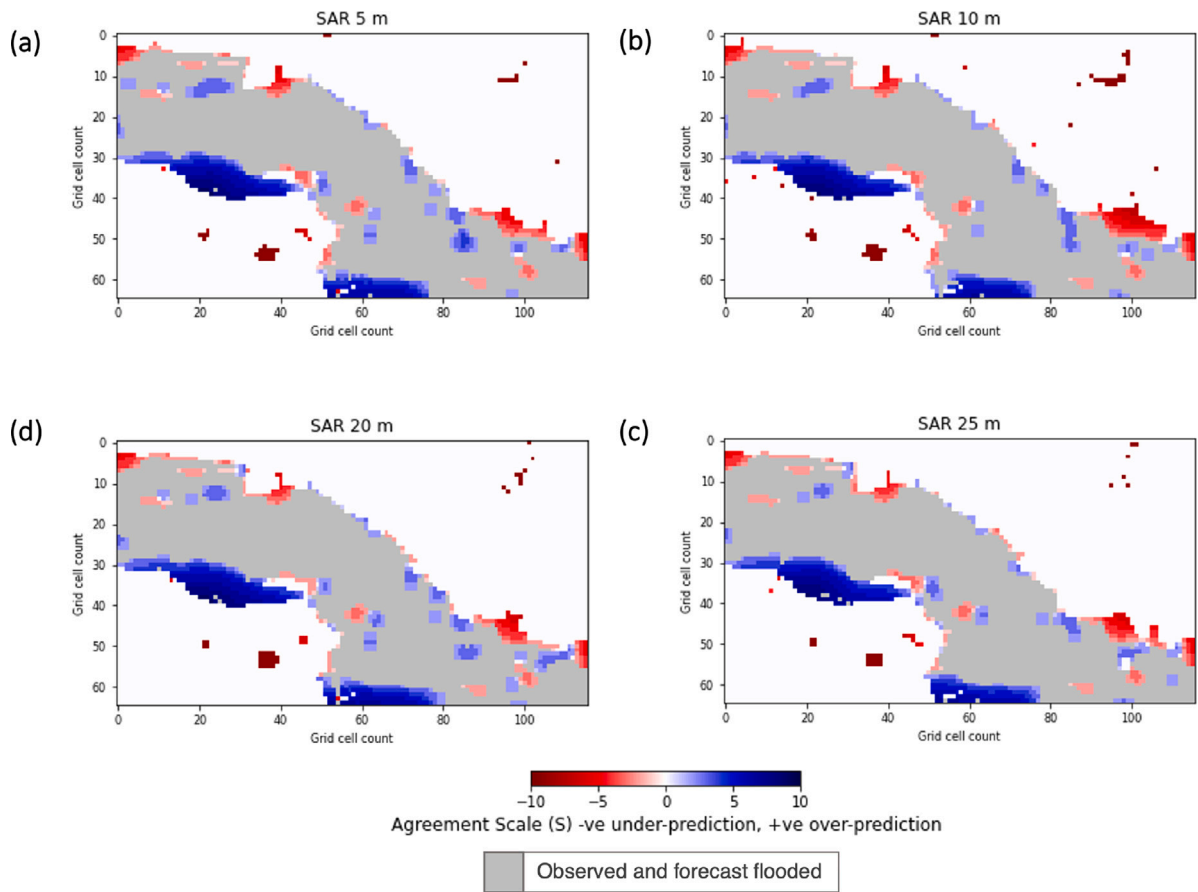


Fig. 9. SAR-derived flood maps produced at different spatial resolutions (5 m to 25 m) are re-scaled to the model grid size (25 m) before categorical scale maps are calculated for the River Lugg (C), run date 12th Feb.

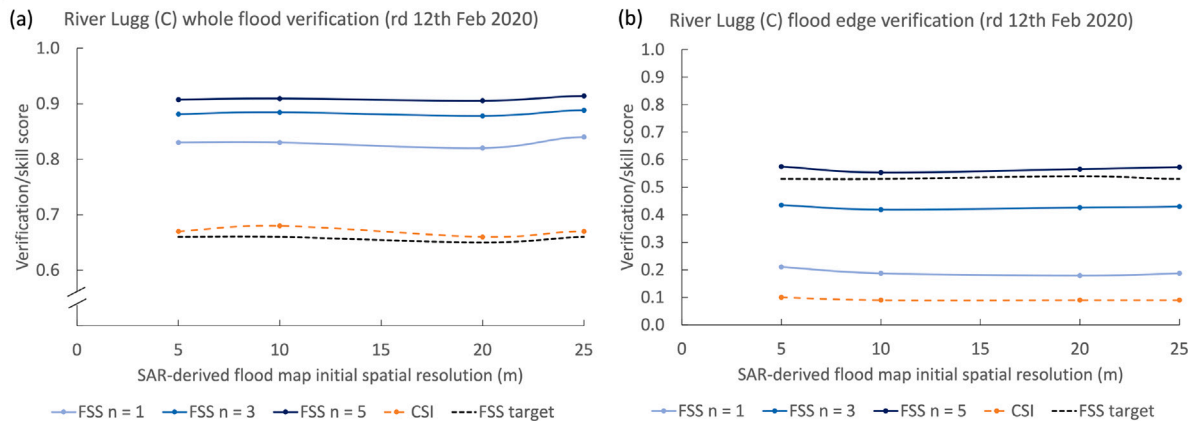


Fig. 10. SAR-derived flood maps at different spatial resolutions (5 m to 25 m) are re-scaled to the model grid size (25 m) before verification scores are calculated for the whole flood (a) and the flood-edge (b). Note that axes in (a) and (b) are on different scales.

regions of over and under-prediction but there are small location-specific variations in skillful scale. The SAR-derived flood map at 25 m, the same spatial scale as the forecast flood maps, shows the best agreement away from the flood edge. This is also evident in the overall FSS score for the 25 m comparison, which marginally outperforms the evaluation after re-scaling finer observation flood maps (Fig. 10). The skillful scale determined for each observation comparison of the whole flooding extent is $n = 1$ or at grid level, and for the flood edge is at $n = 5$. Overall, based on the results from this simulation experiment, the scale-selective approach is not overly sensitive to the observation spatial scale and the skillful scale determined remains the same for each

of the observed SAR-derived flood maps for both the entire flood extent and the flood edge. Small errors are introduced by re-scaling and interpolating finer resolution observations to the model spatial scale which slightly reduce the skill score and change location-specific details on the categorical scale maps. Observation scale selection and re-scaling along with interpolation errors must be considered when evaluating model performance, particularly where model or observation scales vary in space and time, or where comparisons are made across different models.

6. Discussion and conclusions

Overall, the aim of this paper was to introduce and apply a new scale-selective approach to forecast flood map evaluation with an emphasis on providing a physically meaningful verification of the flood-edge location. The skilful spatial scale for comparison of forecast flood inundation maps against SAR-derived observed flood extent has been evaluated by the application of the Fraction Skill Score: this provides a domain-averaged skilful scale. The verification measure has been applied to a forecast of an extreme flood event in the UK on the River Wye and the River Lugg following Storm Dennis in February 2020. Flood Foresight inundation predictions with lead times out to 10 days are evaluated against a Sentinel-1 SAR-derived flood map captured close to the flood peak for three domains, each differing in hydrological characteristics. Conventional binary performance measures were calculated alongside the FSS for comparison. Flood-edge verification shows greater sensitivity to changes in forecast skill and spatial scale, relative to verification of the entire flood extent. The skilful scale determined is physically meaningful and can be used to estimate the average flood-edge discrepancy from the observed flood-edge. The observed flood map spatial resolution relative to the model scale is important and re-scaling and interpolation errors will impact the model verification scores. Ideally, the observed flood map should be derived at the same spatial scale as the forecast model to minimise these errors.

In operational practice the scale at which the forecast flood maps are presented to forecasters and decision makers should reflect the uncertainty within the forecast. Very high resolution flood maps can be presented where a detailed DTM is available. If this is presented as a deterministic forecast to flood risk management teams, it could lead to an over confidence in the forecast, or where the actual observed flood magnitude is different, the forecast may be devalued in the future (Savage et al., 2016; Speight et al., 2021). Application of a spatial-scale approach to forecast evaluation can determine the scale at which it is best to present the forecast flood map. Conversely, if the model is found to be skilful at grid level, there is scope to increase the flood map resolution adding more detail to the flood-edge location. Improvements made to hydrodynamic models, such as through data assimilation to improve inputs, initial conditions or model parameters may not improve the forecast flood-edge location at grid level. However, improvements may be evident through evaluation using FSS across a range of scales. Categorical scale maps are a useful evaluation and forecasting tool, adding location-specific detail. Model improvements can be spatially targeted and as improvements are made, the categorical scale map will highlight location-specific changes. For example, the categorical scale maps for Hereford indicate the local infrastructure (in particular bridges) impact the movement of the flood wave, which suggests a digital surface model (DSM) would be beneficial in urban areas.

The verification approach is presented here in the context of an operational flood forecasting system. The skilful scale determined for each flooding scenario, lead time and at specific locations within a domain depends on the skill of the entire hydrometeorological chain of forecasting models from the meteorological inputs to the hydrodynamic model (run offline in the case study presented here) used to determine the inundation extent for a given river discharge. The scale-selective approach is equally applicable for the validation of flood maps from hydrodynamic models that are not part of an operational system. Here, we focus on the use of SAR-derived flood maps for validation, however the approach would apply to any remotely observed flooding such as from optical satellite data or UAV aerial imagery that can be converted into a gridded dataset. The FSS must be applied to binary data and for this reason it is very easily applicable to flooding extent with grid cells categorised as flooded/unflooded. In operational forecasting, flood depth is also an important metric to verify and by applying a threshold (depths below/above a certain level or percentile), the depth data can be converted for application of FSS. The method for calculating

categorical scale maps does not require binary data and so the depth values can be used directly in the calculations.

Ideally, in operational forecast systems, quantitative validation should run in tandem with the forecast system where observations are available. Over time, a catalogue of skilful scales, flood edge discrepancy distances and categorical scale maps could be built up. This catalogue would enable analysis of scale across different flood event type, season, meteorological scenario, forecast lead time and at specific locations within a catchment or sub-catchment. Such a verification library would enable forecasters to increase intuition and expert judgement on the relevant scales for a given forecast. Based on this analysis and an increased understanding of the predictability of flood inundation, forecast flood maps could be presented at a variable scale. For example, a coarser scale at longer lead times becoming more detailed, closer to the flooding event. Coarse scales can appear jagged or with large steps along the edge and so ideally these would be converted to smooth contours, but with some indication (for example, lighter shading) that the flood edge lies somewhere within the width of the grid cell, rather than exactly at the contour edge. At shorter lead times, as forecast confidence is assumed to increase, the flood edge location would show more detail and a narrower band of uncertainty (grid cell width). This flood edge uncertainty information will prove invaluable for impact-based forecasting practice.

The spatial-scale approach will also prove a useful tool in multi-model performance comparisons where forecast flood maps are presented at different spatial resolutions or to evaluate the performance of an increase in model resolution. Evaluating a skilful scale for each model can be compared directly whereas the skill score values should not be compared across models with different spatial scales (Emerton et al., 2016). These methods will also benefit surface water flooding verification where the flood map is likely to be localised and discrete and accounting for variations in spatial skill more critical. An improved approach to evaluating forecast flood maps will result in improved accuracy in the predictions of flooding. Ultimately, this will benefit disaster management teams and those living in flood prone areas to enable future mitigation of flooding impacts.

CRedit authorship contribution statement

Helen Hooker: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing – original draft, Visualization. **Sarah L. Dance:** Conceptualization, Methodology, Writing – review & editing, Supervision, Project administration, Funding acquisition. **David C. Mason:** Methodology, Writing – review & editing, Supervision. **John Bevington:** Methodology, Resources, Writing – review & editing, Supervision, Project administration. **Kay Shelton:** Methodology, Resources, Writing – review & editing, Supervision.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Code and data availability

The functions used to evaluate the forecast flood maps using a scale-selective approach along with the SAR-derived and forecast flood maps are available on the following Zenodo page: <https://doi.org/10.5281/zenodo.6011882> (Hooker, 2022).

Acknowledgements

Funding

This work was supported in part by the Natural Environment Research Council, UK as part of a SCENARIO funded PhD project with a CASE award from the JBA Trust (NE/S007261/1). Sarah L Dance and David C Mason were funded in part by the UK EPSRC DARE project (EP/P002331/1). Sarah L Dance also received funding from NERC National Centre for Earth Observation, UK.

References

- Bradbrook, K., 2006. JFLOW: A Multiscale two-dimensional dynamic flood model. *Water Environ. J.* <http://dx.doi.org/10.1111/j.1747-6593.2005.00011.x>.
- Breiman, L., Friedman, J., Stone, C.J., Olshen, R., 1984. *Classification and Regression Trees*. Chapman and Hall/CRC.
- Briand, T., Monasse, P., 2018. Theory and practice of image B-spline interpolation. *Image Process. On Line* 8, 99–141. <http://dx.doi.org/10.5201/ipol.2018.221>.
- Briggs, W.M., Levine, R.A., 1997. Wavelets and field forecast verification. *Mon. Weather Rev.* 125 (6), 1329–1373. [http://dx.doi.org/10.1175/1520-0493\(1997\)125<1329:waffv>2.0.co;2](http://dx.doi.org/10.1175/1520-0493(1997)125<1329:waffv>2.0.co;2).
- Casati, B., Wilson, L.J., 2007. A new spatial-scale decomposition of the brier score: Application to the verification of lightning probability forecasts. *Mon. Weather Rev.* 135 (9), 3052–3069. <http://dx.doi.org/10.1175/MWR3442.1>.
- Chini, M., Hostache, R., Giustarini, L., Matgen, P., 2017. A hierarchical split-based approach for parametric thresholding of SAR images: Flood inundation as a test case. *IEEE Trans. Geosci. Remote Sens.* 55 (12), 6975–6988. <http://dx.doi.org/10.1109/TGRS.2017.2737664>.
- Cloke, H.L., Pappenberger, F., 2008. Evaluating forecasts of extreme events for hydrological applications: an approach for screening unfamiliar performance measures. *Meteorol. Appl.* 15 (1), 181–197. <http://dx.doi.org/10.1002/met.58>.
- Cooper, E.S., Dance, S.L., García-Pintado, J., Nichols, N.K., Smith, P.J., 2019. Observation operators for assimilation of satellite observations in fluvial inundation forecasting. *Hydrol. Earth Syst. Sci.* 23 (6), 2541–2559. <http://dx.doi.org/10.5194/hess-23-2541-2019>.
- Copernicus Programme, 2021. Copernicus Emergency Management Service. <https://emergency.copernicus.eu/>. (last Accessed 14 September 2021).
- Dasgupta, A., Grimaldi, S., Ramsankaran, R.A., Pauwels, V.R., Walker, J.P., 2018. Towards operational SAR-based flood mapping using neuro-fuzzy texture-based approaches. *Remote Sens. Environ.* 215 (June), 313–329. <http://dx.doi.org/10.1016/j.rse.2018.06.019>.
- Dasgupta, A., Hostache, R., Ramsankaran, R., Schumann, G.J.-P., Grimaldi, S., Pauwels, V.R.N., Walker, J.P., 2021a. On the impacts of observation location, timing and frequency on flood extent assimilation performance. *Water Resour. Res.* <http://dx.doi.org/10.1029/2020wr028238>.
- Dasgupta, A., Hostache, R., Ramsankaran, R.A., Schumann, G.J., Grimaldi, S., Pauwels, V.R., Walker, J.P., 2021b. A mutual information-based likelihood function for particle filter flood extent assimilation. *Water Resour. Res.* 57 (2), 1–28. <http://dx.doi.org/10.1029/2020WR027859>.
- Davies, P.A., McCarthy, M., Christidis, N., Dunstone, N., Fereday, D., Kendon, M., Knight, J.R., Scate, A.A., Sexton, D., 2021. The wet and stormy UK winter of 2019/2020. *Weather* 76 (12), 396–402. <http://dx.doi.org/10.1002/wea.3955>, URL: <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/wea.3955>.
- Dey, S.R., Leoncini, G., Roberts, N.M., Plant, R.S., Migliorini, S., 2014. A spatial view of ensemble spread in convection permitting ensembles. *Mon. Weather Rev.* <http://dx.doi.org/10.1175/MWR-D-14-00172.1>.
- Dey, S.R., Roberts, N.M., Plant, R.S., Migliorini, S., 2016. A new method for the characterization and verification of local spatial predictability for convective-scale ensembles. *Q. J. R. Meteorol. Soc.* <http://dx.doi.org/10.1002/qj.2792>.
- Di Mauro, C., Hostache, R., Matgen, P., Pelich, R., Chini, M., van Leeuwen, P.J., Nichols, N., Blöschl, G., 2020. Assimilation of probabilistic flood maps from SAR data into a hydrologic-hydraulic forecasting model: A proof of concept. *Hydrol. Earth Syst. Sci. Discuss.* (September), 1–24. <http://dx.doi.org/10.5194/hess-2020-403>.
- ECMWF, 2022. Skill scores of forecasts of weather parameters by TIGGE centres. <https://www.ecmwf.int/en/forecasts/charts>. (last Accessed 20 April 2022).
- Emerton, R.E., Stephens, E.M., Pappenberger, F., Pagano, T.C., Weerts, A.H., Wood, A.W., Salamon, P., Brown, J.D., Hjerdt, N., Donnelly, C., Baugh, C.A., Cloke, H.L., 2016. Continental and global scale flood forecasting systems. *Wiley Interdiscip. Rev.: Water* 3 (3), 391–418. <http://dx.doi.org/10.1002/wat2.1137>.
- Environment Agency, 2021. National LIDAR Programme. <https://data.gov.uk/dataset/f0db0249-f17b-4036-9e65-309148c97ce4/national-lidar-programme>. (last Accessed 29 April 2021).
- García-Pintado, J., Mason, D.C., Dance, S.L., Cloke, H.L., Neal, J.C., Freer, J., Bates, P.D., 2015. Satellite-supported flood forecasting in river networks: A real case study. *J. Hydrol.* 523, 706–724. <http://dx.doi.org/10.1016/j.jhydrol.2015.01.084>.
- GFM, 2021. GloFAS Global Flood Monitoring (GFM). <https://www.globalfloods.eu/technical-information/glofas-gfm/>. (last Accessed 28 October 2021).
- Google Earth Engine CART, 2021. ee.Classifier.smileCart. <https://developers.google.com/earth-engine/apidocs/ee-classifier-smilecart>. (last Accessed 29 April 2021).
- Google Earth Engine Catalog, 2021. Sentinel collection. https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S1_GRD. (last Accessed 4 August 2021).
- Google Earth Engine Scale, 2021. Image Pyramids. <https://developers.google.com/earth-engine/guides/scale>. (last Accessed 16 September 2021).
- Griffith, H.V., Wade, A.J., Lavers, D.A., Watts, G., 2020. Atmospheric river orientation determines flood occurrence. *Hydrol. Process.* 34 (23), 4547–4555. <http://dx.doi.org/10.1002/hyp.13905>.
- Grimaldi, S., Li, Y., Pauwels, V.R., Walker, J.P., 2016. Remote sensing-derived water extent and level to constrain hydraulic flood forecasting models: Opportunities and challenges. *Surv. Geophys.* 37 (5), 977–1034. <http://dx.doi.org/10.1007/s10712-016-9378-y>.
- Hagen, A., 2003. Fuzzy set approach to assessing similarity of categorical maps. *Int. J. Geogr. Inf. Sci.* 17 (3), 235–249. <http://dx.doi.org/10.1080/13658810210157822>.
- Hooker, H., 2022. Spatial scale evaluation of forecast flood inundation maps (v1.0) Zenodo. <http://dx.doi.org/10.5281/zenodo.6011881>. (Accessed 8 February 2022).
- Hostache, R., Chini, M., Giustarini, L., Neal, J., Kavetski, D., Wood, M., Corato, G., Pelich, R.M., Matgen, P., 2018. Near-real-time assimilation of SAR-derived flood maps for improving flood forecasts. *Water Resour. Res.* 54 (8), 5516–5535. <http://dx.doi.org/10.1029/2017WR022205>.
- Hostache, R., Martinis, S., Bauer-Marschallinger, B., Chini, M., Chow, S., Pelich, R., Li, Y., Böhnke, C., Knopp, L., Roth, F., Wieland, M., Wagner, W., Matgen, P., McCormick, N., Salamon, P., 2021. A first evaluation of the future CEMS systematic global flood monitoring product. <https://events.ecmwf.int/event/222/contributions/2274/attachments/1280/2347/Hydrological-WS-Hostache.pdf>. (last Accessed 4 August 2021).
- Janjić, T., Bormann, N., Bocquet, M., Carton, J.A., Cohn, S.E., Dance, S.L., Losa, S.N., Nichols, N.K., Potthast, R., Waller, J.A., Weston, P., 2018. On the representation error in data assimilation. *Q. J. R. Meteorol. Soc.* 144 (713), 1257–1278. <http://dx.doi.org/10.1002/qj.3130>, URL: <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3130>.
- JBA, 2021. Storm Ciara, Dennis and Jorge. <https://www.jbarisk.com/flood-services/event-response/storm-ciara-dennis-and-jorge/>. (last Accessed 14 September 2021).
- Kendon, M., 2020. Storm dennis. https://www.metoffice.gov.uk/binaries/content/assets/metofficegovuk/pdf/weather/learn-about/uk-past-events/interesting/2020_03_storm_dennis.pdf. (last Accessed 29 April 2021).
- Lavers, D.A., Allan, R.P., Wood, E.F., Villarini, G., Brayshaw, D.J., Wade, A.J., 2011. Winter floods in Britain are connected to atmospheric rivers. *Geophys. Res. Lett.* 38 (23), 1–8. <http://dx.doi.org/10.1029/2011GL049783>.
- Mason, D.C., Bevington, J., Dance, S.L., Revilla-Romero, B., Smith, R., Vetrá-Carvalho, S., Cloke, H.L., 2021b. Improving urban flood mapping by merging synthetic aperture radar-derived flood footprints with flood hazard maps. *Water (Switzerland)* 13 (11), <http://dx.doi.org/10.3390/w13111577>.
- Mason, D.C., Dance, S.L., Cloke, H.L., 2021a. Floodwater detection in urban areas using sentinel-1 and WorldDEM data. *J. Appl. Remote Sens.* 15 (03), 1–22. <http://dx.doi.org/10.1117/1.jrs.15.032003>.
- Mason, D.C., Dance, S.L., Vetrá-Carvalho, S., Cloke, H.L., 2018. Robust algorithm for detecting floodwater in urban areas using synthetic aperture radar images. *J. Appl. Remote Sens.* 12 (04), 1. <http://dx.doi.org/10.1117/1.jrs.12.045011>, URL: <http://centaur.reading.ac.uk/80110/8/045011.1.pdf>.
- Met Office, 2020. Record breaking rainfall. <https://www.metoffice.gov.uk/about-us/press-office/news/weather-and-climate/302020/2020-winter-february-stats>. (last Accessed 29 April 2021).
- Met Office, 2022. Weather Observation Website (WOW). <https://wow.metoffice.gov.uk/observations/details/20220503d5q856sksh63xiyyb96sm40y>. (last Accessed 4 May 2022).
- National River Flow Archive, 2021. NRFA. <https://nrfa.ceh.ac.uk/data/station/info/55002>. (last Accessed 29 April 2021).
- ndimage.zoom, 2021. Scipy. <https://docs.scipy.org/doc/scipy/reference/generated/scipy.ndimage.zoom.html>. (last Accessed 21 September 2021).
- Pappenberger, F., Frodsham, K., Beven, K., Romanowicz, R., Matgen, P., 2007. Fuzzy set approach to calibrating distributed flood inundation models using remote sensing observations. *Hydrol. Earth Syst. Sci.* 11 (2), 739–752. <http://dx.doi.org/10.5194/hess-11-739-2007>.
- Revilla-Romero, B., Shelton, K., Wood, E., Berry, R., Bevington, J., Hankin, B., Lewis, G., Gubbin, A., Griffiths, S., Barnard, P., Pinnell, M., Huyck, C., 2017. Flood foresight: A near-real time flood monitoring and forecasting tool for rapid and predictive flood impact assessment. In: *EGU General Assembly Conference Abstracts*. p. 1230.
- riverlevels.uk, 2020. River levels, river wye at hereford bridge. <https://riverlevels.uk/herefordshire-hereford-old-wye-bridge-lvl#.YIFkT31KgUE>. (last Accessed 29 April 2021).
- Roberts, N.M., Lean, H.W., 2008. Scale-selective verification of rainfall accumulations from high-resolution forecasts of convective events. *Mon. Weather Rev.* <http://dx.doi.org/10.1175/2007MWR2123.1>.
- Savage, J.T.S., Bates, P., Freer, J., Neal, J., Aronica, G., 2016. When does spatial resolution become spurious in probabilistic flood inundation predictions? *Hydrol. Process.* 30 (13), 2014–2032. <http://dx.doi.org/10.1002/hyp.10749>.

- Schumann, G.J., 2019. The need for scientific rigour and accountability in flood mapping to better support disaster response. *Hydrol. Process.* 33 (24), 3138–3142. <http://dx.doi.org/10.1002/hyp.13547>.
- Schumann, G., Bates, P.D., Horritt, M.S., Matgen, P., Pappenberger, F., 2009. Progress in integration of remote sensing-derived flood extent and stage data and hydraulic models. *Rev. Geophys.* <http://dx.doi.org/10.1029/2008RG000274>.
- Sefton, C., Muchan, K., Parry, S., Matthews, B., Barker, L.J., Turner, S., Hannaford, J., 2021. The 2019/2020 floods in the UK: A hydrological appraisal. *Weather* 76 (12), 378–384. <http://dx.doi.org/10.1002/wea.3993>, URL: <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/wea.3993>.
- Skok, G., Roberts, N., 2016. Analysis of fractions skill score properties for random precipitation fields and ECMWF forecasts. *Q. J. R. Meteorol. Soc.* 142 (700), 2599–2610. <http://dx.doi.org/10.1002/qj.2849>.
- Speight, L.J., Cranston, M.D., White, C.J., Kelly, L., 2021. Operational and emerging capabilities for surface water flood forecasting. *Wiley Interdiscip. Rev.: Water* 8 (3), 1–24. <http://dx.doi.org/10.1002/wat2.1517>.
- Stein, J., Stoop, F., 2019. Neighborhood-based contingency tables including errors compensation. *Mon. Weather Rev.* 147 (1), 329–344. <http://dx.doi.org/10.1175/MWR-D-17-0288.1>.
- Stephens, E., Schumann, G., Bates, P., 2014. Problems with binary pattern measures for flood model evaluation. *Hydrol. Process.* <http://dx.doi.org/10.1002/hyp.9979>.
- Trepekli, K., Friberg, T., Balstrøm, T., Fog, B., Allotey, A., Kofie, R., Møller-Jensen, L., 2021. UAV-Lidar observations increase the precision of urban flood modelling in accra by detecting critical micro-topographic features. <http://dx.doi.org/10.5194/egusphere-egu21-10457>, EGU General Assembly Online 19–30 Apr 2021.
- UK Water Resources Portal, 2022. NRFA UKceh UK water resources portal. <https://eip.ceh.ac.uk/hydrology/water-resources/>. (last Accessed 4 May 2022).
- Ulbrich, U., Leckebusch, G.C., Pinto, J.G., 2009. Extra-tropical cyclones in the present and future climate: A review. *Theor. Appl. Climatol.* 96 (1–2), 117–131. <http://dx.doi.org/10.1007/s00704-008-0083-8>.