



Evaluating and improving flood inundation forecasts using satellite data

HELEN HOOKER

A thesis submitted in fulfilment of the requirements for the degree of
Doctor of Philosophy

PhD Atmosphere, Oceans and Climate

Department of Meteorology

School of Mathematical, Physical and Computational Sciences

January 2024

University of Reading

Declaration

I confirm that this is my own work and the use of all material from other sources has been properly and fully acknowledged.

Helen Hooker

Publications

Chapters 3, 4 and 5 of this thesis are reproduced from the following publications respectively:

- Hooker, H., Dance, S. L., Mason, D. C., Bevington, J., & Shelton, K. (2022). Spatial scale evaluation of forecast flood inundation maps. *Journal of Hydrology*, 128170. <https://doi.org/https://doi.org/10.1016/j.jhydrol.2022.128170>
- Hooker, H., Dance, S. L., Mason, D. C., Bevington, J., & Shelton, K. (2023). Assessing the spatial spread–skill of ensemble flood maps with remote-sensing observations. *Natural Hazards and Earth System Sciences*, 23(8), 2769–2785. <https://doi.org/10.5194/nhess-23-2769-2023>
- Hooker, H., Dance, S. L., Mason, D. C., Bevington, J., & Shelton, K. (2023). A multi-system comparison of forecast flooding extent using a scale-selective approach. *Hydrology Research*, 54(10), 1115–1133. <https://doi.org/10.2166/nh.2023.025>

All work undertaken in these publications was carried out by Helen Hooker with coauthors providing guidance and review.

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. The maps can be used to inform forecast-based financing schemes to release funds ahead of a predicted flood event. Evaluating the accuracy of forecast flood maps is essential for model development and improving future flood predictions. The goal of this thesis is to develop spatial verification methods for deterministic and ensemble flood map forecasts and to improve forecasts using satellite data. Binary verification measures typically provide a domain-averaged score of forecast skill. The skill score is dependent on the magnitude of the flood and the spatial scale of the flood map. In this thesis, a new scale-selective approach is presented to evaluate both deterministic and ensemble forecast flood maps against remotely observed flood extents. 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. Both the ensemble spatial-skill and spread-skill relationship vary with location and can be linked to the physical characteristics of the flooding event. We find that a scale-selective verification approach can quantify the skill of three systems operating at different spatial scales, so that the benefits and limitations of each system can be evaluated. A new data assimilation framework is presented to update the flood map selection from a static library of flood maps using satellite data, taking account of observation uncertainties. Results show that the flood map selection could be

triggered in four out of five sub-catchments tested. The resultant analysis flood map has the potential to be used to trigger a secondary finance scheme during a flood event and avoid missed financing opportunities for humanitarian action. Overall, sensitive spatial verification methods that are location specific and can evaluate ensemble performance will aid future model development for flood inundation prediction.

Dedication

For Richard, Ebony, Kian and Maple

Acknowledgements

I would like to thank my previous work colleagues for their support and encouragement to apply for and pursue a PhD. Thank you to my lead supervisor Sarah Dance for believing I could do it, trusting in my ideas and correcting my grammar (no more greengrocer's apostrophe)! Thank you to David Mason for his wisdom and support. My work always felt meaningful because of the direct links to my CASE sponsor the JBA Trust and it was great to work with John Bevington and Kay Shelton throughout, thank you for hosting me at your Wallingford office too.

It was always reassuring to be part of something bigger, thank you to SCENARIO DTP and Wendy Neale along with NCEO and Uzma Saeed. I'm grateful for the fun training opportunities both provided. It was great to be part of Cohort 7 PhD students in Meteorology, and thank you to my office buddies Charlie and Natalie for listening to me ramble on. Thank you to DARC for the challenging seminars and to my fellow PhD DARCies Gwyn, Devon, Harriet, Ieuan and Laura. Thank you to Hugo Dalton for sharing his art. I'd also like to thank Hannah Cloke and Liz Stephens for the Water group seminars, journal club meetings and watery activities.

Finally, I'd like to thank my Mum and Dad for teaching me to keep on learning and to Rich, Ebony and Kian for encouraging me to get on with it and keep on going. Thank you to Maple (my English Springer) for her company throughout.

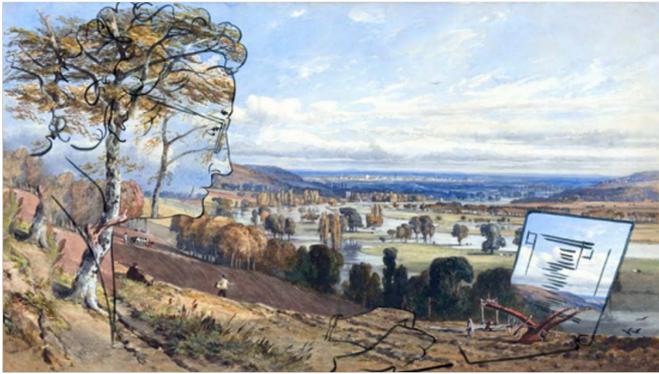
Assimilated Watercolours: Pop up art exhibitions in Care Homes

During my PhD I took part in a sci-art collaboration. The project created a series of exhibitions that explore current research at the University of Reading's Department of Meteorology through historical and contemporary works of art. Scientists at Reading aim to develop new ways to make flood forecasting more accurate by assimilating visual and mathematical data from sources including satellite imagery, aerial photographs and existing flooding maps.

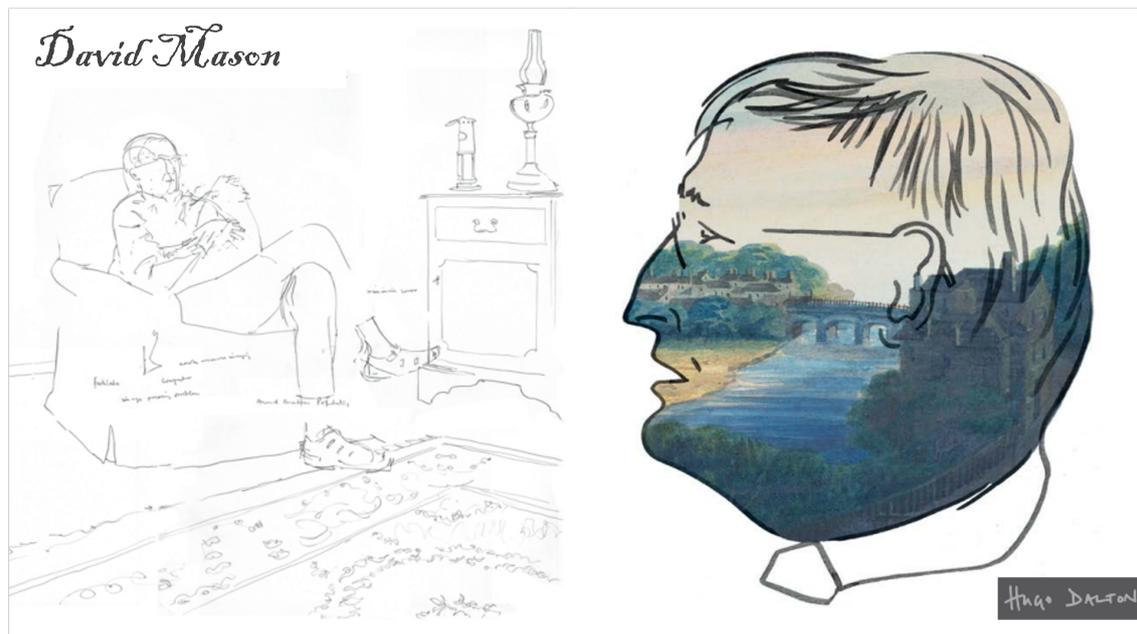
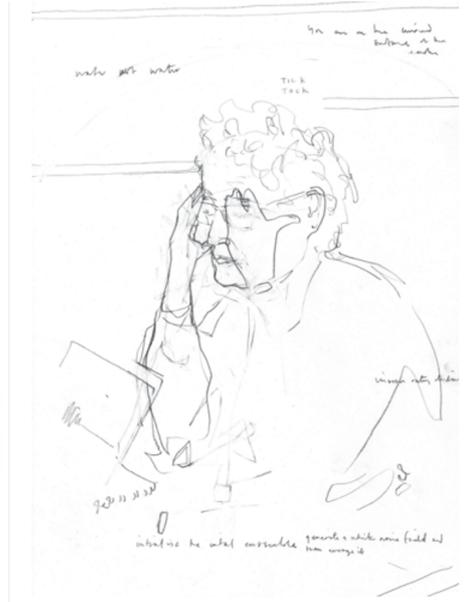
The artist Hugo Dalton has made a series of watercolour paintings of the research locations and the scientists at work. These have been merged with existing reproductions of historic artworks from the Royal Collection at nearby Windsor, which also show views of the surrounding area. The idea is to show how different kinds of information are used to create a better understanding of flooding (Dance et al., 2022). The objective is to offer elderly people in care homes a connection to current science research in their local area. The project provides them with a link to places they may no longer be able to physically access.

A selection of the artwork is included in this thesis as chapter illustrations with permission granted by Hugo Dalton.

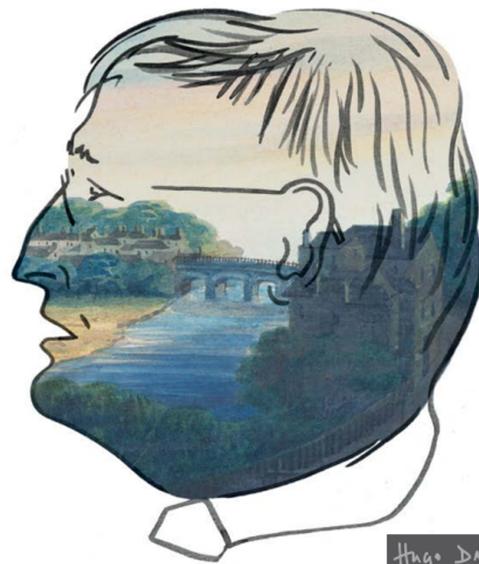
Sarah Dance



Hugo DALTON



David Mason



Hugo DALTON

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Motivation | 1 |
| 1.2 | Thesis aims | 3 |
| 1.3 | Principal new results | 4 |
| 1.4 | Thesis outline | 5 |
| 2 | Background | 7 |
| 2.1 | Fluvial simulation library flood inundation forecasting systems | 7 |
| 2.2 | Observing flooding from Space | 9 |
| 2.3 | Flood extent verification | 10 |
| 2.4 | Introduction to data assimilation | 13 |
| 2.5 | Chapter summary | 14 |
| 3 | Spatial scale evaluation of forecast flood inundation maps | 16 |
| 3.1 | Abstract | 17 |
| 3.2 | Introduction | 19 |
| 3.3 | Flood event | 23 |
| 3.3.1 | February 2020 | 23 |
| 3.3.2 | Catchment location and description | 24 |
| 3.3.2.1 | The River Wye (domains A and B) | 25 |

| | | |
|----------|---|-----------|
| 3.3.2.2 | River Lugg at Lugwardine (domain C) | 25 |
| 3.3.2.3 | Event hydrology | 25 |
| 3.4 | Data | 27 |
| 3.4.1 | Flood Foresight | 27 |
| 3.4.2 | SAR-derived flood maps | 28 |
| 3.5 | Flood map evaluation methods | 30 |
| 3.5.1 | Spatial scale-selective approach | 31 |
| 3.5.2 | Location dependent agreement scales | 34 |
| 3.5.3 | Categorical scale map | 35 |
| 3.5.4 | Binary performance measures | 36 |
| 3.6 | Results | 37 |
| 3.6.1 | Spatial scale variability of forecast flood extent and flood-edge location | 37 |
| 3.6.2 | Comparison of spatial scales at differing lead times and domain lo- cation | 39 |
| 3.6.3 | Categorical scale maps | 43 |
| 3.6.4 | SAR-derived flood map scale variation | 46 |
| 3.7 | Discussion and Conclusions | 48 |
| 3.8 | Chapter summary | 51 |
| 4 | Assessing the spatial spread-skill of ensemble flood maps with remote sensing observations | 53 |
| 4.1 | Abstract | 54 |
| 4.2 | Introduction | 55 |
| 4.3 | Ensemble flood map spatial predictability evaluation methods | 60 |
| 4.3.1 | Fraction Skill Score | 61 |
| 4.3.2 | Location dependent agreement scales | 63 |
| 4.3.3 | Ensemble spatial spread-skill evaluation | 65 |

| | | |
|----------|--|-----------|
| 4.3.4 | Spatial spread-skill visualisation methods | 66 |
| 4.4 | Ensemble forecasting flood event case study | 68 |
| 4.4.1 | Brahmaputra flood, Assam India, August 2017 | 69 |
| 4.4.2 | Ensemble flood forecasting system | 71 |
| 4.4.3 | SAR-derived flood map | 74 |
| 4.4.4 | Forecast data | 75 |
| 4.5 | Results and discussion | 77 |
| 4.5.1 | Ensemble spatial scale evaluation | 78 |
| 4.5.2 | Ensemble spatial predictability | 79 |
| 4.5.3 | Ensemble spatial spread-skill | 80 |
| 4.6 | Conclusions | 85 |
| 4.7 | Chapter summary | 87 |
| 5 | A multi-system comparison of forecast flood extent using a scale-selective approach | 88 |
| 5.1 | Abstract | 89 |
| 5.2 | Introduction | 90 |
| 5.3 | Flood event on the Jamuna River, Bangladesh July 2020 | 93 |
| 5.4 | Flood forecasting systems and data | 96 |
| 5.4.1 | GloFAS Rapid Flood Mapping | 98 |
| 5.4.2 | Flood Foresight | 99 |
| 5.4.3 | Bangladesh Flood Forecasting and Warning Centre | 102 |
| 5.4.4 | Observation data | 103 |
| 5.5 | Scale-selective evaluation methods | 104 |
| 5.6 | Results and discussion | 106 |
| 5.6.1 | Flood Foresight Jamuna River case study | 107 |
| 5.6.2 | Multi-system flood map comparison | 112 |

| | | |
|----------|---|------------|
| 5.6.3 | Discussion | 116 |
| 5.7 | Conclusions | 119 |
| 5.8 | Appendix | 122 |
| 5.9 | Chapter summary | 125 |
| 6 | Updating simulation library flood map selection through assimilation of probabilistic SAR-derived flood extent | 126 |
| 6.1 | Abstract | 128 |
| 6.2 | Introduction | 129 |
| 6.3 | Simulation library forecasting system and observation data | 134 |
| 6.3.1 | Flood Foresight and Forecast-based-Financing | 134 |
| 6.3.2 | Satellite-derived flood likelihood | 135 |
| 6.3.3 | Optical Normalized Difference Water Index (NDWI) | 137 |
| 6.4 | Methods | 137 |
| 6.4.1 | Data assimilation framework | 137 |
| 6.4.2 | Validation methods | 139 |
| 6.5 | Pakistan flood 2022 | 141 |
| 6.5.1 | Event overview | 141 |
| 6.5.2 | Data | 141 |
| 6.6 | Results and discussion | 144 |
| 6.6.1 | Scenario 1 | 144 |
| 6.6.2 | Scenario 2 | 145 |
| 6.6.3 | Scenario 3 | 147 |
| 6.6.4 | Analysis flood map validation | 149 |
| 6.7 | Conclusions | 150 |
| 6.8 | Chapter summary | 152 |

| | | |
|----------|---|------------|
| 7 | Conclusions | 153 |
| 7.1 | Main conclusions | 153 |
| 7.2 | Thesis synthesis | 156 |
| 7.3 | Limitations | 157 |
| 7.4 | Recommendations for future work | 158 |
| | References | 160 |

List of Figures

| | | |
|-----|---|----|
| 2.1 | Data assimilation cycle | 13 |
| 3.1 | 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 x 12.8 km. Domain B: 5.8 km of the River Wye at Hereford, domain size 3.0 x 4.0 km. Domain C: 4 km of the River Lugg at Lugwardine, domain size 2.3 x 2.3 km. Base map from Google Maps. | 24 |
| 3.2 | Daily maximum river levels (m) at Ross-on-Wye, Hereford and Lugwardine. The dashed yellow line indicates Sentinel-1 SAR acquisition date. | 27 |
| 3.3 | Flood Foresight flood map simulation library selection process. Source JBA Consulting. | 29 |
| 3.4 | FSS (see subsection 3.5.1 for calculation details) example applied to a binary flooded (1) / unflooded (0) field at grid scale (yellow box, $n = 1$) and a 3 x 3 neighbourhood (black box, $n = 3$). The observed SAR-derived forecast is in turquoise and the forecast is shown in blue. | 32 |

| | | |
|-----|--|----|
| 3.5 | Left panel: contingency map of a 0-day lead time forecast verses 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. | 38 |
| 3.6 | 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. | 39 |
| 3.7 | 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. | 41 |
| 3.8 | 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 x 25 m for all domains. (Note: rd (forecast run date) varies between location, all dates have been evaluated and the most illustrative maps selected.) | 44 |
| 3.9 | 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. | 47 |

| | | |
|------|---|----|
| 3.10 | 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. | 48 |
| 4.1 | Figure reproduced with permission from Dey et al., (2016) showing results on a binned scatter plot from an idealised experiment indicating the spatial spread-skill relationship between an ensemble forecast and the observation. | 68 |
| 4.2 | Left panel: domain location on the Brahmaputra River in the Assam region of India. Domain size is 57.5 km by 39.3 km. Right panel: Sentinel-1 SAR-derived flood map and permanent water bodies from the JRC Global Surface Water database for the domain of interest (DOI). Base map from ©Google Maps. | 70 |
| 4.3 | Flood Foresight ensemble forecast flood inundation and impact mapping work flow. Prepared by JBA Consulting. | 72 |
| 4.4 | GloFAS grid, permanent water bodies and Flood Foresight sub-catchments for the domain of interest (DOI). | 74 |
| 4.5 | Brahmaputra River, Assam region, August 2017. 51 ensemble member forecast flood maps (labelled 0 to 50), ens_{median} and ens_{all} all at 1-day lead time and the Sentinel-1 SAR-derived flood map. | 76 |
| 4.6 | Brahmaputra River, Assam region, August 2017. Colour shading from white (low) to dark blue (high) indicate the forecast probability of flooding based on a 1-day lead time, 51 ensemble member flood map forecast for the Brahmaputra River in the Assam region, August 2017. (Note map background is grey) | 77 |

- 4.7 The spatial skill of each individual ensemble member forecast flood extent is evaluated along with the ens_{median} (a spatial median where 26 or more members predict flooding at a grid cell location) and ens_{all} (flooded grid cells from all ensemble members are combined). The FSS is calculated at increasing neighbourhood sizes to determine the scale at which the forecast becomes skilful at capturing the observed flood (FSS_T). 79
- 4.8 Brahmaputra River, Assam region, August 2017. Categorical scale maps for (a) ens_{all} (flooded grid cells from all ensemble members are combined), (b) ens_{median} (a spatial median where 26 or more members predict flooding at a grid cell location), (c) individual ensemble member 1 and (d) individual ensemble member 21. Red areas indicate where the forecast is under-predicted and blue regions represent over-prediction. The colour shade gives the scale of agreement (Eq. (7)) between the forecast and the observed flooding with lighter shading indicating a smaller agreement scale is required to reach the agreement criterion (Eq. (6)), a fixed maximum scale S_{lim} is drawn to scale (c). For georeferencing see Figure 4.6, each grid cell is 30 m x 30 m. 81
- 4.9 Brahmaputra River, Assam region, August 2017. (a) The average agreement scale map of each unique pair of forecast ensemble flood maps and (b) between each ensemble member compared against the observed SAR-derived flood map. (c) A binned histogram scatter plot compares (a) and (b) to indicate the spatial spread-skill of the forecast ensemble. (d) indicates the corresponding sub-catchment locations. Areas labelled (1, 2 and 3) are discussed in Section 4.5.3. A fixed maximum scale S_{lim} (Eq. (6)) is drawn to scale (a). Note PWB means permanent water bodies. 83

| | | |
|------|--|-----|
| 4.10 | Brahmaputra River, Assam region, August 2017. (a) The Spatial Spread-Skill (SSS) map shows the difference between the ensemble/ensemble and the ensemble/observed average agreement scales at each grid cell. Negative values (orange) indicate where the ensemble is under-spread and positive values (purple) indicate where the ensemble is over-spread. White areas indicate where the average agreement scales match and indicate good spatial spread-skill. (d) Indicates the corresponding sub-catchment locations. Areas labelled (1, 2 and 3) are discussed in Section 4.5.3. A fixed maximum scale S_{lim} (Eq. (6)) is drawn to scale (a). Note PWB means permanent water bodies. | 84 |
| 5.1 | Four districts (zilas) of interest in the Jamuna catchment in northern Bangladesh. | 95 |
| 5.2 | Flood Foresight/Start Network ensemble flood inundation forecast and population impacts work flow. | 100 |
| 5.3 | Example Flood Foresight forecast domain divided into Impact Zones (shown in colour with black outline where the 2-year return period threshold is exceeded), each linked to a GloFAS grid cell. The Impact Zone colour shows the corresponding return period threshold exceeded, determined by the GloFAS forecast discharge. | 101 |
| 5.4 | An example FSS calculation applied to forecast flood extent, 1 = flooded (in blue), 0 = unflooded (in white) compared to a remotely observed flood extent in pink. The FSS is calculated for two neighbourhood sizes, $n = 1$ (small gold box) and $n = 3$ (large gold box). | 105 |
| 5.5 | Flood Foresight average agreement scale against forecast lead time for each district and updated Jamalpur following reassociation of IZ with GloFAS grid cells, forecast valid for 25 July. | 107 |

| | | |
|------|---|-----|
| 5.6 | (a) GloFAS forecast discharge (control member, 1-day lead time) compared to FFWC observed river water level for the main Jamuna channel at Bahadurabad and (b) the old Brahmaputra distributary in the Jamalpur district. (c) The old Brahmaputra distributary forecast discharge following reassociation of IZ with GloFAS grid cells. The GloFAS RP threshold levels are taken from the nearest GloFAS grid cell to the gauge station location. Station risk levels are provided by FFWC. | 110 |
| 5.7 | (a) Original CSM for Jamalpur. (b) CSM change (updated CSM - original CSM) for Jamalpur following reassociation of IZ with GloFAS grid cells. Run date 20 July, forecast valid for 25 July for (a) and (b). | 111 |
| 5.8 | CSM for flood inundation forecasts from three forecast systems for flood peak 25 July compared to SAR-derived flooding. (a) Flood Foresight run date 18 July, (b) GloFAS RFM run date 18 July and (c) FFWC flood map run date 25 July. | 113 |
| 5.9 | Average FSS plotted against neighbourhood size (n) for each forecast system (run dates as described in Figure 5.8) in Kurigram (a), Gaibandha (b), Sirajganj (c) and Jamalpur (d) and the target skill score for each district. . | 115 |
| 5.10 | Six FFWC station water levels during July and August 2020 across the four districts compared with closest GloFAS grid cell discharge. | 123 |
| 5.11 | Binary performance measures for each district. Flood Foresight run date 18 July (a), GloFAS RFM 18 July (b) and FFWC flood map 25 July (c). . | 124 |
| 6.1 | Flood Foresight/Start Network ensemble flood inundation forecast and population impacts work flow. | 135 |

6.2 The domain of interest (DOI) is located on the Indus Basin, Sindh province, Pakistan (left). The region is divided into sub-catchments or Impact Zones (IZ) in Flood Foresight (right). Satellite-derived flooding (NDWI) from Sentinel-2 data (Section 6.3.3) from 31 August 2022 is highlighted along with permanent water bodies (PWB). 142

6.3 (a) GFM flood likelihood derived from Sentinel-1 SAR data, masked areas (grey) indicate where flooding cannot be reliably detected from SAR data. (b) The maximum return period threshold triggered per IZ by the Flood Foresight system during peak flooding 10-31 August, 2022. Five non-triggered IZ of interest labelled: S - Sukkur, LN - Larkana North, LS - Larkana South, FE1 - flood edge 1, FE2 - flood edge 2. 143

6.4 Scenario 1: Sukkur, a dense urban area. (a) The cost function (eqn. (4)) plotted against the RP value, the yellow star highlights the minimum value of $J(\mathbf{x})$. (b) The analysis RP triggered following DA. (c) The analysis flood extent map (note that no map was triggered for Sukkur) and (d) the CSM comparing the analysis flood extent map against Sentinel-2 NDWI. 145

6.5 Scenario 2: Larkana, a mixed urban and rural area. (a) The cost function (eqn. (4)) plotted against the RP value, the yellow star highlights the minimum value of $J(\mathbf{x})$. (b) The analysis RP triggered following DA. (c) The analysis flood extent map and (d) the CSM comparing the analysis flood extent map against Sentinel-2 NDWI. 146

6.6 Scenario 3: Flood edge location. (a) The cost function (eqn. (4)) plotted against the RP value, the yellow star highlights the minimum value of $J(\mathbf{x})$. (b) The analysis RP triggered following DA. (c) The analysis flood extent map and (d) the CSM comparing the analysis flood extent map against Sentinel-2 NDWI. 148

1 Chapter 1

2 Introduction

3 1.1 Motivation

4 Globally, an estimated 1.8 billion people (23% of the world's population) live in areas that
5 are directly exposed to a 100-year return period flood (Rentschler et al., 2022). The vast
6 majority (89%) of these people live in low- and middle-income countries where infrastruc-
7 ture systems, including flood protection and drainage, and early warning systems, tend
8 to be less developed. Satellite data shows that the number of people exposed to flooding
9 has increased by 20 to 24% globally from 2000 to 2018. Increasing exposure to flooding is
10 predicted to continue with climate change (Tellman et al., 2021). Early warning systems
11 can significantly improve the outcomes following disasters, reducing deaths and damage
12 and enabling faster recovery (UNDRR, 2022a). The 2023 UN Global Assessment Re-
13 port (GAR) on Disaster Risk Reduction (UNDRR, 2023) states that the benefits of early
14 warning systems triple in vulnerable contexts. Despite these benefits, the GAR2022 report
15 shows that just 5.8% (\$5.5 billion USD) of official development assistance contributes to
16 disaster prevention and preparedness compared to 90.1% (\$119.8 billion USD) for emer-
17 gency response. Yet it has been demonstrated (for Europe) that financing for mitigation
18 purposes such as flood forecasting systems can lead to overall cost savings (Pappenberger

1 et al., 2015).

2

3 Global-scale flood forecasting systems can support the early warnings for all initiative
4 (UNDRR, 2022a) by providing flood forecasts for large rivers around the world. Advances
5 in flood forecasting systems link together meteorological and hydrological forecasts to hy-
6 drodynamic models, simulating flood-wave propagation (Emerton et al., 2016; Wu et al.,
7 2020; Apel et al., 2022). A simulation library forecasting system saves computation time
8 by storing static flood extent and depth maps at various return periods. Depending on
9 the forecast river discharge, the maps are looked-up per sub-catchment and mosaicked
10 together. The resulting flood maps can be used to inform disaster risk reduction schemes
11 such as forecast-based financing (FbF). FbF works by quantifying risks in advance of crises
12 or disasters, prepositioning funds, and agreeing in advance how funds will be released based
13 on forecasts, ahead of an event (OCHA, 2020). This results in global-scale models being
14 used to inform local-scale action. Hoch and Trigg (2019) outline a Global Flood Model
15 Validation Framework, which includes a recommendation to routinely validate flood ex-
16 tent. Quantitative performance evaluation forms an important part of fitness-for-purpose
17 assessment and continual system improvement. Currently, there is limited quantitative
18 validation of operational flood forecasting systems producing flood maps.

19

20 The accuracy of forecasts of flood extent can be verified by comparing with obser-
21 vations of flooding from drones or satellite-based sensors. Typically, binary performance
22 measures are calculated and provide an average measure of skill across a region (Stephens
23 et al., 2014). In this thesis we address several limitations of binary performance mea-
24 sures by applying a new scale-selective approach to flood map verification. The approach
25 is developed further to evaluate how well an ensemble flood map forecast represents the
26 uncertainties involved. We apply scale-selective verification methods to a multi-system
27 comparison where each of the systems' forecasts are presented at different spatial scales.

1 Some limitations to the flood forecasting system are addressed by developing a data as-
2 simulation framework using satellite data to improve the flood map analysis.

3

4 **1.2 Thesis aims**

5 The aims of this thesis are to address the following research questions:

6 **1. What are the skilful spatial scales in flood inundation forecasts made**
7 **using a simulation library approach?**

8 How can we determine the skilful spatial scales of forecast flood maps by comparing
9 against satellite-derived observations of flooding? Does the skilful spatial scale vary
10 with location and how can this be visualised? How does validation of the flood edge
11 location alone compare to validation of the entire flood extent? How can the skilful
12 spatial scale results be used in operational practice?

13 **2. How skilfully does an ensemble of forecast flood maps represent the spa-**
14 **tial uncertainty within the flood forecast?**

15 How can we summarise the spatial predictability information in ensemble flood map
16 forecasts? How can we evaluate the spatial spread-skill of an ensemble flood map
17 forecast? How does the spatial spread-skill vary with location and how can this be
18 presented?

19 **3. How useful are scale-selective evaluation approaches when applied to mul-**
20 **tiple flood forecasting systems?**

21 How can we evaluate the performance of flood forecasting systems predicting flood
22 inundation extent at different spatial scales? What can we learn about the flood
23 forecasting system performance and how does each compare?

1 **4. Does a data assimilation framework improve the analysis of flood inun-**
2 **dation from a simulation library system?**

3 Can we incorporate probabilistic information from remotely observed flood inunda-
4 tion into a data assimilation framework to improve the flood map selection within
5 a simulation library flood forecasting system? How does the analysis flood map
6 compare to independent validation data?

7 **1.3 Principal new results**

8 The outcomes of this thesis provide the following answers to the research questions:

- 9 1. A skilful spatial scale for forecast flood maps can be found by calculating the Frac-
10 tion Skill Score, a validation metric, found by comparing a deterministic forecast
11 flood map against a satellite SAR-derived observation of flooding across a range of
12 neighbourhood sizes. A target skill score can be calculated and this depends on the
13 magnitude of the observed flood. The skilful scale determined for the flood edge
14 is more sensitive to changes in spatial accuracy and spatial scale compared to the
15 skilful scale found by evaluating the entire flood extent. Categorical scale maps
16 developed show that the skilful scale varies with location across a domain.
- 17 2. We present a new scale-selective approach to assess the spatial predictability and
18 spread-skill of an ensemble flood map forecast that accounts for the individual spatial
19 prediction of flood extent held within each ensemble member flood map. The method
20 determines, at specific locations within the domain, whether the ensemble forecast
21 is over-, under- or well-spread. The spatial spread-skill relationship can be mapped
22 onto a Spatial spread-skill map.
- 23 3. We investigate a new application of scale-selective verification by evaluating the per-
24 formance of three flood forecasting systems. Two simulation library systems, Flood

1 Foresight (30 m) and GloFAS Rapid Flood Mapping (1000 m) and one hydrody-
2 namically modelled system, the Bangladesh FFWC Super Model (300 m), all made
3 predictions of flood extent at different spatial scales (grid lengths, shown in brack-
4 ets) for the Jamuna River flood, Bangladesh, July 2020. Our results show that the
5 simulation library system accuracy critically depends on the discharge return period
6 threshold set to trigger a flood map selection and the number of hydrological model
7 ensemble members that must exceed it.

8 4. A data assimilation (DA) framework is developed to integrate probabilistic flood
9 extent maps from satellite-based SAR sensors into the simulation library flood map
10 selection process. The method is tested on the severe flood event in Pakistan, 2022,
11 where several sub-catchments resulted in a non-trigger of the forecast-based financing
12 system deployed here, despite significant flooding evident from earth observation
13 data. The DA successfully triggered flood maps in 4 out of 5 sub-catchments tested
14 and we found that evaluating sub-catchments at the flood edge gave the best results.

15 1.4 Thesis outline

16 This thesis is structured as follows:

- 17 • Chapter 2 introduces relevant background information for the thesis including sim-
18 ulation library flood inundation forecasting systems, observing flooding from satel-
19 lites, flood extent verification and data assimilation approaches in flood inundation
20 forecasting.
- 21 • Chapter 3 addresses the first research question in Section 1.2. A new approach to
22 forecast flood map spatial verification against satellite-derived observations of flood-
23 ing is presented. A scale-selective verification method is applied to evaluate the
24 performance of a simulation library flood forecasting system at predicting flood ex-

1 tent on the Rivers Wye and Lugg following storm Dennis in February 2020. The
2 scale-selective approach addresses multiple limitations of conventional binary per-
3 formance measures and we find several applications of the evaluation approach that
4 benefit operational flood forecasting practice. Chapter 3 is reproduced from Hooker
5 et al. (2022).

6 • Chapter 4 addresses the second research question. We present a new approach to
7 evaluate and visualise the spatial spread-skill of an ensemble flood map forecast. The
8 method can be used to assess how well the probabilistic flood maps represent the
9 uncertainty present within the forecast-chain and how this may change with updates
10 to the forecasting system such as including additional observations. Chapter 4 is
11 reproduced from Hooker et al. (2023a).

12 • Chapter 5 addresses the third research question. Through application of scale-
13 selective evaluation methods we can directly compare forecast flood maps from three
14 flood forecasting systems, each predicting flood extent at different spatial scales (grid
15 lengths). This quantitative spatial validation means that the benefits and limitations
16 of the forecast systems can be evaluated. Chapter 5 is reproduced from Hooker et
17 al. (2023b).

18 • Chapter 6 addresses the final research question. We present a new data assimilation
19 framework to incorporate satellite-derived probabilistic flood extent information into
20 the simulation library flood map selection process. This overcomes limitations of the
21 simulation library system and has the potential to improve forecast-based financing
22 schemes.

23 • Chapter 7 summarises the main findings from the thesis and makes recommendations
24 for future work.

1 Chapter 2

2 Background

3 In this chapter we introduce some of the key topics used in this thesis. In Section 2.1 we
4 introduce simulation library flood inundation forecasting systems and their application in
5 disaster risk reduction. Observing flooding from satellite data is explained in Section 2.2
6 and using these observations to verify flood forecasts is discussed in Section 2.3. Data
7 assimilation approaches used in flood inundation forecasting are introduced in Section 2.4.

8 **2.1 Fluvial simulation library flood inundation forecasting** 9 **systems**

10 The current state-of-the-art in operational flood inundation forecasting at national or
11 transnational scales uses a simulation library system (Revilla-Romero et al., 2017). Flood
12 extent and depth maps are precomputed using a hydrodynamic model at a range of return
13 period thresholds. Together, these flood maps form a simulation library for a particular
14 country or river basin.

15

16 The system links together a chain of models that begins with a numerical weather
17 prediction (NWP) model providing meteorological inputs such as forecast precipitation to

1 a hydrological model and subsequent hydrodynamic model. The NWP model, combined
2 with recent observations also provides initial conditions for the hydrological model. The
3 forecast river discharge from the hydrological model, along with the return period thresh-
4 olds set (from historical observations or reanalysis data-sets of river discharge (Grimaldi,
5 2022)) determines which flood map from the simulation library is selected and presented
6 (for an example diagram see Figure 4.3).

7

8 Pre-computing the flood maps reduces the model run-time and means that the flood
9 forecasting system can operate in near real-time. Additionally, the system can handle
10 multiple inputs from ensemble NWP forecasts (Cloke & Pappenberger, 2009; Emerton
11 et al., 2016). An ensemble of meteorological inputs means that some of the uncertainty
12 in the NWP forecast can be accounted for in the hydrological model and ultimately in
13 the forecast flood maps. The forecast flood maps can be presented probabilistically as an
14 ensemble flood map forecast indicating the probability of flooding at a specific location
15 within a catchment. An ensemble NWP forecast lengthens the forecast lead-time where
16 the forecast is deemed skilful compared to a deterministic forecast (Emerton et al., 2016).
17 Information on flood inundation uncertainty is particularly useful for disaster management
18 teams operating before a flood event occurs and can be directly linked to flood impacts
19 such as maps of vulnerable infrastructure. Probabilistic flood maps can also be used to in-
20 form disaster risk reduction schemes such as Forecast-based Financing (FbF). FbF schemes
21 work by quantifying risks in advance of disasters, prepositioning funds, and agreeing in
22 advance how funds will be released based on forecasts, ahead of an event (OCHA, 2020).
23 FbF allows time for local action utilising the insurance funds for flood mitigation purposes.

24

25 In Chapters 3 to 6 we use forecast flood maps from JBA's Flood Foresight system.
26 Flood Foresight is a simulation library flood inundation forecasting system deployed in
27 several international countries for FbF applications.

1 2.2 Observing flooding from Space

2 Flooding events are usually observed using in situ ground-based gauging stations recording
3 river discharge or water level. The recorded flood discharge does not linearly correlate with
4 the observed inundation extent and is highly uncertain due to instrument error and rating
5 curve extrapolation uncertainties (Beven, 2016). Globally, there is limited coverage of
6 catchments with maintained gauging stations and the data is not always openly available.
7 Satellite-derived observations of flood extent have the potential to bring additional spatial
8 information into flood inundation forecasts compared to in situ point gauging stations.
9 Satellite-based Synthetic Aperture Radar (SAR) sensors are well known for their flood
10 detection capability (Grimaldi et al., 2016). Unobstructed flood waters appear dark on
11 SAR images due to the low backscatter return from the relatively smooth water surface.
12 SAR sensors also have an advantage over optical instruments as they can scan at night
13 and are not impacted by cloud and weather, usually associated with a flooding situation.
14 Optical instruments rely on solar energy and cannot penetrate cloud, making them less
15 useful during a flooding situation. Recent studies have investigated the flood detection
16 benefits from combining both optical and SAR imagery (Konapala et al., 2021; Tavus et
17 al., 2020).

18
19 Due to improvements in spatial resolution and more frequent revisit times, SAR data
20 have been used successfully to calibrate and validate hydrodynamic and hydraulic fore-
21 cast models (Schumann et al., 2009; Grimaldi et al., 2016). Dasgupta et al. (2018) detail
22 some of the challenges along with approaches to solutions of flood detection using SAR.
23 Examples of these challenges include: roughening of the water surface by heavy rain and
24 strong wind, emergent or partially submerged vegetation and flood detection in urban
25 areas. Accurate flood detection in urban areas particularly due to surface water flooding
26 has become increasingly important (Speight et al., 2021) and recent techniques have led

1 to improved flood detection (Mason et al., 2018, 2021a, 2021b).

2

3 The Copernicus Emergency Management Service (CEMS) (Copernicus Programme,
4 2021) offers freely available, open access Sentinel-1 SAR data. Currently (due to the
5 malfunction of Sentinel-1B in December, 2021) one satellite is in orbit, at 10 m ground
6 resolution and a six day revisit time (for the mid-latitudes). Sentinel-1C is due to launch
7 in April 2024 to replace Sentinel-1B. Nevertheless, Sentinel-1 data offers good coverage of
8 a potential flood event. For a major flood event CEMS can be triggered to offer additional
9 rapid flood mapping. Since late 2021, SAR-derived flood maps are produced for every
10 Sentinel-1 image detecting flooding around the world by the Global Flood Monitoring
11 (GFM) service (EU Science Hub, 2021; GFM, 2021; Hostache et al., 2021). Within eight
12 hours of the Sentinel-1 image acquisition, three flood detection algorithms are combined to
13 give the flood class (flooded or unflooded), the likelihood of a flood class representing the
14 uncertainty estimation per grid cell along with an exclusion mask where flooding cannot
15 be reliably detected from SAR. In Chapters 3, 4 and 5 we make use of SAR-derived flood
16 maps from Sentinel-1 satellite data for deterministic and ensemble forecast flood map
17 verification. In Chapter 6 we assimilate the GFM flood likelihood data to improve the
18 flood extent analysis from a simulation library flood forecasting system.

19 **2.3 Flood extent verification**

20 Verification is an essential part of model understanding and improvement, but has only
21 received limited attention over the past decade (Schumann, 2019). Forecast flood maps
22 can be verified against satellite-derived observations of flood extent. Validation of fore-
23 cast flood maps against remotely observed flood extent is typically carried out by labelling
24 each grid cell using a contingency table with categories: correctly predicted flooded, under-
25 prediction (miss), over-prediction (false alarm) and correctly predicted unflooded. Follow-

1 ing this categorisation, a variety of conventional binary performance measures such as the
2 Critical Success Index (CSI) can be calculated (see Chapter 3). Within a domain of inter-
3 est, a flood covering a significant area of the domain will be easier to accurately predict
4 (by chance) compared to a flood of a smaller extent. Thus, flood magnitude can create
5 a bias in the skill scores. It has been suggested by Stephens et al. (2014); Pappenberger
6 et al. (2007) that it is less important to validate all flooded cells, when only cells that are
7 close to the flood margin are difficult to predict. The flood edge location is an important
8 consideration for flood risk mitigation and response activities.

9
10 Flood maps at different spatial scales (grid lengths) will also impact conventional skill
11 scores. A high resolution, fine scale forecast flood map will show greater detail of the flood
12 extent and the flood-edge location compared to a low resolution, coarse scale flood map.
13 At a high resolution the discrepancy between the forecast and observed flood maps may be
14 closer in terms of distance, however a small mismatch will lead to a double penalty impact
15 on forecast verification. The model is penalised twice for the over-prediction (false alarm)
16 and the under-prediction (miss) (Stein & Stoop, 2019). When high resolution forecasts
17 are verified against observations at grid level, the predictability can appear to worsen and
18 the high resolution forecast would need to perform better than the low resolution forecast
19 to achieve the same verification score. Hence, it is not meaningful to compare verifi-
20 cation scores across different spatial scales. Conventional binary performance measures
21 (reviewed in Chapter 3) give a single, domain averaged, skill score. The averaged score
22 is less sensitive to changes in accuracy and does not indicate where location specific im-
23 provements could be made. The Fraction Skill Score (FSS; see Chapter 3 for details) uses
24 a neighbourhood approach to overcome the double penalty impact problem in convective
25 precipitation verification (Roberts & Lean, 2008). In Chapter 3 we apply the FSS scale-
26 selective verification approach to evaluate the whole flood, and the flood edge of forecast
27 flood maps and in Chapter 5 we use the same approach to compare three flood forecast-

1 ing systems, each producing forecast flood maps at different spatial scales. In Chapter 6
2 we use scale-selective verification methods to compare the analysis flood map, following
3 DA, against independent observations of flood extent from Sentinel-2 optical satellite data.

4

5 Ensemble flood maps require an extra dimension of verification, a measure of *spread*
6 as well as *skill*. A perfect ensemble should encompass forecast uncertainties such that
7 the ensemble spread is correlated to the RMSE of the forecast (Hopson, 2014). The
8 verification of ensemble forecasts usually involves comparing the RMSE of the ensemble
9 mean against an observed quantity to assess the skill of the forecast with the ensemble
10 standard deviation used as a measure of spread. To evaluate the accuracy of an ensemble
11 forecast, a number of verification measures have been proposed. Anderson et al. (2019)
12 developed a joint verification framework for end-to-end assessment of the England and
13 Wales Flood Forecasting Centre (FFC) ensemble flood forecasting system. Anderson et al.
14 (2019) describe verification metrics such as the continuous rank probability score (CRPS),
15 rank histograms, Brier Skill Score (BSS) and the relative operative characteristics (ROC)
16 diagrams that are commonly applied to assess the main ensemble attributes desirable in
17 both precipitation and streamflow ensemble forecasts (e.g., Renner et al., 2009). These
18 metrics refer to flooding events as part of a time series evaluated against a reference
19 benchmark, such as climatology, to produce an average skill score. In contrast, in Chapter
20 4 we consider ensemble *spatial* verification at a single time point. The spatial spread-skill
21 of the ensemble forecast is determined by evaluating the full ensemble against remote
22 observations of flooding. For a flood map ensemble to be considered spatially well-spread,
23 the spread or variation between ensemble members should equal the spatial predictability,
24 or skill of the ensemble members (Dey et al., 2014).

2.4 Introduction to data assimilation

Data assimilation (DA) typically finds an optimal *state* (such as river water level) and/or model parameter values (such as floodplain roughness coefficients) of a dynamical system, taking account of the previous forecast, the observations available, and both of their associated uncertainties. The updated state, the *analysis*, and/or parameter values are used to initiate the next forecast in a feedback loop or cycle (Figure 2.1).

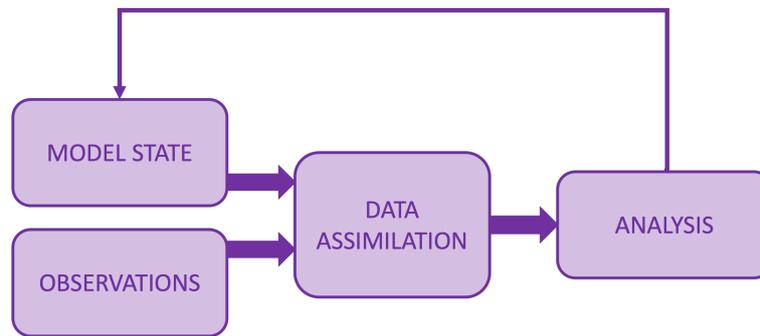


Figure 2.1: Data assimilation cycle

Bayesian estimation forms the basis of most data assimilation techniques and this assumes that errors in the forecast and observations can be represented by a Gaussian distribution. The uncertainties associated with the *prior* or *background*, usually the previous forecast, are represented by the covariance of the prior error distribution. Similarly, the uncertainties in the observations depend on the covariance of the observation *likelihood*. The aim of DA is to find the optimal state that maximises the *posterior* probability and thus minimises the variance of the posterior error distribution (Bouttier & Courtier, 2002). The optimal state is found by minimising a cost function (derived from the probability distributions), or in other words finding the state or parameter variables where the gradient of the cost function is equal to zero.

1 Variational DA methods use numerical minimisation methods such as iterative gra-
2 dient descent to minimise the cost function. Three-dimensional variational DA (3D-Var)
3 methods assimilate observations at a fixed point in time (Lorenz et al., 2000), whereas
4 4D-Var methods assimilate observations over a window of time (Bannister, 2017). Filter-
5 ing methods such as particle filtering (PF) and ensemble Kalman filter (EnKF) methods
6 sample (using Monte Carlo methods) and weight the prior distribution according to the
7 observation likelihood (van Leeuwen, 2009; Evensen, 1994).

8

9 Previously, SAR data have been used in several different ways to improve hydraulic
10 models and flood prediction through data assimilation (DA). A review of approaches used
11 to assimilate satellite-derived data into hydraulic models (from 2007 until 2015) can be
12 found in Table 7 of Grimaldi et al. (2016) and Table 1 of Revilla-Romero et al. (2016).
13 In Chapter 6 we review more recent DA approaches used to improve flood inundation
14 forecasts. A new DA framework is presented to improve the flood inundation analysis
15 from a simulation library forecasting system. The analysis is evaluated using scale selective
16 verification methods described in Chapter 3.

17 **2.5 Chapter summary**

18 In this chapter we have introduced and discussed a number of topics relevant to the thesis.
19 In Section 2.1 we described a simulation library flood forecasting system and its application
20 to disaster risk reduction. Forecast data from a simulation library system is used in all
21 of the case studies in this thesis. Satellite-derived observations of flooding outlined in
22 Section 2.2 are used to compare against flood forecasts. The spatial verification methods
23 introduced in Section 2.3 are calculated alongside new scale-selective methods in Chapters
24 3, 4, and 5. Data assimilation methods are introduced as background for Chapter 6 in
25 Section 2.4. In the next chapter we develop a scale-selective verification approach to

- 1 evaluate forecast flood extent maps.

1 Chapter 3

2 Spatial scale evaluation of forecast 3 flood inundation maps

4 In this chapter we address the first research question outlined in Chapter 1; What are
5 the skilful spatial scales in flood inundation forecasts made using a simulation library
6 approach? In particular we wish to find out:

- 7 • How can we determine the skilful spatial scales of forecast flood maps by comparing
8 against satellite-derived observations of flooding?
- 9 • Does the skilful spatial scale vary with location and how can this be visualised?
- 10 • How does validation of the flood edge location alone compare to validation of the
11 entire flood extent?
- 12 • How can the skilful spatial scale results be used in operational practice?

13 The remainder of this chapter (except for the chapter summary, Section 3.8), has been
14 published and is reproduced from (Hooker et al., 2022).

The Scale Puzzle



1 **3.1 Abstract**

2 Flood inundation forecast maps provide an essential tool to disaster management teams
3 for planning and preparation ahead of a flood event in order to mitigate the impacts of
4 flooding on the community. Evaluating the accuracy of forecast flood maps is essential
5 for model development and improving future flood predictions. Conventional, quantita-
6 tive binary verification measures typically provide a domain-averaged score, at grid level,
7 of forecast skill. This score is dependent on the magnitude of the flood and the spatial
8 scale of the flood map. Binary scores have limited physical meaning and do not indicate
9 location-specific variations in forecast skill that enable targeted model improvements to
10 be made. A new, scale-selective approach is presented here to evaluate forecast flood inun-
11 dation maps against remotely observed flood extents. A neighbourhood approach based
12 on the Fraction Skill Score is applied to assess the spatial scale at which the forecast be-
13 comes skilful at capturing the observed flood. This skilful scale varies with location and
14 when combined with a contingency map creates a novel categorical scale map, a valuable
15 visual tool for model evaluation and development. The impact of model improvements
16 on forecast flood map accuracy skill scores are often masked by large areas of correctly

1 predicted flooded/unflooded cells. To address this, the accuracy of the flood-edge location
2 is evaluated. The flood-edge location accuracy proves to be more sensitive to variations
3 in forecast skill and spatial scale compared to the accuracy of the entire flood extent.
4 Additionally, the resulting skilful scale of the flood-edge provides a physically meaningful
5 verification measure of the forecast flood-edge discrepancy. The methods are illustrated
6 by application to a case study flood event (with an estimated return period of 120 to 550
7 years) of the River Wye and River Lugg (UK) in February 2020.

8

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

18

19 Highlights

- 20 • A novel spatial scale-selective approach to evaluate forecast flood maps against Syn-
21 thetic Aperture Radar data.
- 22 • Validation of the flood edge gives a physically meaningful measure of prediction
23 accuracy.
- 24 • Conventional contingency flood maps are improved by including a location-specific
25 skilful spatial scale.

3.2 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 (Trepekli 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 pe-

1 nalisied twice for the over-prediction (false alarm) and the under-prediction (miss) (Stein
2 & Stoop, 2019). When high resolution forecasts are verified against observations at grid
3 level, the predictability can appear to worsen and the high resolution forecast would need
4 to perform better than the low resolution forecast to achieve the same verification score.
5 It is not meaningful to compare verification scores across different spatial scales. Spatial
6 verification methods for flood inundation mapping have only received limited attention
7 over the past decade (Schumann, 2019).

8

9 Verification approaches that account for uncertainties in observations and small dis-
10 crepancies in gridded data using a fuzzy set approach (Hagen, 2003) have previously been
11 applied to flood mapping (Pappenberger et al., 2007; Dasgupta et al., 2018). However, the
12 fuzzy set method does not incorporate variations in spatial scale (Cloke & Pappenberger,
13 2008). In atmospheric sciences, verification approaches that account for changes in spatial
14 scale are well established. These approaches include the Fraction Skill Score (FSS), which
15 applies a neighbourhood approach to assess a useful/skilful scale (Roberts & Lean, 2008)
16 of a precipitation forecast. Dey et al. (2014); Dey, Roberts, et al. (2016) developed the FSS
17 approach to produce location-specific agreement scales between the forecast and observed
18 fields to understand the spatial predictability of an ensemble forecast. Other spatial scale
19 approaches include the wavelet method of scale decomposition, where the forecast and
20 observed fields are decomposed into maps at different scales by wavelet transformation
21 and subsequently verified (Briggs & Levine, 1997; Casati & Wilson, 2007). Cloke and
22 Pappenberger (2008) note that this method is extremely sensitive to offsetting of maps.

23

24 In general, the performance of forecast flood maps are evaluated for the entire flood ex-
25 tent, regardless of flood magnitude, adding bias to binary performance measures (Stephens
26 et al., 2014). Stephens et al. (2014) question whether it is important to validate all
27 flooded cells, when only cells that are close to the flood margin are difficult to predict.

1 Pappenberger et al. (2007) evaluated model performance only on cells that were subject
2 to change between differing model runs to address the issue of large areas of correctly
3 predicted flooded/unflooded cells masking variations in forecast skill scores.

4

5 Satellite based Synthetic Aperture Radar (SAR) sensors are well known for their flood
6 detection capability. Unobstructed flood waters appear dark on SAR images due to the
7 low backscatter return from the relatively smooth water surface. SAR sensors also have
8 an advantage over optical instruments as they can scan at night and are not impacted by
9 cloud and weather, usually associated with a flooding situation. Due to improvements in
10 spatial resolution and more frequent revisit times, SAR data has been used successfully
11 to calibrate and validate hydrodynamic and hydraulic forecast models (Schumann et al.,
12 2009; Grimaldi et al., 2016). Further model improvements have been shown through the
13 assimilation of SAR data (e.g., García-Pintado et al., 2015; Hostache et al., 2018; Cooper
14 et al., 2019; Di Mauro et al., 2020; Dasgupta et al., 2018, 2021a, 2021b). Recent tech-
15 niques have improved the flood detection in urban areas using medium and high resolution
16 SAR (Mason et al., 2018, 2021a, 2021b). The Copernicus Emergency Management Service
17 (CEMS) (Copernicus Programme, 2021) offers freely available, open access Sentinel-1 SAR
18 data. Currently (due to the malfunction of Sentinel-1B in December, 2021) one satellite
19 is in orbit, at 10 m ground resolution and a six day revisit time (for the mid-latitudes).
20 Nevertheless, Sentinel-1 data offers good coverage of a potential flood event. For a major
21 flood event CEMS can be triggered to offer additional rapid flood mapping. From 2022,
22 the new Global Flood Monitoring (GFM) product (GFM, 2021; Hostache et al., 2021) of
23 the Copernicus Emergency Management Service (CEMS) (Copernicus Programme, 2021)
24 produces Sentinel-1 SAR-derived flood inundation maps using three flood detection al-
25 gorithms providing uncertainty and population affected estimates within 8 hours of the
26 image acquisition.

27

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

9
10 The objective of this paper is to present a scale-selective approach to evaluate flood
11 inundation forecast maps and to develop a physically meaningful measure of flood-edge
12 location accuracy that can be automated and easily applied in practice. The method has
13 been developed with operational forecast verification in mind, but it is applicable to all
14 flood inundation maps. A new approach is described and applied here to evaluate the
15 spatial scale at which the forecast becomes useful/skilful at capturing the remotely ob-
16 served flood extent and specifically the flood-edge location. The spatial skill of a forecast
17 flood map varies with location. We aim to improve the conventional contingency map by
18 incorporating the skilful scale to create a new *categorical scale map*. Also, we address how
19 representation errors arising from observation spatial scale variations and interpolation
20 have an impact on model evaluation.

21
22 In the rest of this paper we explore the features of a novel scale-selective evaluation
23 approach illustrated through application to a case study. In Section 3.3 we describe the
24 case study, a recent flooding event in the UK following Storm Dennis, February 2020, along
25 with catchment descriptions for three chosen domains. The flood inundation forecasting
26 system developed by JBA Consulting, Flood Foresight, (Revilla-Romero et al., 2017) is
27 used to produce forecast flood maps for the event and is detailed in Section 3.4.1. Section

1 3.4.2 explains two methods that are used to derive remotely observed flood maps from SAR
2 imagery. Our new approach to the spatial evaluation of flood maps is detailed in Section
3 3.5 along with descriptions of other binary performance measures. The novel categorical
4 scale map is applied to the case studies in Section 3.6, and the evaluation results are
5 discussed. We conclude in Section 3.7 and discuss the wider applications of a spatial scale
6 approach to flood map skill evaluation.

7 **3.3 Flood event**

8 This extreme flooding event is chosen here as a case study to demonstrate the features
9 of a spatial scale approach to forecast flood map evaluation. During February 2020,
10 three named Storms, Ciara, Dennis and Jorge, arrived in quick succession delivered by
11 a powerful and ideally positioned jet-stream that enabled rapid cyclogenesis (Davies et
12 al., 2021). Each storm rapidly intensified and deepened bringing damaging winds and
13 exceptionally heavy rainfall across the UK (Met Office, 2020). This led to the River Wye
14 reaching its highest ever recorded water level at the Old Bridge in Hereford (riverlevels.uk,
15 2020). The annual exceedance probability (AEP) for the recorded peak flow of the Lugg
16 and Wye rivers was 0.2 - 0.8 % (return period 120-550 years) and 0.6 - 2.0 % (160-550
17 years) respectively (Sefton et al., 2021).

18 **3.3.1 February 2020**

19 February 2020 was the UK's wettest February on record and the fifth wettest month
20 ever recorded. The UK average rainfall total exceeded the 1981 – 2010 average by 237%
21 (Kendon, 2020). Locally, in northwest England and north Wales the rainfall exceedance
22 was three to four times the typical monthly average rainfall. During this period around
23 4000 to 5000 properties were flooded in the UK, with significant river water levels recorded
24 in Wales, west and northwest England (Sefton et al., 2021). With six days between Ciara

1 and Dennis, groundwater and river levels were high and soils saturated. The Environment
2 Agency issued a record number of over 600 flood alerts and warnings for England (JBA,
3 2021).

4 3.3.2 Catchment location and description

5 Three domains, each differing in hydrological characteristics, have been selected for fore-
6 cast flood map evaluation during the storm Dennis flooding event. Two domains (A and
7 B) have been chosen from the Wye catchment (Fig. 3.1), a 28.4 km length centred upon
8 Ross-on-Wye (A) and the Wye at Hereford (B), a 5.8 km section. A third domain (C)
9 includes 4 km of the River Lugg.

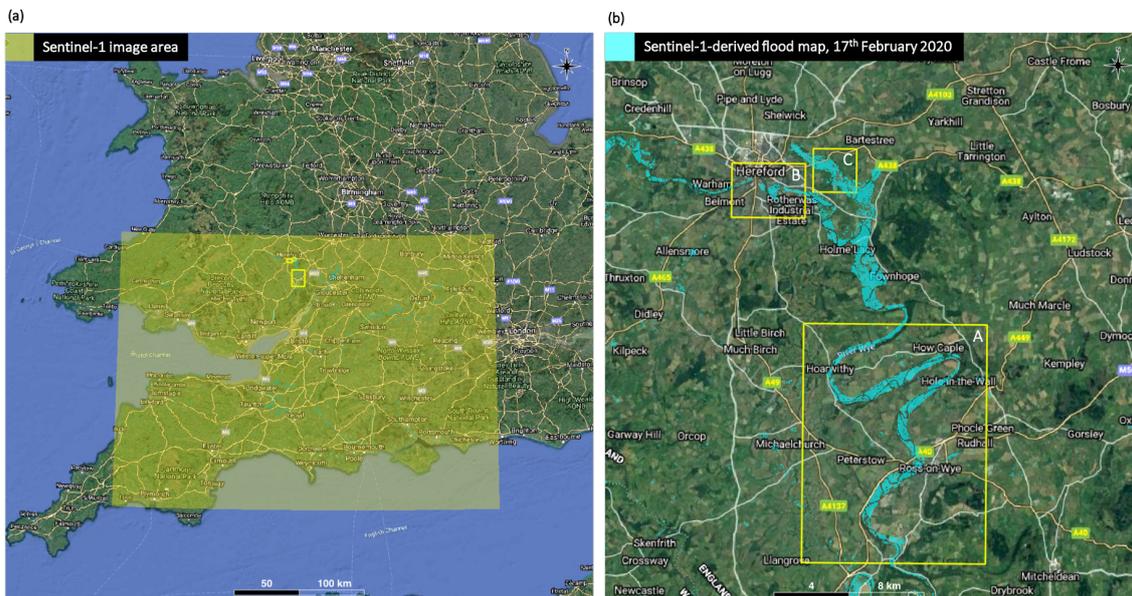


Figure 3.1: Location of Sentinel-1 image acquisition over southeast UK (a) and flood map evaluation domains (b). Domain A: 28.4 km length of the River Wye centred at Ross-on-Wye, domain size 9.8 x 12.8 km. Domain B: 5.8 km of the River Wye at Hereford, domain size 3.0 x 4.0 km. Domain C: 4 km of the River Lugg at Lugwardine, domain size 2.3 x 2.3 km. Base map from Google Maps.

1 **3.3.2.1 The River Wye (domains A and B)**

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

14 **3.3.2.2 River Lugg at Lugwardine (domain C)**

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

22 **3.3.2.3 Event hydrology**

23 The observed catchment rainfall (which also includes a downstream section of the River
24 Wye) shows that 50 mm fell on the 15th, 10 mm on the 16th and 1 mm on the 17th February

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

8

9 Daily maximum river levels recorded at Ross-on-Wye, the Old Bridge, Hereford and
10 Lugwardine for January to March 2020 are plotted in Figure 3.2 (riverlevels.uk, 2020). The
11 impact of the three storms on the River Wye is indicated by a very sharp rise in water
12 levels from the 8th to the 10th February following storm Ciara. Further heavy showers
13 maintained high water levels before storm Dennis brought an exceptional rise in water
14 levels, peaking on the morning of the 17th February with record levels recorded at Hereford
15 (6.11 m at 9.30 am UCT) and Ross-on-Wye (4.77 m at 5.45 am UTC). Unfortunately there
16 are two days of missing data at Ross-on-Wye following the flood event. By analysing the
17 trend between the Hereford and Ross-on-Wye river levels, the peak level at Ross-on-Wye
18 was likely higher and later than recorded. The response of the Wye at Hereford is faster
19 than at Ross-on-Wye, most likely due to the upstream location of Hereford and a more
20 constrained embankment with the city center located either side of the river. In comparison
21 to the fast, rapid response of the Wye, the River Lugg displays a distinctively dampened
22 response. Whilst the Lugg initially responded quickly to the heavy rainfall, once bankfull
23 was reached and overtopping occurred the water levels remained consistently high, with
24 floodwaters extending across the relatively flat flood plain.

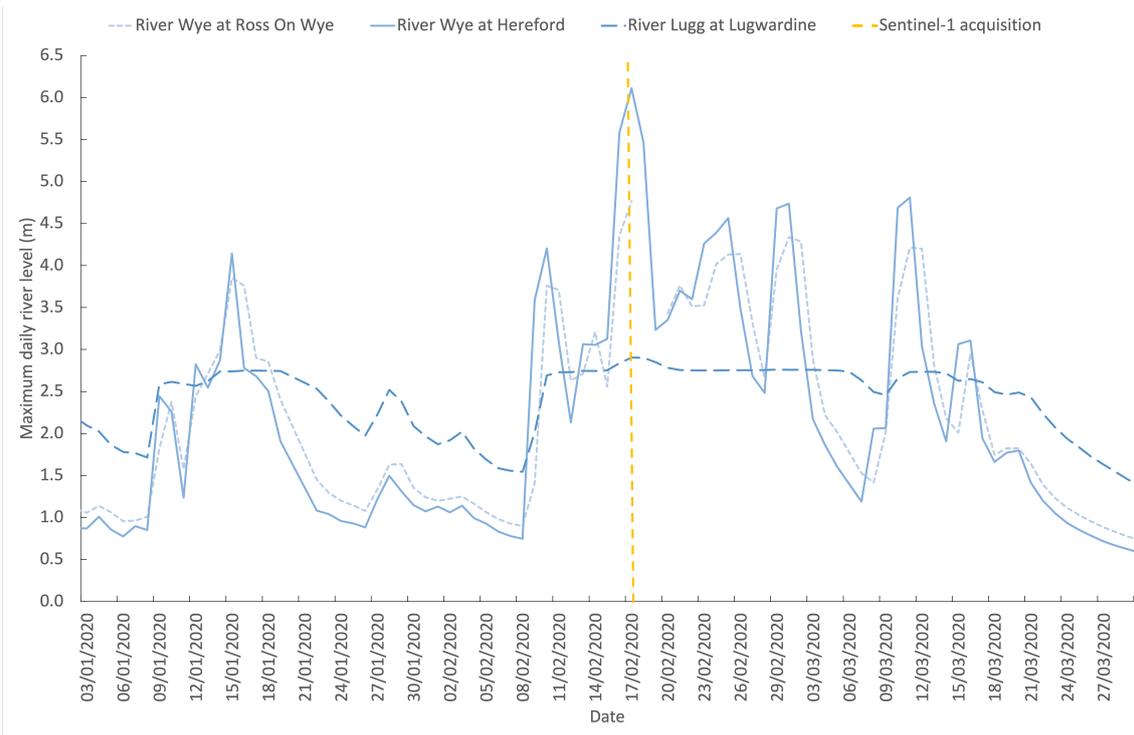


Figure 3.2: Daily maximum river levels (m) at Ross-on-Wye, Hereford and Lugwardine. The dashed yellow line indicates Sentinel-1 SAR acquisition date.

1 3.4 Data

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

4 3.4.1 Flood Foresight

5 Flood Foresight (Fig. 3.3), developed and run routinely by JBA Consulting, is a flu-
 6 vial flood inundation mapping system that can be implemented in any catchment around
 7 the globe. Flood Foresight utilises a simulation library approach to generate maps of real
 8 time and forecast flood inundation and water depth. The simulation library approach saves
 9 valuable computing time and allows the application of Flood Foresight in near continuous
 10 real-time at national and international scales. A library of flood maps is pre-computed

1 using JFlow[®], a 2D hydrodynamic model (Bradbrook, 2006). Note that in this study the
2 flood maps are undefended i.e. temporary flood defences are not included. JFlow uses a
3 raster-based approach with a detailed underlying DTM and a simplified form of the full 2D
4 hydrodynamic equations that capture the main controls of the flood routing for shallow,
5 topographically driven flow. Five flood maps at 5 m resolution are created for 20, 75, 100,
6 200 and 1000 year return period flood events (corresponding to annual exceedance prob-
7 abilities (AEPs) of 5%, 1.3%, 1%, 0.5% and 0.1% respectively). These are interpolated
8 to derive five intermediate maps between each adjacent pair of the JFlow maps, equally
9 spaced in return period creating a total library of thirty flood maps. Flood Foresight
10 takes inputs of rainfall from numerical weather prediction (NWP) models, river gauge
11 data (both historical and real-time) and forecast streamflow and uses these to select the
12 most appropriate flood map for the location and forecast time period. The UK and Ireland
13 configurations of the Flood Forecasting Module use deterministic streamflow forecast data
14 from the Swedish Meteorological and Hydrological Institute (SMHI) European HYdrolog-
15 ical Predictions for the Environment (E-HYPE). The meteorological input data for the
16 E-HYPE model is the European Centre for Medium-range Weather Forecasts (ECMWF)
17 Atmospheric Model high resolution (HRES) numerical weather prediction (NWP) model
18 on a $0.1^\circ \times 0.1^\circ$ grid with forecasts issued daily out to 10 days lead time. Forecast flood
19 maps for the UK are produced on a 25 m grid length out to 10 days ahead (see Mason et
20 al. (2021b) Section 2.1 for additional details).

21 **3.4.2 SAR-derived flood maps**

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

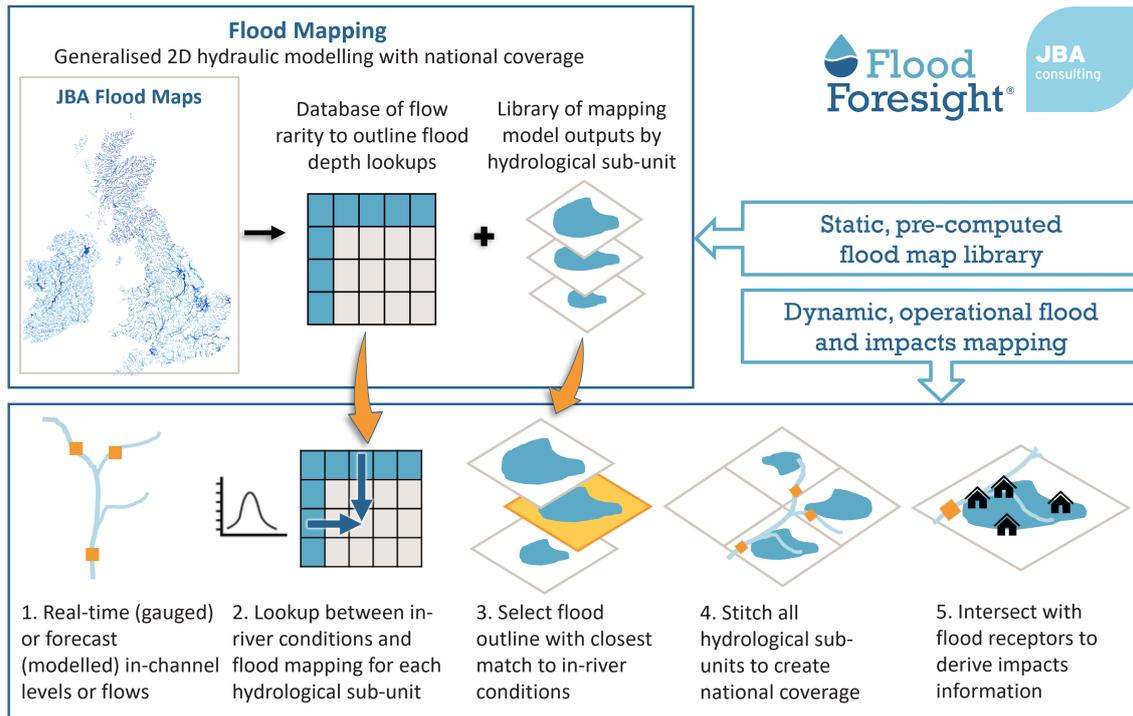


Figure 3.3: Flood Foresight flood map simulation library selection process. Source JBA Consulting.

1 track was used to derive the flood map in both methods.

2

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

1 of the three domains out to 10 days lead time.

2

3 In the second method, the same Sentinel-1 SAR image (in this case using both VV and
4 VH) was processed using Google Earth Engine (GEE) to derive flood maps at a range of
5 spatial resolutions (5 m to 25 m). GEE holds a catalogue of level-1 preprocessed Sentinel-1
6 SAR images (Google Earth Engine Catalog, 2021). A smoothing filter is applied to reduce
7 speckle and a pre and post flood image are used to train a Classification And Regres-
8 sion Tree (CART) classifier (Breiman et al., 1984; Google Earth Engine CART, 2021).
9 The classifier is applied to the whole image to produce a flood map at a specified scale.
10 GEE uses an image pyramid approach to scale, or pixel resolution, analysis. This means
11 variations in the scale selected are determined from the scale of the input image (Google
12 Earth Engine Scale, 2021). The variation of the flood extent detected at a range of spatial
13 resolutions and the impact of re-scaling and interpolation errors on performance measures
14 are investigated.

15

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

20 **3.5 Flood map evaluation methods**

21 The following subsections detail a new spatial scale-selective approach to forecast flood
22 map evaluation. The Fraction Skill Score (FSS) developed by Roberts and Lean (2008) for
23 validation of convective precipitation forecasts in atmospheric science uses a neighbour-
24 hood approach to determine the scale at which the forecast becomes skilful. Dey, Roberts,
25 et al. (2016) developed this approach to determine an agreement scale between an ensem-

1 ble forecast and observations at each grid cell to add location-specific information. Here
 2 we extend the technique to apply it to the new application of flood inundation mapping,
 3 and further develop a novel categorical scale map that combines an agreement scale map
 4 with a conventional contingency map.

5 3.5.1 Spatial scale-selective approach

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

13

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

$$MSE_n = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [O_{nij} - M_{nij}]^2. \quad (3.1)$$

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

is calculated

$$MSE_{n(ref)} = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [O_{nij}^2 + M_{nij}^2]. \quad (3.2)$$

The FSS is given by

$$FSS_n = 1 - \frac{MSE_n}{MSE_{n(ref)}}. \quad (3.3)$$

- 1 Figure 3.4 illustrates an example of the FSS application at grid level ($n = 1$) and at the
- 2 next neighbourhood size $n = 3$. In this simple example, there is no agreement between
- 3 the model and observation at grid level but at $n = 3$, the skill score improves to 0.92.

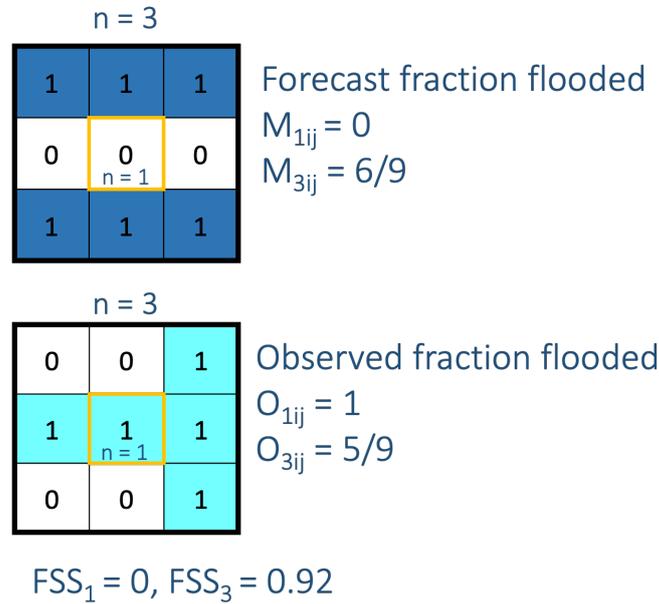


Figure 3.4: FSS (see subsection 3.5.1 for calculation details) example applied to a binary flooded (1) / unflooded (0) field at grid scale (yellow box, $n = 1$) and a 3 x 3 neighbourhood (black box, $n = 3$). The observed SAR-derived forecast is in turquoise and the forecast is shown in blue.

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 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

$$FSS_T \geq 0.5 + \frac{f_o}{2}, \quad (3.4)$$

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

6

7 When the FSS is plotted against spatial scale (neighbourhood size), we can identify a
 8 spatial scale when the FSS first equals or exceeds FSS_T (Fig. 3.6 shows an example of
 9 this plot). The spatial scale (neighbourhood size) reached at FSS_T can tell us the dis-
 10 placement distance (D_T) between the observed and forecast flood, or more meaningfully
 11 the flood-edge locations. As the flood-edge represents a very small fraction of the domain,
 12 the scale at FSS_T will tend to $2D_T$, meaning the displacement distance is half of this
 13 scale (see Figure 4 in Roberts and Lean (2008)).

14

15 It has been shown by Skok and Roberts (2016) that care must be taken when calculating
 16 the FSS near to the domain boundary since increasingly larger neighbourhood sizes would

1 extend further beyond the boundary edge. Skok and Roberts (2016) concluded that as
 2 long as the domain was sufficiently large, relative to the spatial errors, then the boundary
 3 effect could be considered to be insignificant. For flood mapping verification purposes the
 4 domain area should be selected to include the area of interest (e.g. the floodplain) with the
 5 neighbourhoods considered extending beyond the domain at the boundary. This assumes
 6 that the observations available allow this. If this is not that case then another boundary
 7 method could be applied, such as cropping at the domain edge.

8 3.5.2 Location dependent agreement scales

9 The FSS gives an overall domain-averaged measure of forecast performance and an average
 10 minimum scale at which the forecast is deemed skilful. Dey, Roberts, et al. (2016) describe
 11 a method for calculating an agreement scale at each grid cell located at coordinate position
 12 (i, j) . A brief summary of the method is presented here. Two fields are considered f_{1ij}
 13 and f_{2ij} . In this application these are the forecast and observed fields. In alternative
 14 applications the method could be applied to measure similarity between members of an
 15 ensemble. The fields in this instance are not required to be thresholded and can be applied
 16 to flood depths. The aim is to find a minimum neighbourhood size (or scale) for every
 17 grid point such that there is an agreement between f_{1ij} and f_{2ij} . This is known as the
 18 agreement scale S_{ij} . The relationship between the agreement scale and the neighbourhood
 19 size described in Section 3.5.1 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, $S = 0$ ($n = 1$),

$$D_{ij}^S = \frac{(f_{1ij}^S - f_{2ij}^S)^2}{(f_{1ij}^S)^2 + (f_{2ij}^S)^2}. \quad (3.5)$$

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:

$$D_{ij}^S \leq D_{crit,ij}^{S_{ij}} \quad \text{where} \quad D_{crit,ij}^S = \alpha + (1 - \alpha) \frac{S}{S_{lim}} \quad (3.6)$$

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

7 3.5.3 Categorical scale map

8 Currently, the agreement scale map proposed by Dey, Roberts, et al. (2016) provides a
 9 location-specific scale of agreement between the forecast and observed flood map. However,
 10 it does not show whether the model is over- or under-predicting the flood extent. In our
 11 work, we develop the agreement scale map further by combining with a contingency map
 12 for the forecast to create a new *categorical scale map*. This highlights the agreement scale
 13 for areas of over- or under-prediction. In a contingency map, each cell in the forecast and
 14 observed flood map are compared and classified using a contingency table (Table 3.1).
 15 The categories are re-classified numerically in the array for automated updating of the
 16 agreement scale map. Over-predicted cells (B) are set to -1, under-predicted cells (C) are
 17 set to +1, correctly predicted flooded cells (A) are assigned NaN and correctly predicted
 18 unflooded cells are set to 0. The array element-wise product of the agreement scale map
 19 and the numerical contingency map produces the new categorical scale map.

Table 3.1: Contingency table

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

1 **3.5.4 Binary performance measures**

2 It has been suggested by Cloke and Pappenberger (2008) that a range of performance
 3 measures should be applied so that a forecast can be assessed as rigorously as possible. A
 4 selection of commonly applied binary performance measures, each focusing on a different
 5 aspect of performance have been included here for comparison with the Fraction Skill
 6 Score results. Following the application of a contingency table (Table 3.1) to the forecast
 7 flood map, a number of binary performance measures can be calculated (Table 3.2). Table
 8 3.2 describes the range of performance value, the ideal score and a description of which
 9 aspects of the forecast flood map performance each binary measure assesses.

Table 3.2: Binary performance measures and formula based on contingency Table 3.1.

| 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 |

1 **3.6 Results**

2 We illustrate and discuss our new method applied to the flood event in subsection 3.6.1
3 and 3.6.2. The scale-selective approach is applied to an extreme flooding event in the
4 UK to determine a useful/skilful spatial scale for both the entire flood extent and the
5 flood-edge location for three domains out to 10-days lead time. An example forecast flood
6 map for 0-day lead time compared with the SAR-derived flood map is presented as a
7 contingency map in Figure 3.5. The zoomed in perspective shows the double penalty
8 impact described in Section 3.2. The discrepancy at the flood-edge depends on the spatial
9 scale of the forecast flood maps along with the model performance. Next, in subsection
10 3.6.3 location-specific agreement scales are presented on categorical scale maps. The final
11 subsection 3.6.4 addresses the question of the impact of representation error caused by
12 variations in SAR-derived flood map spatial resolution on the evaluation results.

13 **3.6.1 Spatial scale variability of forecast flood extent and flood-edge** 14 **location**

15 An evaluation of the spatial skill of the Flood Foresight forecast flood maps against the
16 SAR-derived flood map for the flood peak on the 17th February 2020 has been calculated
17 for each domain (Fig. 3.1) for both the entire flood extent and the flood-edge location. The
18 Fraction Skill Score (FSS) is applied to increasing neighbourhood sizes (n) to determine
19 the spatial scale at which the forecast becomes skilful at capturing the observed flood.
20 Figure 3.6 shows FSS against n for one example, the River Lugg (domain C) for the entire
21 flood (a) and the flood-edge (b). Each line represents a different model run date from the
22 10/02/2020 (7-day lead time) to the 17/02/2020 (0-day lead time). With the exception
23 of the 7-day lead time, all forecasts for the whole flood (Fig. 3.6a) exceed the FSS_T at
24 grid level ($n = 1$) with gradually improving skill as n increases. In contrast to this, the
25 FSS applied to the flood-edge (Fig. 3.6b) shows all forecasts below FSS_T at grid level

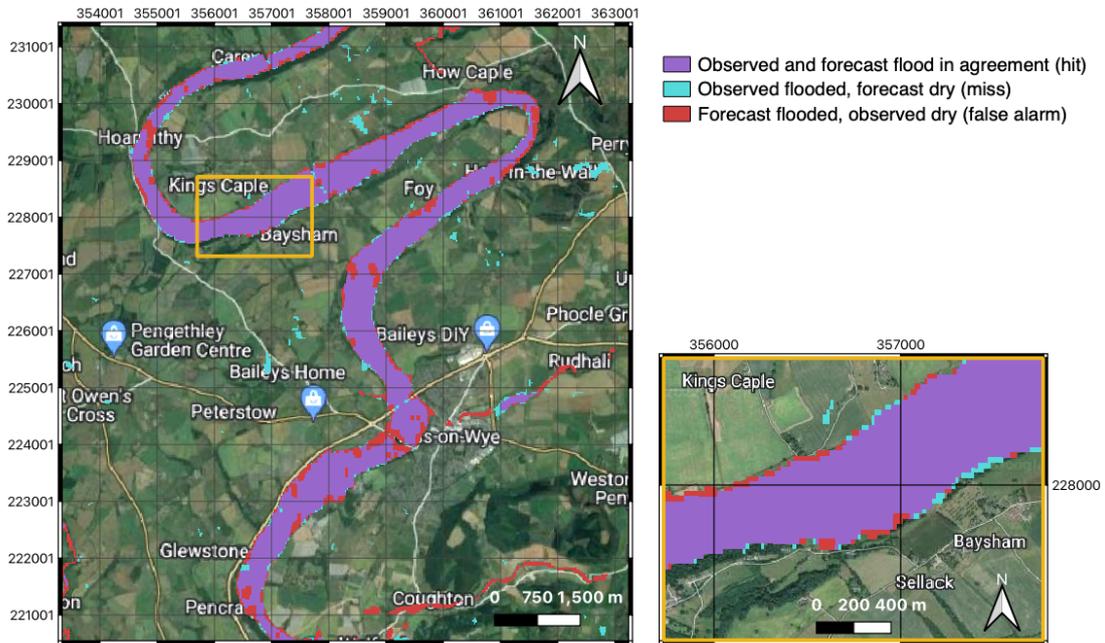


Figure 3.5: Left panel: contingency map of a 0-day lead time forecast versus the HASARD SAR-derived flood map for the Wye valley indicates the model is predicting the flood extent accurately, including the position of the flood-edge. Right panel: Zoom of yellow box on the left panel. On closer inspection, at grid level, the flood-edge in many places is over- or under-predicted by around one grid length. Base map from Google Maps.

1 and $n = 3$ with the skill increasing more rapidly compared with the whole flood to reach
 2 FSS_T at $n = 5$ for all run dates within a 5-day lead time (except for 16/02/2020, which
 3 is just below FSS_T). This indicates that the flood-edge is forecast to be around 62.5 m
 4 from the observed flood-edge, on average, for a 5-day lead time. The difference between
 5 the gradients of the plots indicate the flood-edge is more sensitive to changes in spatial
 6 scale compared with evaluation of the whole flooded area. The whole flood verification
 7 here indicates a strong model performance. However, verifying the whole flood alone could
 8 mask the flood-edge location performance, which in this case has a coarser scale at FSS_T .
 9 Similar trends in FSS with neighbourhood size and comparisons between the entire flood

1 and the flood-edge verification scales are found for all domains. The rate of FSS increase,
 2 or FSS gradient with n , tells us how quickly the forecast skill improves with increasing
 3 scale. A more spatially accurate forecast of the flood-edge will demonstrate a steeper
 4 gradient, reaching FSS_T at a smaller neighbourhood size.

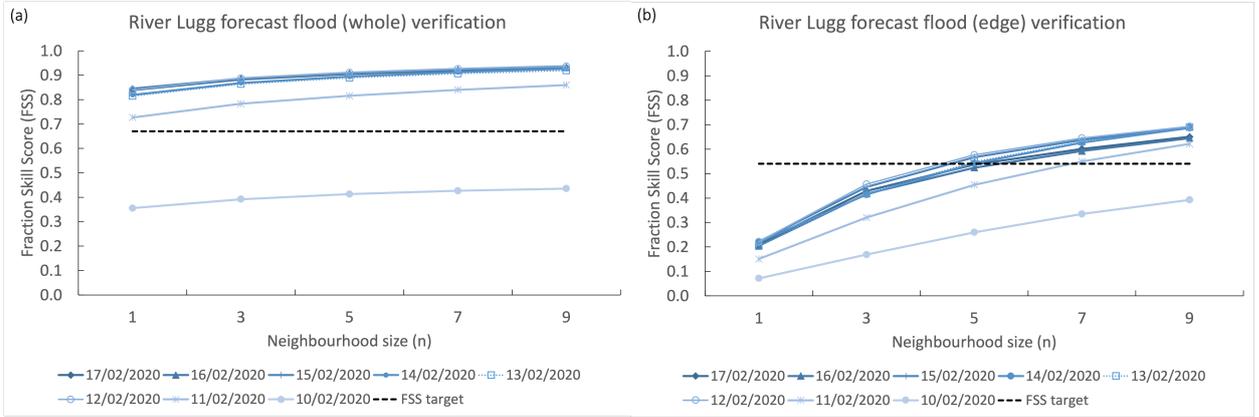


Figure 3.6: FSS calculated for the River Lugg at Lugwardine for (a) entire flood extent and (b) the flood-edge for increasing neighbourhood sizes for daily forecast lead times up to 7 days.

5 3.6.2 Comparison of spatial scales at differing lead times and domain 6 location

7 The performance measures for each domain for daily lead times out to 10 days are pre-
 8 sented in Figure 3.7. The FSS at $n = 1, 3$, and 5 are shown along with Critical Success
 9 Index (CSI), Hit Rate (HR), Pierce Skill Score (PSS) and the Bias (see Table 3.2 for def-
 10 initions). The Bias score is an indicator of over- or under-prediction of the flood extent
 11 and is plotted on a separate axis to account for the larger range. For lead times within
 12 5-days of the flood peak, $FSS > 0.8$ for the entire flooded area at grid level for the River
 13 Wye (domain A) indicates a strong model performance (Fig. 3.7a). There is a dip in
 14 the FSS on the 16/02/2020 where the forecast over-predicts the flood extent. This is
 15 also reflected in the CSI score. In contrast to this the HR and PSS increase, despite the
 16 over-prediction, as more observed flood cells are correctly predicted wet. We note that

1 the PSS (HR - FAR) does account for over-prediction, however the FAR is the fraction of
2 the dry area incorrectly predicted wet, which is very small relative to the HR (0.03 versus
3 0.90). Validation of the River Wye flood-edge (Fig. 3.7b) is more sensitive to changes in
4 neighbourhood size compared with the whole flood validation. Here the flood-edge is very
5 well forecast in terms of spatial location and exceeds FSS_T at $n = 3$ (on average, 37.5 m
6 displacement) for a 5-day lead time (except for 1-day lead time where FSS_T is exceeded
7 at $n = 5$). As shown previously in Subsection 3.6.1, the forecast of the River Lugg flood-
8 edge is skilful at $n = 5$ (Fig. 3.7f) (on average, 62.5 m displacement) for a 5-day lead
9 time. Differences in the hydrological characteristics might explain differences in model
10 performance. The Wye valley flood plain is well defined with distinctive valley sides and
11 this event proved to be valley filling in contrast to the Lugg flood plain which is relatively
12 flat and extensive. This could explain the increased skill shown for the prediction of the
13 Wye flood-edge. The average observed flood top width for the Lugg (domain C) is 740 m
14 and for the Wye (domain A) 430 m. This gives a flood-edge displacement as a fraction of
15 the flood top width of 7.4% for the Lugg and 7.8% for the Wye.

16

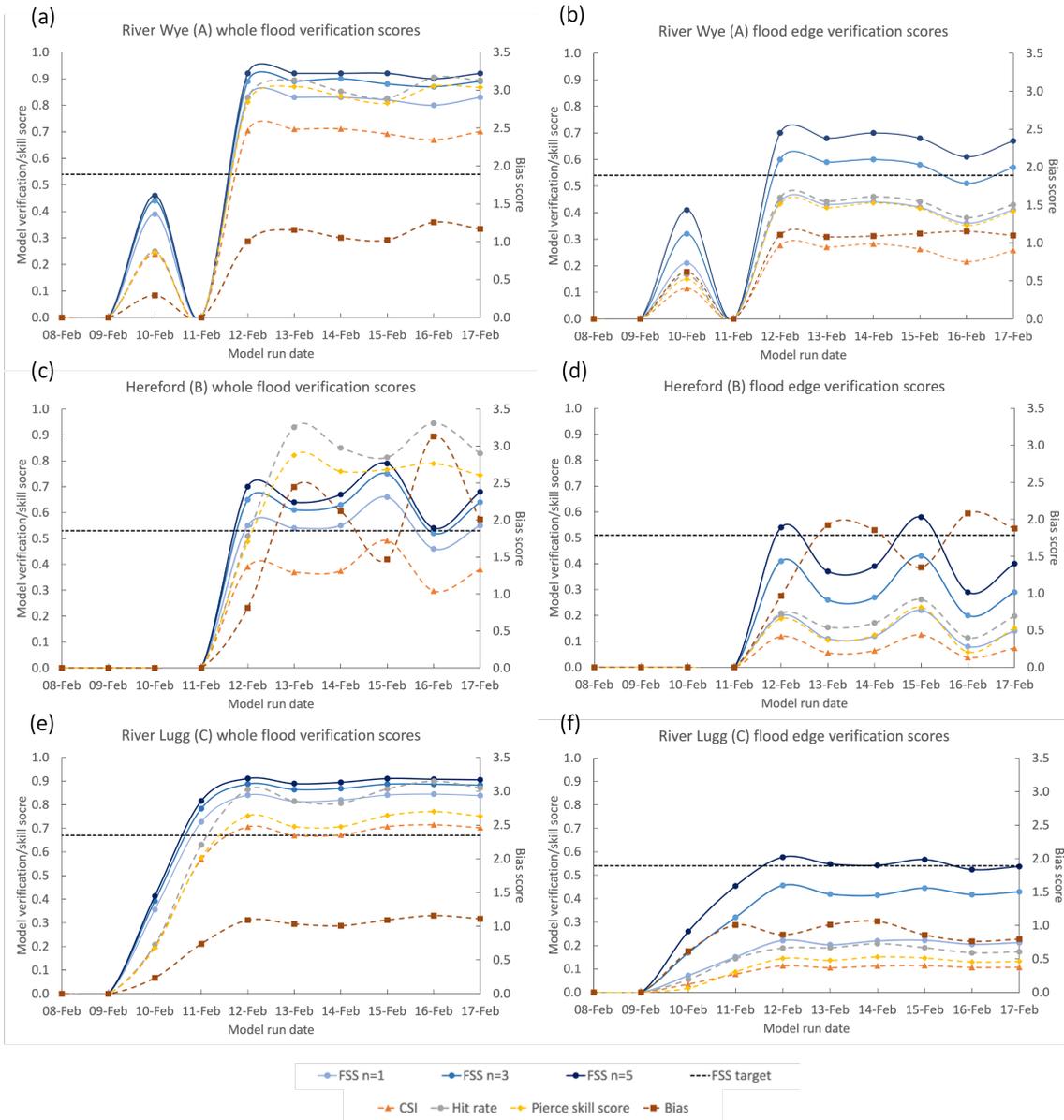


Figure 3.7: Conventional binary performance measures (dashed lines) and FSS (solid lines) at $n = 1, 3,$ and 5 for each domain for both the whole flooded area and the flood-edge for daily lead times out to 10 days for the River Wye (domain A, (a) and (b), Hereford (domain B, (c) and (d)) and the River 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.

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

17
18 There is more variation in skilful scale with lead time evident for the Wye at Hereford
19 (domain B) in Figure 3.7c and d compared with domain A and C. To achieve the same
20 FSS for the whole flood as domain A and C up to a 5-day lead time, the neighbourhood
21 size would need to exceed $n = 5$. The model is over-predicting the flood extent, in par-
22 ticular on the 16/02/2020 (1-day) lead time. This overprediction at 1-day lead time is
23 evident for all domains as can be seen in the Bias scores but the impact of this is most
24 noticeable at Hereford. Hereford has more complex topography compared to the other
25 domains, particularly along the river bank with bridges, buildings, permanent and tempo-
26 rary flood defences deployed during the event affecting the flow of the flood wave through
27 the city. The maps used in the simulation library of Flood Foresight are produced using a

1 bare-earth DTM. Despite this, the model performs well, exceeding FSS_T at $n = 5$ at the
2 5-day and 2-day lead times for the flood-edge forecast.

3

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

10 **3.6.3 Categorical scale maps**

11 Location dependent categorical scale maps (Subsection 3.5.3) have been calculated for all
12 run dates for both the entire flooded area and the flood-edge. Figure 3.8 shows categor-
13 ical scale maps for the whole flood for three different lead times for each domain, longer
14 lead times are on the left. The run dates vary with domain to present the most informa-
15 tive maps such that variation in forecast skill can be seen across the different lead times.
16 The colours on the map indicate grid cell specific agreement scales (Subsection 3.5.2)
17 between the forecast flood map and the SAR-derived flood map. Grey/white regions
18 indicate correctly predicted flooded/unflooded cells, red shows the forecast flood extent
19 is under-predicted (miss) and blue indicates over-prediction (false alarm). Increasingly
20 darker shades of red/blue show that larger scales were needed for the agreement criteria
21 to be met. The darkest blue at $S = 10$ indicates a total mismatch between forecast and
22 observed flooding. The addition of the agreement scale information in comparison to a
23 conventional contingency map (for an example, see Fig. 3.5) quickly highlights regions
24 of total mismatch through the darkest shading, with areas that are slightly misaligned
25 in lighter shades. The agreement scale indicated gives a physical measure of distance at
26 specific locations between the forecast and the observed flood map (where $S < S_{lim}$).

1

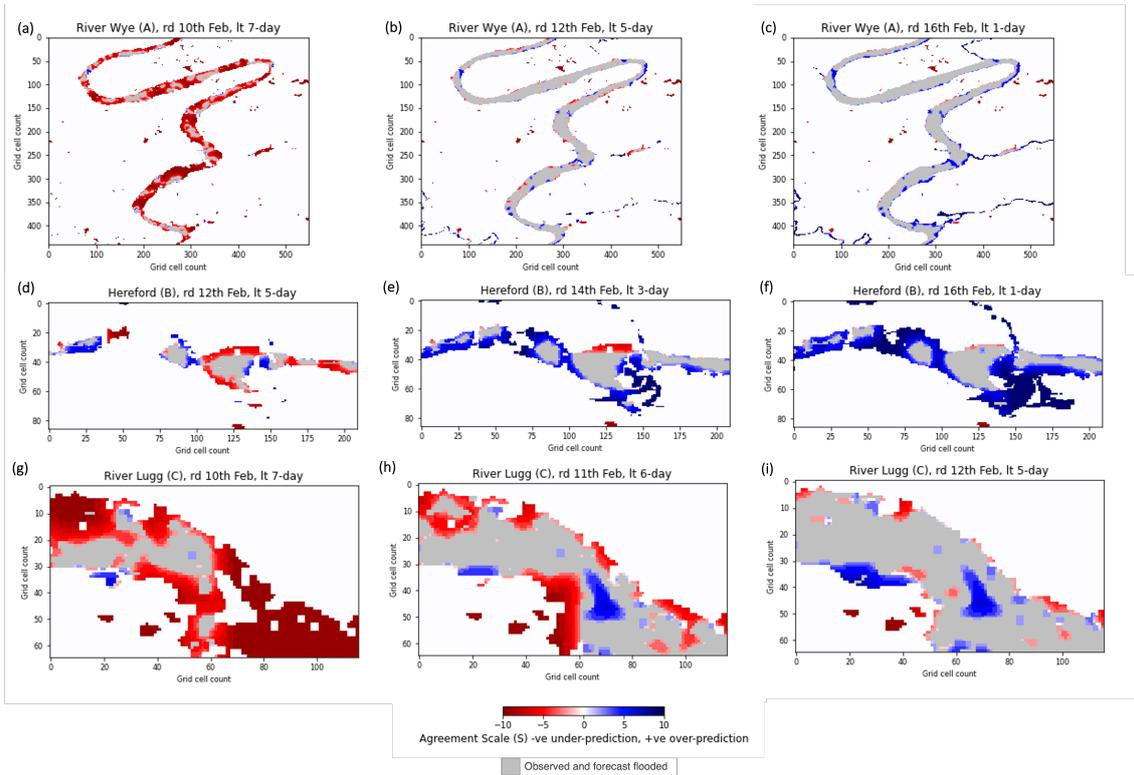


Figure 3.8: Categorical scale maps for each domain at various lead times (lt). Red indicates where the forecast flood extent is under-predicted, blue indicates over-prediction. The shading indicates the agreement scale, a measure of distance between the forecast and observed flood maps. Grey areas are correctly predicted flooded, white areas are correctly predicted unflooded. Each grid cell represents 25 m x 25 m for all domains. (Note: rd (forecast run date) varies between location, all dates have been evaluated and the most illustrative maps selected.)

2 The location-specific skilful scale varies with location and lead time as indicated on
 3 the categorical scale maps. For a 7-day lead time forecast for the River Wye (Fig. 3.8a),
 4 the model is indicating some flooding could occur, although under-estimating the total
 5 extent as show by the darkest red areas, which show the limits of the agreement scale
 6 have been reached. By 5-days lead time the forecast is in very close agreement with the
 7 observed flood at grid level (in grey) with larger agreement scales indicated by red/blue
 8 shading along some of the flood-edge locations (Fig. 3.8b) and a balance between under-
 9 and over-prediction. Over-prediction is more evident by 1-day lead time for the River

1 Wye (Fig. 3.8c) and flooding is also over-predicted along smaller tributaries. There are
2 several detached areas of flooding observed remotely that are most likely due to ponding
3 of surface water flooding, which were not predicted by the fluvial flood forecasting system.

4

5 The Hereford forecast is most skilful on the 12th February (Fig. 3.8d) with over-
6 prediction, particularly towards the southwest at 3-day and 1-day lead times (Fig.3.8e
7 and f). A small stream running southwards to the Wye, the Eign Brook, could be con-
8 tributing to the over-prediction seen here. It is also worth mentioning that SAR will
9 struggle to detect flood waters where buildings are closer together when the distance be-
10 tween them is less than the ground resolution of the SAR. Shadow and layover effects due
11 to the side-looking nature of the SAR also mean flood detection is more difficult in urban
12 areas (Mason et al., 2021a). This will likely only impact a small area of the Hereford
13 domain but this observation uncertainty should be considered when interpreting these re-
14 sults. There is an area of under-prediction of the flood extent in the centre of the Hereford
15 domain visible at all lead times. This could be due to surface water flooding, which most
16 likely occurred due to the very high intensity rainfall observed. This combined with the
17 urban area and steeply sloping gradient to the north of this area most likely contributed
18 to rapid surface water runoff towards the river. Since Flood Foresight is a fluvial flood-
19 ing forecast system we would not expect surface water flooding such as this to be predicted.

20

21 Flood Foresight selects multiple flood maps and stitches them together when the return
22 period threshold is exceeded for a given area. The Hereford section of the Wye does not
23 trigger a flood map selection until a 5-day lead time, this area also influences part of the
24 River Lugg flood map and can be seen as a mismatch on the lower left hand side of Figure
25 3.8g and h. Once this is included the forecast flood map is in very good agreement from a
26 5-day lead time. There are areas that could be further improved, indicated by the lighter
27 shading (Fig. 3.8i). An acceptable level of agreement scale could be determined for a given

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

4 **3.6.4 SAR-derived flood map scale variation**

5 In practice, particularly where a flood event is prolonged or the flood extent covers a wide
6 area, there may be multiple sources of SAR data available for model evaluation, usually
7 at higher spatial resolutions compared to the model grid size (e.g. ICEYE in spot mode
8 at 1 m and strip mode at 3 m ground resolution). It is important to consider the impact
9 of using observations at different spatial scales on the scale-selective approach results. By
10 conducting a simulation experiment we address the question of how re-scaling and inter-
11 polating three higher spatial resolution SAR-derived flood maps (relative to the forecast
12 flood maps) affects the scale selective skill scores and location-specific forecast skill. In
13 order to simulate a range of observation spatial scales, SAR-derived flood maps are pro-
14 duced using method two described in Section 3.4.2 at spatial resolutions from 5 m to 25
15 m. These are re-scaled by 0-order spline interpolation (ndimage.zoom, 2021; Briand &
16 Monasse, 2018) to match the model resolution (25 m) and compared to the forecast flood
17 map for the River Lugg (5-day lead time). A comparison of the GEE flood map against
18 the HASARD flood map, both at 20 m spatial scale produce almost identical verification
19 scores for all performance measures for the River Lugg ($\Delta FSS < 0.01$).

20

21 The categorical scale maps for the comparison between the forecast flood map and
22 the re-scaled simulated SAR-derived flood maps are shown in Figure 3.9. The resulting
23 domain-averaged skill scores for the same forecast flood map against the four SAR-derived
24 flood maps are displayed in Figure 3.10. The scores are calculated for the whole flood and
25 the flood edge alone. In general, the categorical scale maps show similar regions of over
26 and under-prediction but there are small location-specific variations in skilful scale. The

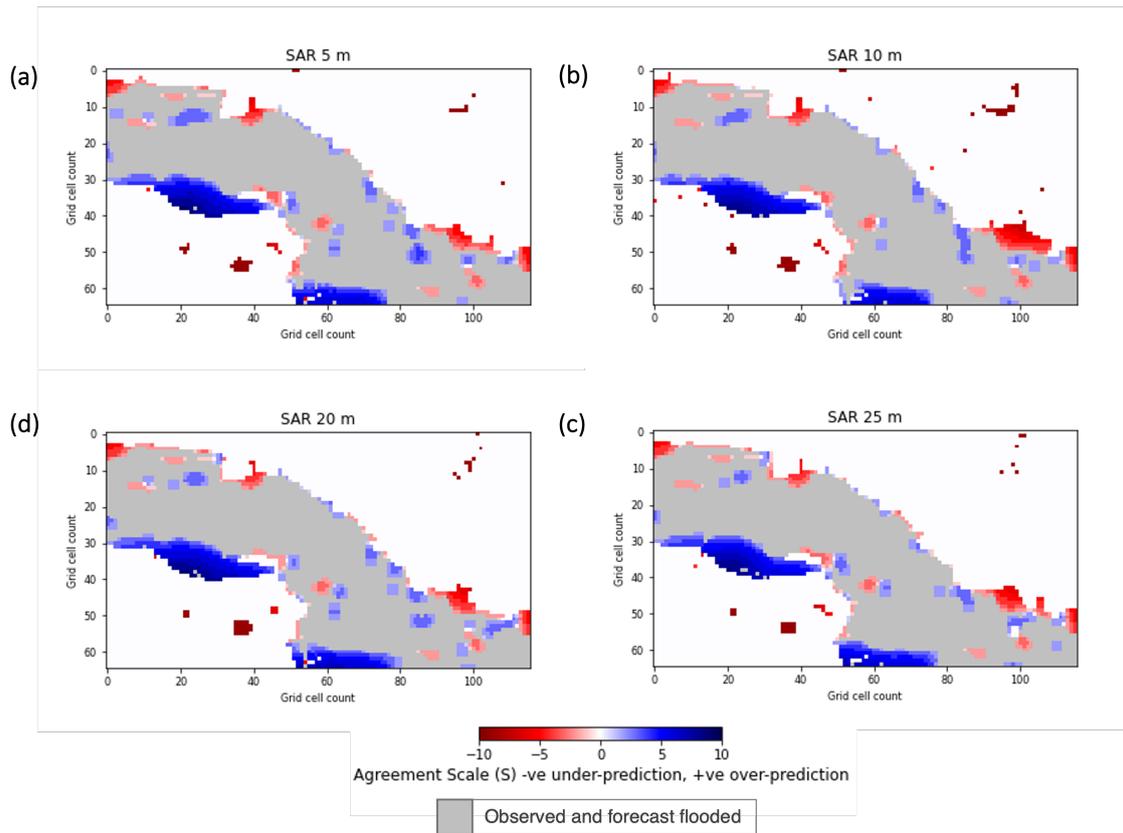


Figure 3.9: SAR-derived flood maps produced at different spatial resolutions (5 m to 25 m) are re-scaled to the model grid size (25 m) before categorical scale maps are calculated for the River Lugg (C), run date 12th Feb.

1 SAR-derived flood map at 25 m, the same spatial scale as the forecast flood maps, shows
 2 the best agreement away from the flood edge. This is also evident in the overall FSS score
 3 for the 25 m comparison, which marginally outperforms the evaluation after re-scaling finer
 4 observation flood maps (Fig. 3.10). The skilful scale determined for each observation com-
 5 parison of the whole flood extent is $n = 1$ or at grid level, and for the flood edge is at $n = 5$.
 6

7 Overall, based on the results from this simulation experiment, the scale-selective ap-
 8 proach is not overly sensitive to the observation spatial scale and the skilful scale deter-
 9 mined remains the same for each of the observed SAR-derived flood maps for both the

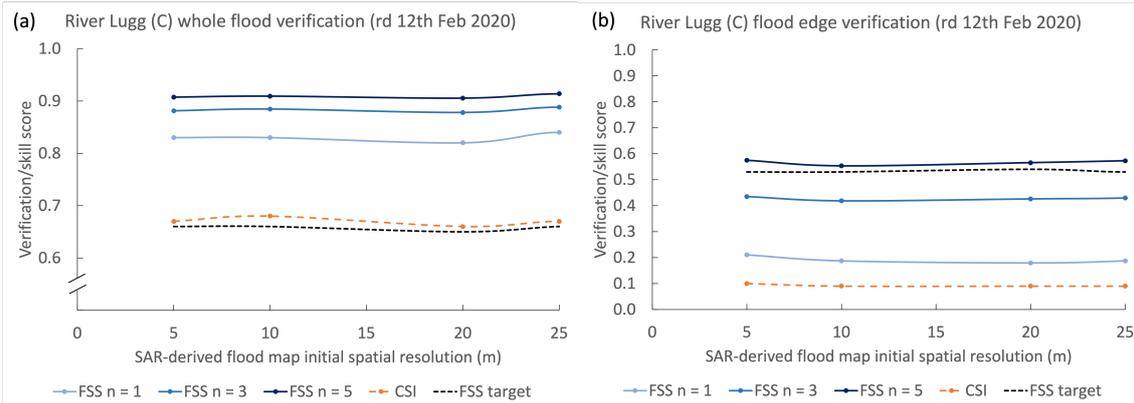


Figure 3.10: SAR-derived flood maps at different spatial resolutions (5 m to 25 m) are re-scaled to the model grid size (25 m) before verification scores are calculated for the whole flood (a) and the flood-edge (b). Note that axes in (a) and (b) are on different scales.

1 entire flood extent and the flood edge. Small errors are introduced by re-scaling and in-
 2 terpolating finer resolution observations to the model spatial scale which slightly reduce
 3 the skill score and change location-specific details on the categorical scale maps. Obser-
 4 vation scale selection and re-scaling along with interpolation errors must be considered
 5 when evaluating model performance, particularly where model or observation scales vary
 6 in space and time, or where comparisons are made across different models.

7 3.7 Discussion and Conclusions

8 Overall, the aim of this paper was to introduce and apply a new scale-selective approach
 9 to forecast flood map evaluation with an emphasis on providing a physically meaningful
 10 verification of the flood-edge location. The skilful spatial scale for comparison of forecast
 11 flood inundation maps against SAR-derived observed flood extent has been evaluated by
 12 the application of the Fraction Skill Score: this provides a domain-averaged skilful scale.
 13 The verification measure has been applied to a forecast of an extreme flood event in the
 14 UK on the River Wye and the River Lugg following Storm Dennis in February 2020. Flood
 15 Foresight inundation predictions with lead times out to 10 days are evaluated against a

1 Sentinel-1 SAR-derived flood map captured close to the flood peak for three domains, each
2 differing in hydrological characteristics. Conventional binary performance measures were
3 calculated alongside the FSS for comparison. Flood-edge verification shows greater sen-
4 sitivity to changes in forecast skill and spatial scale, relative to verification of the entire
5 flood extent. The skilful scale determined is physically meaningful and can be used to
6 estimate the average flood-edge discrepancy from the observed flood-edge. The observed
7 flood map spatial resolution relative to the model scale is important and re-scaling and
8 interpolation errors will impact the model verification scores. Ideally, the observed flood
9 map should be derived at the same spatial scale as the forecast model to minimise these
10 errors.

11

12 In operational practice the scale at which the forecast flood maps are presented to
13 forecasters and decision makers should reflect the uncertainty within the forecast. Very
14 high resolution flood maps can be presented where a detailed DTM is available. If this is
15 presented as a deterministic forecast to flood risk management teams, it could lead to an
16 over confidence in the forecast, or where the actual observed flood magnitude is different,
17 the forecast may be devalued in the future (Savage et al., 2016; Speight et al., 2021).
18 Application of a spatial-scale approach to forecast evaluation can determine the scale at
19 which it is best to present the forecast flood map. Conversely, if the model is found to
20 be skilful at grid level, there is scope to increase the flood map resolution adding more
21 detail to the flood-edge location. Improvements made to hydrodynamic models, such as
22 through data assimilation to improve inputs, initial conditions or model parameters may
23 not improve the forecast flood-edge location at grid level. However, improvements may
24 be evident through evaluation using FSS across a range of scales. Categorical scale maps
25 are a useful evaluation and forecasting tool, adding location-specific detail. Model im-
26 provements can be spatially targeted and as improvements are made, the categorical scale
27 map will highlight location-specific changes. For example, the categorical scale maps for

1 Hereford indicate the local infrastructure (in particular bridges) impact the movement of
2 the flood wave, which suggests a digital surface model (DSM) would be beneficial in urban
3 areas.

4

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

21

22 Ideally, in operational forecast systems, quantitative validation should run in tandem
23 with the forecast system where observations are available. Over time, a catalogue of skil-
24 ful scales, flood edge discrepancy distances and categorical scale maps could be built up.
25 This catalogue would enable analysis of scale across different flood event type, season,
26 meteorological scenario, forecast lead time and at specific locations within a catchment or
27 sub-catchment. Such a verification library would enable forecasters to increase intuition

1 and expert judgement on the relevant scales for a given forecast. Based on this analysis and
2 an increased understanding of the predictability of flood inundation, forecast flood maps
3 could be presented at a variable scale. For example, a coarser scale at longer lead times
4 becoming more detailed, closer to the flooding event. Coarse scales can appear jagged
5 or with large steps along the edge and so ideally these would be converted to smooth
6 contours, but with some indication (for example, lighter shading) that the flood edge lies
7 somewhere within the width of the grid cell, rather than exactly at the contour edge. At
8 shorter lead times, as forecast confidence is assumed to increase, the flood edge location
9 would show more detail and a narrower band of uncertainty (grid cell width). This flood
10 edge uncertainty information will prove invaluable for impact-based forecasting practice.

11

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

22 **3.8 Chapter summary**

23 In this chapter, we described and applied a scale-selective approach as a new verification
24 tool in flood inundation forecasting. The scale-selective approach overcomes issues with
25 conventional binary performance measures such as spatial scale dependency and flood

1 magnitude biases. We found that verification of the flood edge gave more sensitive skill
2 scores compared to verifying the whole flood extent. The resultant skilful scale can be
3 converted into a discrepancy distance between the forecast and the observed flood edge
4 location. Finally, a quantitative location specific agreement scale can be plotted on a
5 categorical scale map, which also indicates whether the forecast is accurate, over- or under-
6 predicting the flood extent at each grid cell location. This enables targeted improvements
7 to be made in operational practice. In the next chapter, we develop the method further
8 by considering the evaluation of the spatial spread-skill of an ensemble flood map forecast.

1 Chapter 4

2 Assessing the spatial spread-skill 3 of ensemble flood maps with 4 remote sensing observations

5 In this chapter we address the second research question outlined in Chapter 1; How skilfully
6 does an ensemble of forecast flood maps represent the spatial uncertainty within the flood
7 forecast?:

- 8 • How can we summarise the spatial predictability information in ensemble flood map
9 forecasts?
- 10 • How can we evaluate the spatial spread-skill of an ensemble flood map forecast?
- 11 • How does the spatial spread-skill vary with location and how can this be presented?

12 The remainder of this chapter (except for the chapter summary, Section 4.7), has been
13 published and is reproduced from (Hooker et al., 2023a).

The Brahmaputra River, Assam, India



Hugh DALTON



1 4.1 Abstract

2 An ensemble of forecast flood inundation maps has the potential to represent the uncer-
3 tainty in the flood forecast and provide a location specific likelihood of flooding. Ensemble
4 flood map forecasts provide probabilistic information to flood forecasters, flood risk man-
5 agers and insurers and will ultimately benefit people living in flood prone areas. Spatial
6 verification of the ensemble flood map forecast against remotely observed flooding is impor-
7 tant to understand both the skill of the ensemble forecast and the uncertainty represented
8 in the variation or spread of the individual ensemble member flood maps. In atmospheric
9 sciences, a scale-selective approach has been used to evaluate a convective precipitation
10 ensemble forecast. This determines a skilful scale (agreement scale) of ensemble perfor-

1 mance by locally computing a skill metric across a range of length scales. By extending
2 this approach through a new application, we evaluate the spatial predictability and the
3 spatial spread-skill of an ensemble flood forecast across a domain of interest. The spatial
4 spread-skill method computes an agreement scale at every grid cell between each unique
5 pair of ensemble flood maps (ensemble spatial spread) and between each ensemble flood
6 map with a SAR-derived flood map (ensemble spatial skill). These two are compared
7 to produce the final spatial spread-skill performance. These methods are applied to the
8 August 2017 flood event on the Brahmaputra River in the Assam region of India. Both
9 the spatial-skill and spread-skill relationship vary with location and can be linked to the
10 physical characteristics of the flooding event such as the location of heavy precipitation.
11 During monitoring of flood inundation accuracy in operational forecasting systems, valida-
12 tion and mapping of the spatial spread-skill relationship would allow better quantification
13 of forecast systematic biases and uncertainties. This would be particularly useful for un-
14 gauged catchments where forecast streamflows are uncalibrated and would enable targeted
15 model improvements to be made across different parts of the forecast chain.

16 **4.2 Introduction**

17 Forecast flood maps indicating the extent and depth of fluvial flooding within an action-
18 able lead time, are a useful tool for flood risk managers and emergency response teams
19 prior to and during a flood event. Typically, forecast flood maps are presented as deter-
20 ministic forecasts showing precisely where flooding will occur. This can lead to incidents
21 of false alarms or missed warnings and subsequent recriminations causing mistrust in the
22 system (Arnal et al., 2020; Savage et al., 2016). A timely prediction of exactly where
23 and when fluvial flooding caused by intense or prolonged rainfall will occur is virtually
24 impossible due to the chaotic nature of the atmosphere (Lorenz, 1969). The ensemble
25 forecasting approach aims to address the sensitivity of the atmospheric dynamics to initial

1 conditions and through multiple model runs these initial condition uncertainties can be
2 quantified (Leutbecher & Palmer, 2008). The ensemble forecast results in a probabilistic
3 weather forecast that indicates the predictability of the atmosphere at a given space and
4 time. State-of-the-art operational ensemble flood forecasting systems link together a chain
5 of forecast models to produce probabilistic streamflow and flood inundation forecasts at
6 national and global scales (Cloke & Pappenberger, 2009; Emerton et al., 2016; Wu et al.,
7 2020). Ensemble Numerical Weather Prediction models provide meteorological inputs into
8 land-surface, hydrological and hydraulic models, cascading the atmospheric uncertainty
9 through to the flood forecast. Throughout this chain of models, multiple sources of un-
10 certainty exist that have been investigated in numerous studies (Beven, 2016; Matthews
11 et al., 2022; Pappenberger et al., 2005; Zappa et al., 2011). As discussed by Boelee et al.
12 (2019), these uncertainties include those arising from meteorological inputs, measurements
13 and observations, initial conditions, unresolved physics within the models and parameter
14 estimates. A probabilistic flood inundation forecast should present a meaningful predic-
15 tion of the likelihood of flooding so that there is confidence in the forecast, given the
16 uncertainties represented in the system (Alfonso et al., 2016).

17

18 The accuracy of the location of flooding, predicted in advance, is defined as spatial
19 predictability. The spatial predictability of ensemble forecasts of flood inundation could be
20 verified by comparing with a remote observation of the flood from satellite or unmanned
21 aerial vehicle (UAV) based sensors. Satellite-based optical and Synthetic Aperture Radar
22 (SAR) sensors are well known for their flood detection capability (e.g., Horritt et al., 2001;
23 Mason, Davenport, et al., 2012; Mason, Schumann, et al., 2012). SAR sensors are active,
24 which enables them to scan the Earth through weather and clouds, and at night. The SAR
25 backscatter intensity detected depends on the roughness of the surface, with unobstructed
26 flooded areas and other surface water bodies appearing relatively smooth and returning
27 low backscatter values. Dasgupta et al. (2018) detail some of the challenges along with ap-

1 proaches to solutions of flood detection using SAR, examples of these challenges include:
2 roughening of the water surface by heavy rain and strong wind, emergent or partially
3 submerged vegetation and flood detection in urban areas. Accurate flood detection in
4 urban areas particularly due to surface water flooding has become increasingly important
5 (Speight et al., 2021) and recent techniques have led to improved flood detection (Mason
6 et al., 2018, 2021a, 2021b). Optical instruments rely on solar energy and cannot pen-
7 etrate cloud, making them less useful during a flooding situation. Recent studies have
8 investigated the flood detection benefits from combining both optical and SAR imagery
9 (Konapala et al., 2021; Tavus et al., 2020). Improvements in the spatial-temporal resolu-
10 tion of SAR images and their open source availability mean that they are an increasingly
11 valuable tool for hydraulic and hydrodynamic model improvements through calibration,
12 validation and data assimilation (e.g., García-Pintado et al., 2015; Grimaldi et al., 2016;
13 Cooper et al., 2018, 2019; Di Mauro et al., 2021; Dasgupta et al., 2018, 2021a, 2021b). The
14 Global Flood Monitoring (GFM) product (EU Science Hub, 2021; GFM, 2021; Hostache
15 et al., 2021) of the Copernicus Emergency Management Service (CEMS) (Copernicus Pro-
16 gramme, 2021) produces SAR-derived flood inundation maps for every Sentinel-1 image
17 detecting flooding. Three flood detection algorithms provide uncertainty estimation and
18 population affected estimates within 8 hours of the image acquisition. The European
19 Space Agency (ESA) Copernicus Programme have recently included the ICEYE constel-
20 lation of small satellites into the fleet of missions contributing to Europe’s Copernicus
21 environmental monitoring programme (ESA, 2021). ICEYE captures very high resolution
22 (spot mode ground range resolution = 1 m) SAR images which brings the potential for
23 increased accuracy of flood detection, particularly in urban areas.

24

25 To evaluate the accuracy of an ensemble forecast, a number of verification measures
26 have been proposed. Anderson et al. (2019) developed a joint verification framework for
27 end-to-end assessment of the England and Wales Flood Forecasting Centre (FFC) ensem-

1 ble flood forecasting system. Anderson et al. (2019) describe verification metrics such as
2 the continuous rank probability score (CRPS), rank histograms, Brier Skill Score (BSS)
3 and the relative operative characteristics (ROC) diagrams that are commonly applied to
4 assess the main ensemble attributes desirable in both precipitation and streamflow ensem-
5 ble forecasts (e.g., Renner et al., 2009). These metrics refer to flooding events as part of
6 a time series evaluated against a reference benchmark, such as climatology, to produce an
7 average skill score. In contrast, here we consider ensemble *spatial* verification at a single
8 time point. The verification of ensemble forecasts usually involves comparing the RMSE of
9 the ensemble mean against an observed quantity to assess the *skill* of the forecast with the
10 ensemble standard deviation used as a measure of *spread*. A perfect ensemble should en-
11 compass forecast uncertainties such that the ensemble spread is correlated to the RMSE of
12 the forecast (Hopson, 2014). This *spread-skill* relationship was assessed by Buizza (1997)
13 to investigate the predictability limits of the European Centre for Medium-Range Weather
14 Forecasts (ECMWF) Ensemble Prediction System (EPS). This approach to ensemble ver-
15 ification is based on point values and makes the assumption that the ensemble mean is the
16 forecast state with the highest probability and that the forecast distribution is Gaussian.
17 Significant flooding events are, in their nature, a rare occurrence and in certain circum-
18 stances a few ensemble members can indicate a low probability of an extreme flood. Also,
19 in particular atmospheric scenarios the ensemble forecast may result in a multi-modal
20 forecast where two clusters of ensemble members are each equally likely (Galmiche et al.,
21 2021). For example, both clusters may indicate flooding events but at different magni-
22 tudes. In both of these instances the individual ensemble member details are important
23 and evaluation of the ensemble mean alone would not be meaningful. When mapping the
24 flood extent prediction, the ensemble mean field alone does not retain the spatial detail of
25 the individual member forecasts.

26

27 The spatial spread-skill of the ensemble forecast is determined by evaluating the full

1 ensemble against observations of flooding. For a flood map ensemble to be considered
2 spatially well-spread, the spread or variation between ensemble members should equal the
3 spatial predictability, or skill of the ensemble members (Dey et al. (2014), see Section
4 4.3). Presently, to the best of our knowledge, quantitative evaluation methods assessing
5 the spatial spread-skill of ensemble forecast flood maps do not exist. However, previous
6 work in numerical weather prediction by Ben Bouallègue and Theis (2014) investigated
7 the application of spatial techniques to ensemble precipitation forecasts using a neighbour-
8 hood, or fuzzy approach that allowed comparisons at larger scales than grid level (native
9 resolution). A location dependent approach to the spatial spread-skill evaluation of a
10 convective precipitation ensemble forecast was developed by Dey, Roberts, et al. (2016).
11 This method compares every ensemble member across a range of scales on a spatial field
12 against an observation field to assess whether the ensemble forecast is spatially over-,
13 under- or well-spread on average across a domain of interest (Chen et al., 2018). In a
14 recent study, a scale-selective approach was developed and applied to evaluate a deter-
15 ministic flood map forecast where comparisons were made against conventional binary
16 performance measures (Hooker et al., 2022). A scale-selective approach to flood map eval-
17 uation was found to have several benefits over conventional binary performance measures.
18 These include over-coming the double penalty impact problem when validating at higher
19 spatial resolutions and accounting for the impact of the flood magnitude on the skill score.
20 The work described here extends and applies this scale-selective approach to assess the
21 spatial predictability and the spatial spread-skill of an ensemble flood map forecast.

22

23 In this paper we aim to address the following questions:

- 24 • How can we summarise the spatial predictability information in ensemble flood map
25 forecasts?
- 26 • How can we evaluate and visualise the spatial spread-skill of an ensemble flood map

1 forecast?

2 • How does the spatial spread-skill vary with location and how can this be presented?

3 In Section 4.3 we present a new approach to the evaluation of spatial predictability and
4 the spatial spread-skill of an ensemble flood map forecast by comparing against a remotely
5 observed flood extent. We illustrate the features of the methods through an example case
6 study of an extreme flooding event of the Brahmaputra River which impacted India and
7 Bangladesh in August 2017, with focus on the Assam region of India. The flood event
8 details are described in Section 4.4.1. The international ensemble version of the JBA
9 Consulting Flood Foresight system provides forecast flood maps for the study and is
10 described in Section 4.4.2. Observations of the flood are derived from satellite based SAR
11 sensors and the method is explained in Section 4.4.3. The results including the Spatial
12 spread-skill (SSS) map are discussed in Section 4.5. Our results show that individual
13 ensemble member spatial predictions of flooding are meaningful and that the full ensemble
14 spatial detail should be evaluated. We conclude in Section 4.6 that the spatial spread-
15 skill of the ensemble forecast varies with location across the domain and can be linked to
16 physical characteristics of the flooding event.

17 **4.3 Ensemble flood map spatial predictability evaluation** 18 **methods**

19 In this Section we present new methods for evaluating and visualising the spatial-spread
20 skill of an ensemble flood map forecast. Hooker et al. (2022) described and applied a new
21 scale-selective approach to evaluate the spatial skill of a *deterministic* flood map forecast
22 relative to an observed SAR-derived flood map. Here, we apply this same measure to
23 evaluate different aspects of an *ensemble* forecast. The scale-selective Fraction Skill Score
24 (FSS) method is outlined in Section 4.3.1. Agreement scale maps indicating forecast ac-

1 curacy are defined for location-specific comparisons between forecast and observed flood
 2 maps in Section 4.3.2. These are used to assess the spatial relationship between each
 3 unique pair of ensemble member flood maps (member-member) and between every ensem-
 4 ble member flood map and the observed SAR-derived flood map (member-SAR, Section
 5 4.3.3). Visualisation methods of the spatial spread-skill relationship including the *Spatial*
 6 *Spread-Skill* (SSS) map are presented in Section 4.3.4.

7 4.3.1 Fraction Skill Score

The FSS is a scale-selective verification measure that can determine the skilful scale of a modelled flood map, when compared against a remotely sensed observation of flooding (Roberts & Lean, 2008; Hooker et al., 2022). We will call these flood maps the *model array* and the *observed array* respectively. For an ensemble forecast, the model array could be an individual ensemble member, or a summarised flood estimate derived from a combination of ensemble members such as a combined ensemble or the ensemble median (see Section 4.4.4). Both the model and observed arrays are converted into binary fields using a situation dependent threshold (e.g. depths greater than 0.2 m are labelled flooded). For this ensemble application of the FSS we evaluate the entire flood extent across the domain. Each grid cell is labelled as inundated (1) or dry (0). All grid cells are numbered according to their spatial locations (i, j) , $i = 1 \dots N_x$ and $j = 1 \dots N_y$ where N_x is the number of columns and N_y is the number of rows. Surrounding each grid cell, a square of length n creates an $n \times n$ neighbourhood. The fraction of 1s (inundated cells) in the square neighbourhood area is calculated for every grid cell. This creates two arrays of fractions across the domain for both the observed O_{nij} and modelled M_{nij} data. The mean squared error (MSE) for the fraction arrays is calculated for the domain and a given neighborhood size, n :

$$MSE_n = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [O_{nij} - M_{nij}]^2. \quad (4.1)$$

A potential maximum $MSE_{n(ref)}$ depends on the fraction of flooding in the domain for the modelled and observed fields and is calculated as:

$$MSE_{n(ref)} = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [O_{nij}^2 + M_{nij}^2]. \quad (4.2)$$

Finally, the FSS is

$$FSS_n = 1 - \frac{MSE_n}{MSE_{n(ref)}}. \quad (4.3)$$

The FSS is initially calculated at grid level ($n = 1$) followed by the smallest neighbourhood size ($n = 3$) before increasingly larger neighbourhood sizes ($n = 5, n = 7...$) are considered. The FSS ranges between 0 (no skill) and 1 (perfect skill). Increasing the neighbourhood size typically leads to an improved FSS as the fractions are calculated over a larger area. Plotting FSS against the neighbourhood size can indicate a range of scales where the model is deemed to be the most skilful. A target FSS score (FSS_T) can be determined from the fraction of observed flooding across the whole domain (f_0):

$$FSS_T \geq 0.5 + \frac{f_0}{2}. \quad (4.4)$$

1 The point where the FSS_n exceeds FSS_T can be viewed as being equidistant between the
2 skill of a random forecast and perfect skill (Roberts & Lean, 2008). A recent study by
3 Skok and Roberts (2018) investigated the sensitivity of the calculated skilful scale to the
4 constant value (0.5) in Eq. (4), and found that 0.5 gave meaningful results compared with
5 the measured displacement. The magnitude of the observed flood, relative to the domain
6 area, determines the value of FSS_T . This allows the comparison of the skilful scale
7 (neighbourhood size) where FSS_T is reached across different domain sizes and floods of
8 different magnitudes.

1 4.3.2 Location dependent agreement scales

2 The FSS (Section 4.3.1) gives a domain average measure of forecast performance and
3 a minimum spatial scale at which the forecast is deemed skilful. To enable the spatial
4 spread-skill of the full ensemble to be evaluated at specific locations, we first define an
5 agreement scale (see Dey et al. (2014); Dey, Roberts, et al. (2016); Hooker et al. (2022)
6 for full methodology). The agreement scale is calculated and mapped for every grid cell
7 in the domain and shows a measure of similarity between two arrays of data. In contrast
8 to the FSS method the arrays are not required to be thresholded. The agreement scale
9 method can be applied to both binary flood extent maps as well as flood depth fields.
10 These could both be ensemble member flood maps or an ensemble member flood map
11 and an observed flood map. Two data arrays are compared F_{1ij} and F_{2ij} and the aim
12 is to find a minimum neighbourhood size (or spatial scale) for every grid cell such that
13 there is a predetermined acceptable minimum level of agreement between F_{1ij} and F_{2ij} .
14 This is known as the agreement scale $S_{ij}^{A(F_1 F_2)}$. (Note that the relationship between the
15 agreement scale and the neighbourhood size described previously in section 4.3.1 is given
16 by $S_{ij}^{A(F_1 F_2)} = (n - 1)/2$.) The agreement scale (now defined S for simplicity in the
17 following equations) is determined individually for every grid cell by testing and meeting
18 a chosen criteria.

A relative MSE, D_{ij}^S is calculated for all grid cells, initially at grid level, $S = 0$ ($n = 1$),

$$D_{ij}^S = \frac{(F_{1ij}^S - F_{2ij}^S)^2}{(F_{1ij}^S)^2 + (F_{2ij}^S)^2}. \quad (4.5)$$

If $F_{1ij} = 0$ and $F_{2ij} = 0$ (both dry) then $D_{ij}^S = 0$ (correct at grid level). The value of D_{ij}^S ranges between zero and 1. The arrays are deemed to be in agreement at the scale being tested if:

$$D_{ij}^S \leq D_{crit,ij}^S \quad \text{where} \quad D_{crit,ij}^S = \alpha + (1 - \alpha) \frac{S}{S_{lim}} \quad (4.6)$$

The parameter value α indicates an acceptable bias at grid level such that $0 \leq \alpha \leq 1$. Additional historical forecast data of flood events is not available for the region in this study, so we assume there is no background bias between the forecast and the observations and set $\alpha = 0$. A fixed maximum scale S_{lim} is predetermined using human judgement considering the physical characteristics of the flood event. The value chosen for S_{lim} depends on the magnitude of the flood extent relative to the size of the sub-catchment. For the case study presented here, we set $S_{lim} = 80$ (2400 m), which is approximately $\frac{1}{4}$ to $\frac{1}{2}$ of the sub-catchment widths in the domain. If $D_{ij}^S \geq D_{crit,ij}^S$ then the next neighbourhood size up is considered ($S = 1, n = 3$, a 3 by 3 square) where F_{1ij}^1 and F_{2ij}^1 are arrays containing the average value of each neighbourhood surrounding the grid cell at position (i, j) for each array. The process continues by comparing increasingly larger neighbourhoods (e.g. $S = 2, n = 5$, a 5 by 5 square) until the agreement criterion:

$$S_{ij}^{A(F_1 F_2)} \text{ or } S_{lim} \text{ at } D_{ij}^S \leq D_{crit,ij}^S \quad (4.7)$$

1 is met for every cell in the domain. The agreement scale at which the agreement criterion
 2 is met will usually vary from grid cell to grid cell and these values ($S = 0, S = 1, S = 2$ and
 3 so on up to S_{lim}), each specific to each grid cell location can be mapped onto the domain
 4 of interest to provide a location specific measure of agreement between the two data arrays
 5 that are compared. A small value for the agreement scale means that the two arrays being
 6 compared are very similar (spatially) at a specific location, whereas a large value for the
 7 agreement scale means that the two arrays being compared are dissimilar. Note that the
 8 skilful scale determined by the FSS (Section 4.3.1) differs from the agreement scale defined
 9 here. The former links directly with the spatial differences between objects e.g. Skok and
 10 Roberts (2018), whereas the latter reflects a pre-defined acceptable bias at different scales.

11

12 Validation of forecast flood maps against remotely observed flood extent is typically

1 carried out by labelling each grid cell using a contingency table with categories: cor-
2 rectly predicted flooded, under-prediction (miss), over-prediction (false alarm) and cor-
3 rectly predicted unflooded. In the contingency table under-predicted cells are set to +1,
4 over-predicted cells are set to -1, correctly predicted flooded cells are assigned NaN and
5 correctly predicted unflooded cells are set to 0. Mapping these categories creates a con-
6 ventional contingency map, which combined (by element-wise array product) with an
7 agreement scale map (Eq. (7)) creates a categorical scale map made by plotting the ab-
8 solute agreement scale values coloured according to the contingency class. A categorical
9 scale map shows a measure of spatial accuracy between two data arrays (Hooker et al.,
10 2022). Categorical scale maps may be used as a basis for comparison between ensemble
11 members and observations, as we illustrate with our case study in Section 4.5.3.

12

13 4.3.3 Ensemble spatial spread-skill evaluation

14 We assume that each ensemble forecast flood map represents an equally likely future sce-
15 nario and the evaluation of the full ensemble is needed to quantify the uncertainty and to
16 evaluate the spatial spread-skill relationship. The ensemble flood map spatial character-
17 istics vary with location and in order to preserve the location dependent information, we
18 utilise a method developed to evaluate a convective ensemble precipitation forecast (Dey,
19 Roberts, et al., 2016; Dey, Plant, et al., 2016). Here, we outline the method and describe
20 a new application to evaluate an ensemble forecast flood map.

21

A neighbourhood approach (Section 4.3.2) is used to assess the spatial agreement scale $S_{ij}^{A(F_1 F_2)}$ or measure of similarity at each grid cell location (i, j) between each unique pair of ensemble flood maps. For an ensemble of M members, there are

$$M_p = \frac{M(M-1)}{2}, \quad (4.8)$$

unique pairs (e.g., 1275 pairs for a 51 member ensemble). For an ensemble, the skillful scale can be renamed as a *believable scale*, which is the scale where ensemble members become sufficiently similar to observations such that they are a useful prediction. Every paired ensemble agreement scale field is averaged at each grid cell to produce a mean field, from the agreement scale field defined in Eq. (7)

$$S_{ij}^{A(\overline{mm})} = \frac{1}{M^p} \sum_{F_1=1}^{M-1} \sum_{F_2=F_1+1}^M S_{ij}^{A(F_1 F_2)} \quad (4.9)$$

indicating the location specific believable scales of the forecast flood map ensemble. Maps of $S_{ij}^{A(\overline{mm})}$ summarise the spatial spread of the full ensemble. Each of the agreement scale fields between the ensemble members and the observations are also averaged at each grid cell to give

$$S_{ij}^{A(\overline{m\bar{o}})} = \frac{1}{M} \sum_{f=1}^M S_{ij}^{A(F_f)}. \quad (4.10)$$

1 A measure of the spatial spread-skill of the ensemble can be found by comparing the
 2 average agreement scale between the ensemble members $S_{ij}^{A(\overline{mm})}$ representing the ensemble
 3 *spread* with the average agreement scale between the ensemble members and the observed
 4 flood field $S_{ij}^{A(\overline{m\bar{o}})}$ representing the ensemble *skill*.

5 4.3.4 Spatial spread-skill visualisation methods

6 To evaluate the spatial spread-skill relationship, $S_{ij}^{A(\overline{mm})}$ (representing the ensemble *spread*)
 7 must be compared in the same location as $S_{ij}^{A(\overline{m\bar{o}})}$ (representing the ensemble *skill*). Data
 8 arrays can be visually compared using a binned scatter plot that averages across a selected
 9 bin of cells at the same location within the domain. Dey, Roberts, et al. (2016) demon-
 10 strated for an idealised example that by plotting $S_{ij}^{A(\overline{mm})}$ against $S_{ij}^{A(\overline{m\bar{o}})}$ as a binned scatter
 11 plot in order to preserve the spatial location of the comparison (Fig. 4.1), the ensemble
 12 can be classified as over-, under- or well-spread. The ensemble is deemed to be *well-spread*

1 at a specific location in the domain of interest when the spread of the individual members
2 represented at each grid cell by $S_{ij}^{A(\overline{mm})}$ equals the skill of the ensemble represented at
3 each grid cell by $S_{ij}^{A(\overline{m\bar{o}})}$, i.e. $S_{ij}^{A(\overline{mm})} - S_{ij}^{A(\overline{m\bar{o}})} = 0$. The result would lie on a 1:1 line
4 on the binned scatter plot. Where the spread between the ensemble members exceeds
5 the skill of the ensemble forecast i.e. $S_{ij}^{A(\overline{mm})} > S_{ij}^{A(\overline{m\bar{o}})}$ the ensemble is considered to be
6 *over-spread* and the binned scatter plot will lie beneath the 1:1 line. The converse is true
7 for an *under-spread* ensemble forecast where the agreement between members, the spread,
8 is less than the agreement between the ensemble and the observations, the skill. Here,
9 $S_{ij}^{A(\overline{mm})} < S_{ij}^{A(\overline{m\bar{o}})}$ and the binned scatter plot would lie above the 1:1 line.

10

11 To summarise the spread-skill relationship we develop this visualisation further by
12 plotting a hexagonal binned 2D histogram plot (an example hexbin plot is presented
13 in Section 4.5.3). The domain is divided into a (pre-determined) number of hexagons.
14 Hexagons minimize the perimeter to area ratio and therefore minimize the edge effects.
15 The hexbin histogram plot colour shade represents the number of data points within each
16 bin.

17

18 Whilst the hexbin plot is useful for gaining an understanding of the general spread-skill
19 relationship of the ensemble flood map forecast, it does not tell us specifically where in
20 the domain the ensemble spatial predictability is better or worse. The *Spatial Spread-Skill*
21 (SSS) map plots $S_{ij}^{A(\overline{mm})} - S_{ij}^{A(\overline{m\bar{o}})}$ at every grid cell location so that the spread-skill is
22 mapped across the domain and can be linked directly to different sub-catchments and sur-
23 face features such as tributaries, embankments, bridges and importantly the underlying
24 topography or DTM, which influence the derivation of the ensemble flood maps. Regions
25 on the SSS map where the ensemble is over-spread are positive with negative areas in-
26 dicating where the ensemble is under-spread, zero values show a well-spread ensemble.
27 Note that this does not necessarily mean that the entire ensemble is in agreement with

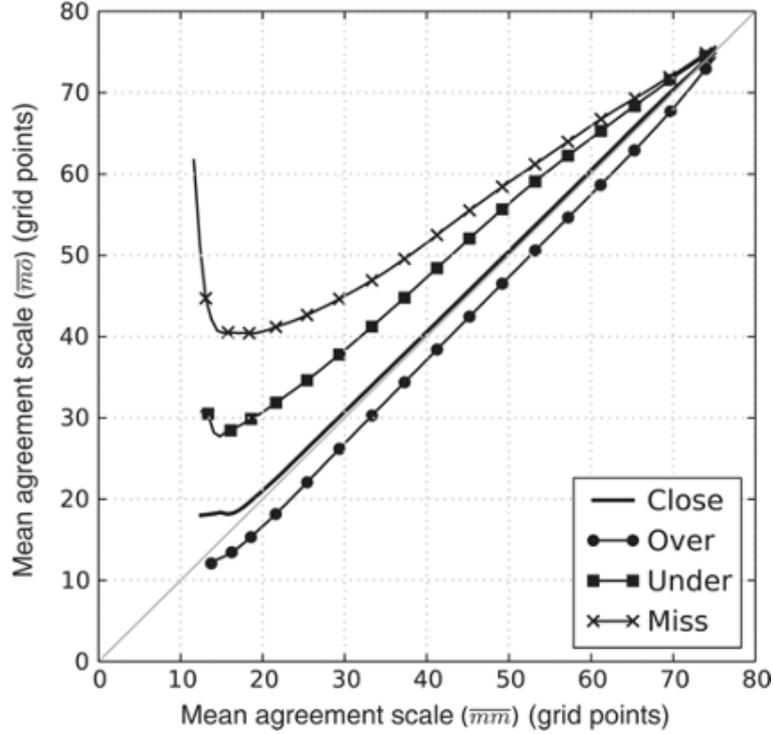


Figure 4.1: Figure reproduced with permission from Dey et al., (2016) showing results on a binned scatter plot from an idealised experiment indicating the spatial spread-skill relationship between an ensemble forecast and the observation.

- 1 observations at grid level, but that the agreement scales between $S_{ij}^{A(\overline{mm})}$ and $S_{ij}^{A(\overline{mo})}$ are
- 2 equal. (An example SSS map is presented in Section 4.5.3).

3 4.4 Ensemble forecasting flood event case study

- 4 In this section we describe an example flooding event used to demonstrate the application
- 5 of the spatial spread-skill evaluation approach. We evaluate a 1-day lead time flood inun-
- 6 dation 51 ensemble member forecast from the Flood Foresight system (Section 4.4.2) for
- 7 the domain area against a satellite SAR-derived flood map (Section 4.4.3).

1 **4.4.1 Brahmaputra flood, Assam India, August 2017**

2 The origin of the Brahmaputra River (also known as the Yarlung Tsangpo in Tibetan,
3 the Siang/Dihang River in Arunachali, Luit in Assamese, and the Jamuna River in
4 Bangladesh) lies in the Himalayan Kailas Range of southwestern Tibet, China. Draining
5 an area of 543,000 km², the Brahmaputra flows for 2000 km across the Tibetan Plateau
6 and a further 1000 km parallel to the Himalayan foothills through the Assam Valley, India
7 before entering Bangladesh where the Brahmaputra joins the Ganges River (Palash et
8 al., 2020). The Brahmaputra baseflow originates from the upstream glacial snow melt,
9 however the streamflow rates are dominated by the summer monsoon precipitation. The
10 basin receives up to 95% of its annual rainfall during the pre-monsoon and monsoon season,
11 which usually runs from April to September and causes annual flooding of the Brahma-
12 putra. The Assam region typically records on average 2300 mm of annual rainfall and up
13 to 5000 mm in the Himalayan foothills (Dhar & Nandargi, 2000, 2003).

14

15 For this example case we focus on the third wave of flooding that occurred during the
16 monsoon season in August 2017, peaking around the 12th. Figure 4.2 shows the location
17 of the Brahmaputra and of a chosen domain centred upon some of the worst flooding
18 that occurred. This area includes a confluence zone where the Subansiri River meets the
19 Brahmaputra. The monsoon flooding impacted an estimated 40 million people across India
20 and Bangladesh. Locally in the Assam region, the flooding in August affected over 3.3
21 million people and approximately 3200 villages, river embankments were damaged in 11
22 districts. Over 14,000 people were evacuated to one of around 700 relief camps that were
23 also needed to house over 180,000 people relocated (Floodlist, 2017). The local Assam
24 State Disaster Management Authority (ASDMA, 2017) flood early warning system issued
25 a low warning alert (disasters that can be managed at the district level) on the 10th August
26 for the district.

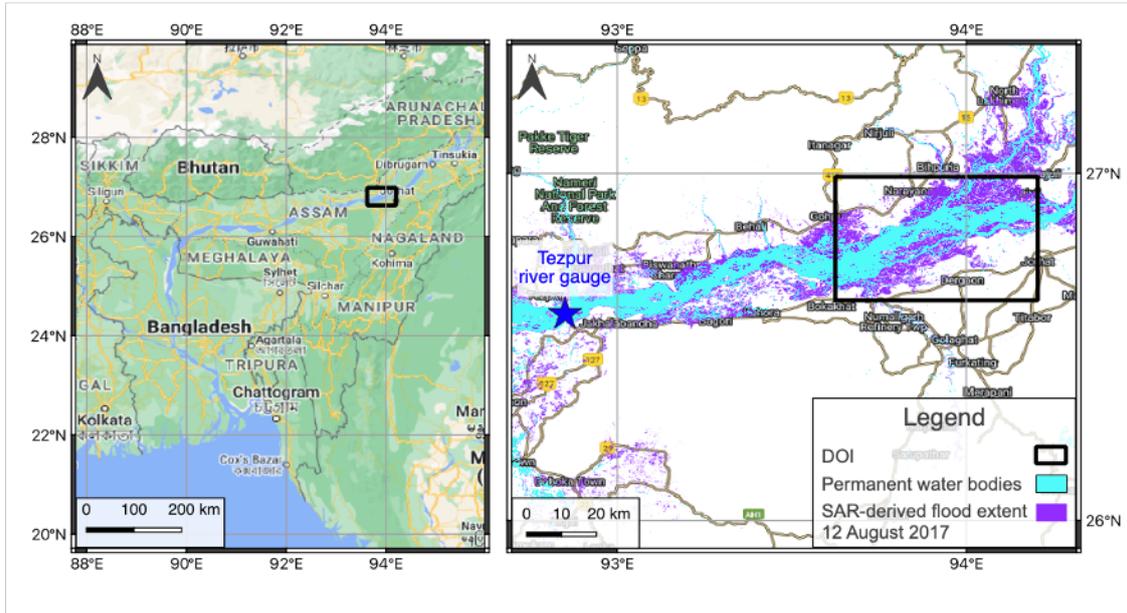


Figure 4.2: Left panel: domain location on the Brahmaputra River in the Assam region of India. Domain size is 57.5 km by 39.3 km. Right panel: Sentinel-1 SAR-derived flood map and permanent water bodies from the JRC Global Surface Water database for the domain of interest (DOI). Base map from ©Google Maps.

1 In 2017, the southwest monsoon season rainfalls were predicted to be *normal* by the
 2 South Asian Climate Outlook Forum (WMO, 2017). However, the pre-monsoon season
 3 began early in the year with heavy thunderstorms affecting the region from March on-
 4 wards. In the Assam region, June and July were 60% wetter than the previous three years
 5 and during August more locally intense rainfall was recorded compared with historical
 6 observations (Palash et al., 2020). In higher latitude areas, 30 km to the north of the
 7 domain at North Lakhimpur, 215.8 mm rainfall was recorded in the three days prior to
 8 the flood peak (Floodlist, 2017; Hossain et al., 2021). An above normal flood situation
 9 is declared in India where the river water level exceeds the Warning Level, a severe flood
 10 occurs where the water level exceeds the Danger Level, and an extreme flood occurs where
 11 the previous Highest Flood Level is exceeded (Central Water Commission, 2023). The
 12 peak water level recorded downstream at Tezpur (Danger Level 65.23 m) on August 14th
 13 was 66.12 m. There are regional variations in maximum water levels reported, with upland

1 regions to the north of the Assam valley recording water levels that exceed the previous
2 Highest Flood Level indicating an extreme flood level (Floodlist, 2017).

3

4 **4.4.2 Ensemble flood forecasting system**

5 The Flood Foresight system (Fig. 4.3), developed and operationally run by JBA Con-
6 sulting, is a fluvial flood inundation mapping system that can be implemented at any
7 river basin around the world. Flood Foresight utilises a simulation library approach to
8 generate real-time and forecast flood inundation and water depth maps. The simulation
9 library approach saves valuable computing time and allows the application of Flood Fore-
10 sight in near continuous real-time at national and international scales. A pre-computed
11 library of flood maps for a river basin or country are created using JFlow[®](where a DTM
12 is available), (Bradbrook, 2006) and RFlow (where a DTM is unavailable). JFlow uses
13 a raster-based approach with a detailed underlying digital terrain model (DTM) and a
14 diffusion wave approximation of the full 2D hydrodynamic shallow water flow equations.
15 RFlow combines a 1D model based upon Normal Depth calculations, optimised for use on
16 a Digital Surface Model (DSM, NEXTmap (2016)) with rapid 2D flood spreading (cre-
17 ated by spreading Normal Depth from upstream to downstream) and is calibrated against
18 JFlow. These equations capture the main controls of the flood routing for shallow, topo-
19 graphically driven flow. Six flood maps at 30 m resolution are created for 20, 50, 100, 200,
20 500 and 1500 year return period flood events (corresponding to annual exceedance prob-
21 abilities (AEPs) of 5%, 2.5%, 1%, 0.5% and 0.2% and 0.07% respectively). Between each
22 adjacent pair of modelled return period maps, five additional intermediate flood maps are
23 created by linear interpolation of both flood depth and extent. An additional five flood
24 maps are also created beneath the lowest return period flood map. This gives, in total, a
25 library of 36 flood maps. Note that these flood maps are undefended and local temporary
26 flood defences are not included. Flood foresight is set up for a region by dividing the river

1 basin into sub-catchments using the HydroBASINS data-set (level 12) (Lehner, 2014a).
 2 Flood Foresight takes gridded inputs of ensemble forecast streamflow and uses these to
 3 select the most appropriate flood map for each sub-catchment. These are mosaicked to-
 4 gether and forecasts of ensemble flood maps are produced daily, out to ten days ahead.

5

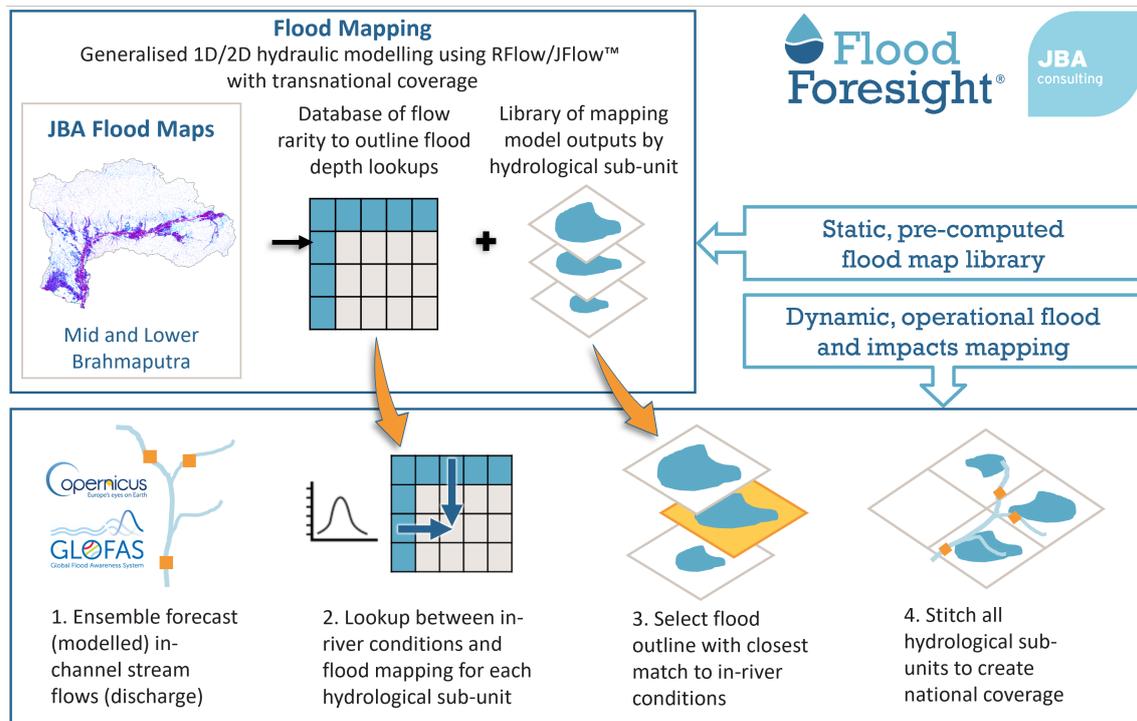


Figure 4.3: Flood Foresight ensemble forecast flood inundation and impact mapping work flow. Prepared by JBA Consulting.

6 The global (non UK and Ireland) configuration of Flood Foresight uses ensemble
 7 streamflow forecast data from the Global Flood Awareness System (GloFAS) (Alferi et
 8 al., 2013; GloFAS, 2021). GloFAS was jointly developed by the European Commission
 9 and the European Centre for Medium-Range Weather Forecasts (ECMWF) and is com-
 10 posed of an integrated hydro-meteorological forecasting chain that couples state-of-the-art
 11 weather forecasts with a land surface and hydrological model. With its continental scale
 12 set-up, GloFAS provides downstream countries with forecasts of upstream river conditions

1 up to one month ahead as well as continental and global overviews for large world river
2 basins. Meteorological forecast data are provided by the ECMWF Ensemble (IFS) model,
3 the operational (51 member) ensemble weather forecasting product of the ECMWF. The
4 meteorological forecast data provide inputs to the land surface module, HTESSSEL (Hydro-
5 logical Tiled ECMWF Scheme for Surface Exchange over Land). HTESSSEL simulates the
6 land surface response to the meteorological data, based on simulated interactions with soil
7 conditions, idealised vegetation cover and land cover. From these simulations, HTESSSEL
8 outputs forecast global surface and sub-surface flows per grid cell. These simulated flows
9 are then used by a simplified version of the hydrological model LISFLOOD, a 1D routing
10 model which simulates the movement of the surface and sub-surface flows. The runoff data
11 produced is routed through a representation of the river network using a double kinematic
12 wave approach, which includes bankfull and over bankfull routing. The river network used
13 is taken from the HydroSHEDS data-set (Lehner & Grill, 2013).

14

15 GloFAS outputs a gridded (approximately 10 km spatial resolution) ensemble forecast
16 of river streamflow (Fig. 4.4). Each of the GloFAS grid cells are linked to the sub-
17 catchments in the Flood Foresight system. The simulation library flood maps are selected
18 when the forecast streamflow exceeds a return period threshold level within each sub-
19 catchment. The RP threshold levels are calculated using ERA5 reanalysis data (Harrigan
20 et al., 2020). Each ensemble member flood map forecast is created by aggregating the
21 individual sub-catchment maps. In summary, the meteorological IFS 51 member ensem-
22 ble input to the flood forecasting chain allows atmospheric evolution uncertainties to be
23 represented within the ensemble streamflow forecast and the ensemble of inundation flood
24 maps, thus creating a probabilistic flood map forecast, indicating the likelihood of flood-
25 ing. Flood foresight produces daily ensemble flood depth and extent forecasts at 30 m
26 spatial resolution out to 10-days.

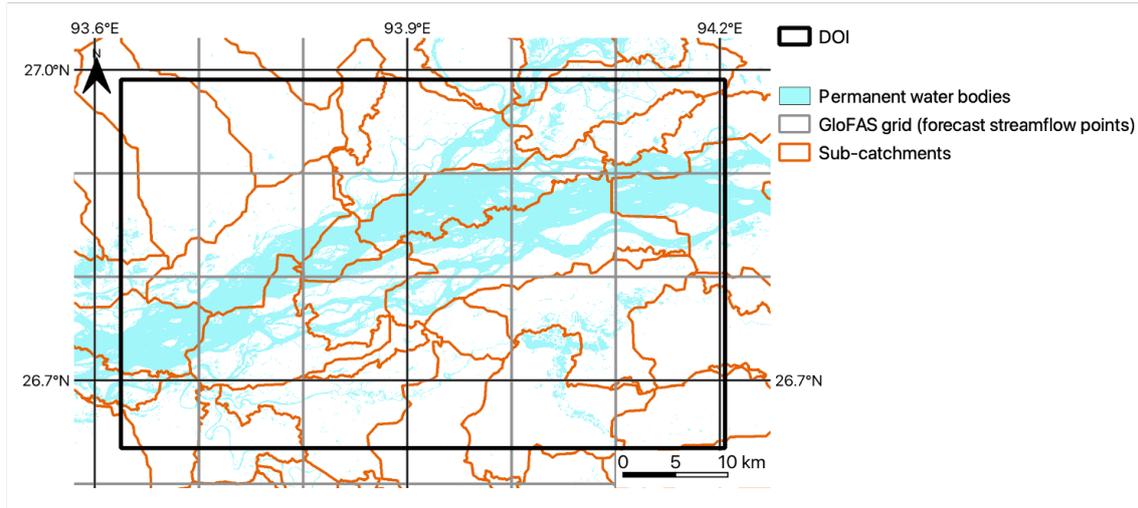


Figure 4.4: GloFAS grid, permanent water bodies and Flood Foresight sub-catchments for the domain of interest (DOI).

1 4.4.3 SAR-derived flood map

2 A Sentinel-1 (S1A) image was acquired in interferometric wide swath mode (swath width
3 250 km) around the time of the flood peak at 17:18 (IST) on the 12th August 2017. The
4 ESA Grid Processing on Demand (GPOD) HASARD service (service terminated June
5 2021) was utilised to map the flooding. The flood mapping algorithm (Chini et al., 2017)
6 uses an automated, statistical, hierarchical split-based approach to distinguish between
7 two classes (background and flood) using a pre-flood and flood image. A pre-flood image
8 (February 2017) from the same satellite sensor and track was used to derive the flood map
9 (Fig. 4.2). Original SAR images (VV polarisation) were pre-processed, which involved:
10 precise orbit correction, radiometric calibration, thermal noise removal, terrain correction,
11 speckle reduction and re-projection to the WGS84 coordinate system. The HASARD
12 mapping algorithm removes permanent water bodies that are detected on the pre-flood
13 image, such as the unflooded river water, lakes and reservoirs by applying a thresholding
14 approach. Flooded areas beneath vegetation, bridges and near to buildings will not be
15 detected using this method. Flood Foresight forecast flood maps include the river channel

1 and exclude surface features such as vegetation and buildings. To smooth the HASARD
2 flood maps and allow a fairer comparison we apply a morphological closing operation
3 (without impacting the location of the flood extent) to flood fill vegetation and buildings.
4 The wide and braided Brahmaputra River in the Assam region covers a significant area of
5 the selected domain. In order to evaluate the flood prediction accuracy alone, the pre-flood
6 occurrence of surface water using the JRC Global Surface Water database (Pekel et al.,
7 2016) has been removed from the Flood Foresight forecast inundation maps. The observed
8 flood extent derived from satellite based SAR data at 20 m grid size is re-scaled to match
9 the forecast flood map grid size (30 m) using average aggregation. The closest available
10 (cloud free) optical image available was a Sentinel-2 image on the 17th August 2017, 5 days
11 after the SAR image acquisition. During this time the flood waters had receded from their
12 peak, which makes this unsuitable for comparison with the SAR-derived flood map. Since
13 no other validation sources are available, for the purposes of this study we assume that the
14 SAR-derived observation of flooding represents the true flood extent. From October 2021,
15 Sentinel-1 SAR images are processed by CEMS GFM (GFM, 2021) to derive flood extent
16 and provide an uncertainty estimate of the grid cell classification. This means uncertainty
17 information in the SAR-derived flood map could be accounted for in future evaluation
18 studies by verifying across different levels of observation uncertainty. Additionally, a flood
19 mask, indicating areas where flood detection using SAR data is not currently possible
20 (at the Sentinel-1 spatial resolution) could be used to exclude areas from the evaluation
21 process (note that this was not possible for this case study, since this information was not
22 available in 2017).

23 **4.4.4 Forecast data**

24 Flood Foresight was set-up for the Brahmaputra basin in India and Bangladesh using the
25 simulation library approach to flood mapping described in Section 4.4.2. Flood maps were
26 pre-computed for the domain of interest (Fig. 4.2) using a DSM and RFlow. The forecast

1 data for the Brahmaputra flood event contains a 51 member ensemble of flood maps
2 indicating flood extent, produced at a 1-day lead time. Vertically stacking each individual
3 ensemble member flood map and adding vertically across every grid cell combines all
4 ensemble members into a single flood map (all flooded grid cells are set to 1) showing
5 where flooding is possible across all members (ens_{all}). A spatial median flood map is
6 created (ens_{median}) where 26 members or more predict flooding at a particular grid cell
7 location. Each of the ensemble member flood maps for the domain are plotted in Figure
8 4.5 along with ens_{all} , ens_{median} and the SAR-derived flood map.

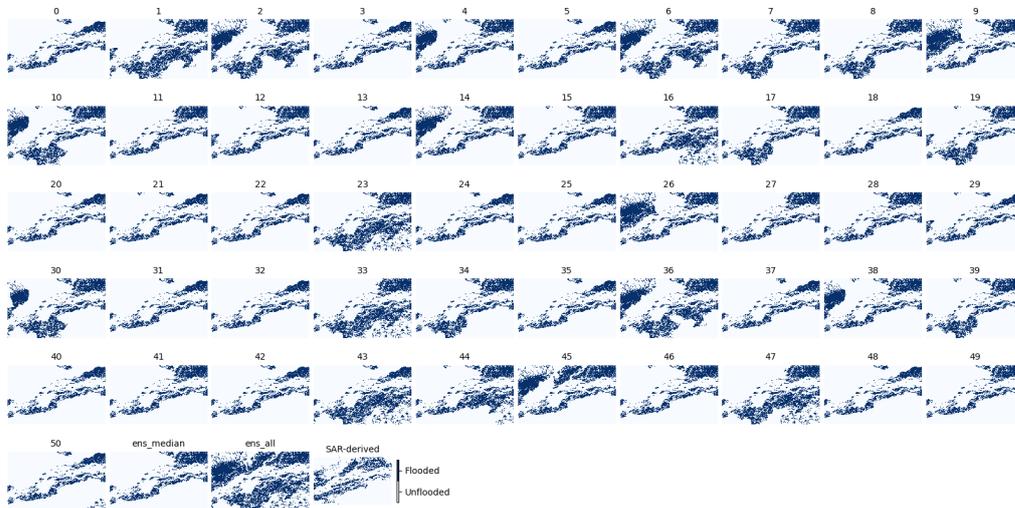


Figure 4.5: Brahmaputra River, Assam region, August 2017. 51 ensemble member forecast flood maps (labelled 0 to 50), ens_{median} and ens_{all} all at 1-day lead time and the Sentinel-1 SAR-derived flood map.

9 Figure 4.6 shows the amalgamated probabilistic ensemble forecast indicating the prob-
10 ability of flooding at each grid cell location. This was produced by vertically stacking each
11 ensemble member flood map and adding vertically the number of flooded cells at each
12 grid cell location across all ensemble members. The total is divided by 51 to calculate

1 the probability. Dark blue colours near to the central river channel indicate agreement
2 between all ensemble members and 100% forecast probability of flooding, lighter colours
3 to the north of the river indicate a low probability of flooding.

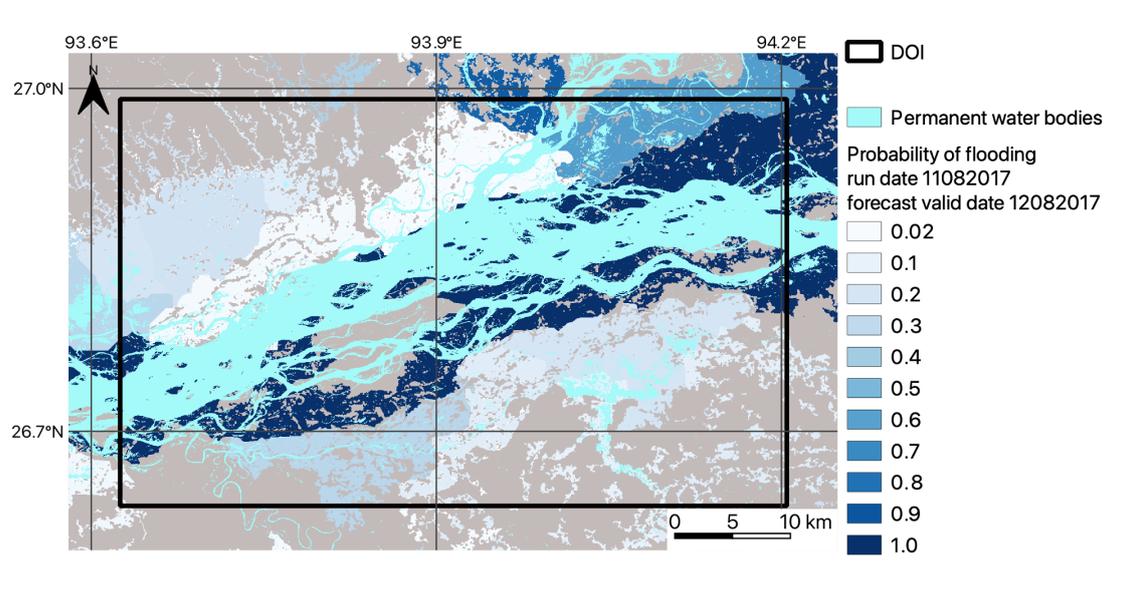


Figure 4.6: Brahmaputra River, Assam region, August 2017. Colour shading from white (low) to dark blue (high) indicate the forecast probability of flooding based on a 1-day lead time, 51 ensemble member flood map forecast for the Brahmaputra River in the Assam region, August 2017. (Note map background is grey)

4 4.5 Results and discussion

5 To demonstrate an application of the spatial scale approach to both ensemble forecast
6 flood map evaluation of forecast skill and the spatial spread-skill relationship, we apply
7 the methods outlined in Section 4.3 to the flooding case described in Section 4.4.1. First,
8 in Section 4.5.1 we verify the full ensemble using a spatial scale approach to calculate a
9 skilful scale of agreement between each ensemble member and the SAR-derived flood map
10 (Fig. 4.2) along with the combined ensemble (ens_{all}) and the ensemble spatial median
11 (ens_{median}). We evaluate the location specific spatial skill of the ensemble by calculating

1 categorical scale maps (Section 4.5.2) for ens_{all} , ens_{median} and a best and worst case
2 ensemble members determined by the skilful scale calculated in Section 4.5.1. In Section
3 4.5.3 we evaluate the spatial predictability of the full ensemble and show this on the
4 *Spatial Spread-Skill* (SSS) map, indicating regions where the ensemble is over-, under- or
5 well-spread.

6 4.5.1 Ensemble spatial scale evaluation

7 Here we investigate how a scale-selective approach can be useful for extracting meaning-
8 ful information from a flood map ensemble forecast where multiple forecast flood maps
9 represent equally likely flooding scenarios (Fig. 4.5). A minimum skilful scale (where
10 $FSS > FSS_T$) has been calculated for each individual member flood map, ens_{all} and
11 ens_{median} . The results in Figure 4.7 show that individual ensemble member spatial skill
12 varies considerably with FSS at grid level ranging from 0.35 to 0.59. One member ens_1 ,
13 which would usually be disregarded as an outlier due to its low probability, outperformed
14 all other members significantly with a skilful scale achieved at a neighbourhood size of
15 $n = 3$. The combined ens_{all} showed more skill at grid level ($n = 1$) and smaller neigh-
16 bourhood sizes compared with ens_{median} , both however exceeded FSS_T at $n = 41$, or 615
17 m. At neighbourhood sizes greater than $n = 41$, ens_{median} outperformed ens_{all} . There is a
18 cluster of members showing similar skill to ens_{median} and ens_{all} and a second cluster, with
19 more ensemble variation but indicating lower skill than the first cluster. The ens_{median} and
20 ens_{all} flood maps outperform the second cluster, however there are individual members
21 with a higher spatial skill score compared to ens_{median} and ens_{all} . These results show that
22 all ensemble member flood maps, including outliers, should be considered individually as
23 possible future flooding scenarios. Spatial variations across individual ensemble members
24 (see Fig. 4.5 ens_1 compared to ens_{median}) indicate that it is not meaningful to consider
25 only the ensemble median flood map to represent the information within the full ensemble.

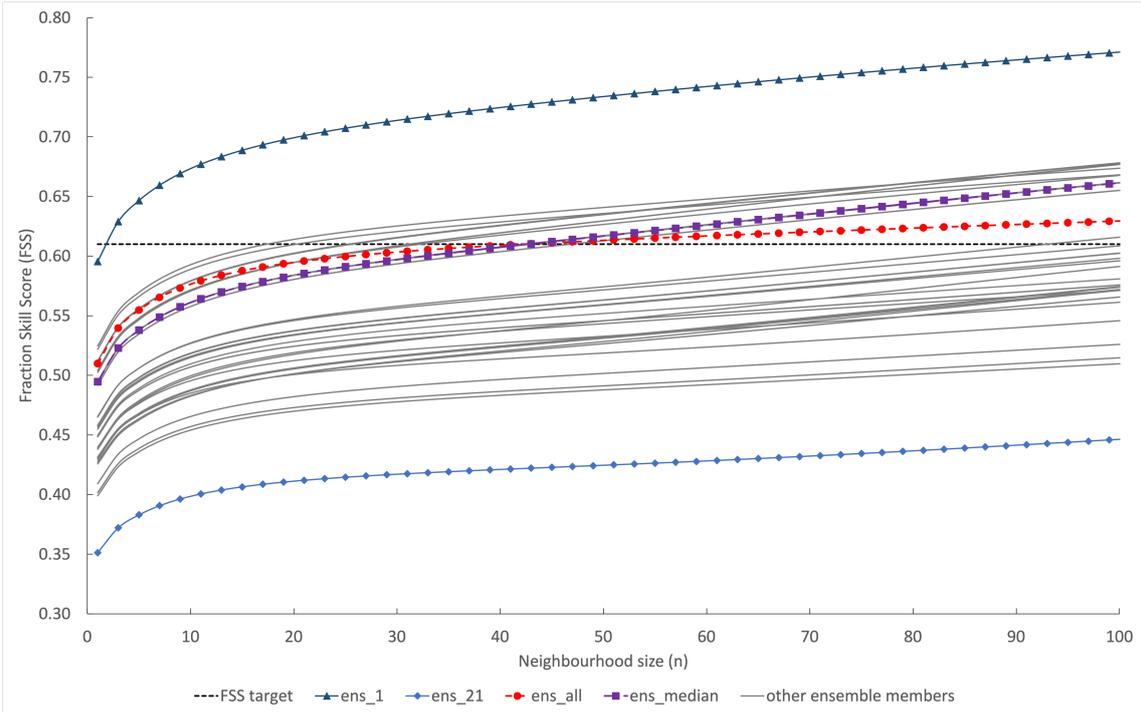


Figure 4.7: The spatial skill of each individual ensemble member forecast flood extent is evaluated along with the ens_{median} (a spatial median where 26 or more members predict flooding at a grid cell location) and ens_{all} (flooded grid cells from all ensemble members are combined). The FSS is calculated at increasing neighbourhood sizes to determine the scale at which the forecast becomes skilful at capturing the observed flood (FSS_T).

1 4.5.2 Ensemble spatial predictability

2 The scale-selective skill scores calculated for different aspects of the ensemble forecast
3 give a domain-averaged score and skilful scale. To understand location specific spatial
4 predictability of the ensemble forecast, categorical scale maps are calculated and presented
5 in Figure 4.8. These show how the agreement scale (Section 4.3.2) varies with location for
6 (a) ens_{all} , (b) ens_{median} , (c) ens_1 , the ‘best’ performing ensemble member and (d) ens_{21} ,
7 the ‘worst’ performing ensemble member. The ensemble summary map, ens_{all} (Fig. 4.8
8 (a)) captures most of the observed flooding (in grey) with small regions of under-prediction
9 (red). However, as you might expect to see by including every potential flooding realisation
10 there are significant regions of over-prediction (blue) or false alarm. The region of over-

1 prediction to the south of the river is less evident in the ens_{median} categorical scale map
2 (Fig. 4.8 (b)) which performs worse to the north by under-predicting flooding here. This
3 flooding is captured well by ens_1 (Fig. 4.8 (c)) and in particular close to a confluence zone
4 where the Subansiri River joins the Brahmaputra (grid cell location (1100, 250)). This
5 ties in with the high rainfall totals accumulated just to the north of this region associated
6 with localised very heavy rainfall (Floodlist, 2017). A region of under-prediction at grid
7 cell location (750, 750) is missed by all members. In future work, a closer inspection
8 of the DTM or surface features included/excluded in the hydraulic modelling, such as
9 embankment heights, may indicate how this modelling could be improved. The ‘worst’
10 performing ensemble member ens_{21} (Fig. 4.8 (d)) accurately predicts flooding closer to
11 the river channel, however under-prediction to the north along with over-prediction to the
12 south show where the forecast was inaccurate. Categorical scale maps enable different
13 ensemble flood map presentations to be evaluated so that the most useful presentation
14 method can be determined for a particular flooding situation.

15 4.5.3 Ensemble spatial spread-skill

16 To evaluate the location specific skill of the full ensemble, one option would be to calculate
17 51 categorical scale maps from each individual member flood map (Fig. 4.5). This ap-
18 proach maintains the spatial detail held within each of the ensemble member flood maps,
19 although does require multiple visual comparisons to be made by the flood forecaster or
20 modeller, which takes time and effort. Making comparisons across the different ensemble
21 member flood maps in Figure 4.5 provides a demonstration of these forecasting difficulties.
22 Further, the categorical scale maps do not evaluate the ensemble spatial spread. To ad-
23 dress this, we develop a Spatial Spread-Skill (SSS) map (derived from Fig. 4.9, presented
24 in Fig. 4.10) showing the spread-skill of the full ensemble forecast and keeping the loca-
25 tion specific detail. All ensemble members are included in this analysis which evaluates
26 both the spatial skill and the ensemble spatial spread of the forecast against the remotely

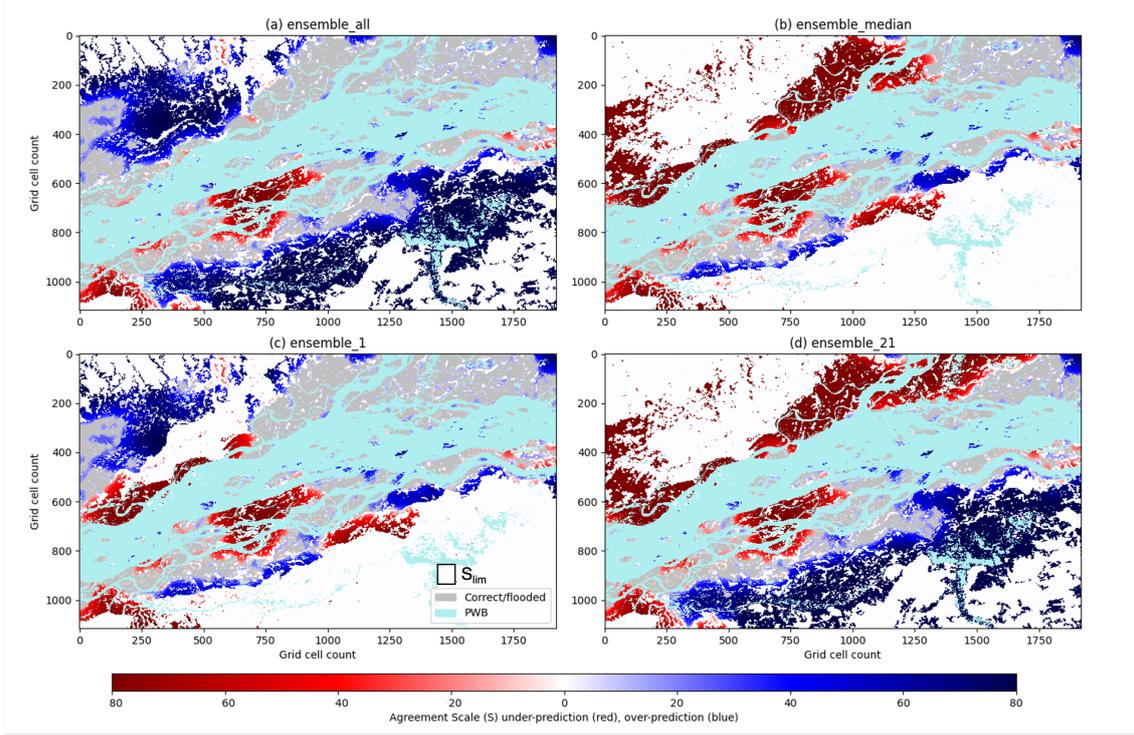


Figure 4.8: Brahmaputra River, Assam region, August 2017. Categorical scale maps for (a) ens_{all} (flooded grid cells from all ensemble members are combined), (b) ens_{median} (a spatial median where 26 or more members predict flooding at a grid cell location), (c) individual ensemble member 1 and (d) individual ensemble member 21. Red areas indicate where the forecast is under-predicted and blue regions represent over-prediction. The colour shade gives the scale of agreement (Eq. (7)) between the forecast and the observed flooding with lighter shading indicating a smaller agreement scale is required to reach the agreement criterion (Eq. (6)), a fixed maximum scale S_{lim} is drawn to scale (c). For georeferencing see Figure 4.6, each grid cell is 30 m x 30 m.

1 observed flood extent.

2

3 Figure 4.9 shows how the average-ensemble/ensemble-agreement scale in (a) $S_{ij}^{A(\overline{mm})}$
 4 calculated at each grid cell (representing ensemble *spread*) compares with the average en-
 5 semble/observed scale in (b) $S_{ij}^{A(\overline{mo})}$ (representing ensemble *skill*) along with the hexbin
 6 scatter plot in (c) which compares (a) and (b) to indicate the spatial spread-skill of the
 7 forecast. The hexagonal tessellation is used so that the distances along the hexbin di-
 8 agonal are on the same scale as those along the x and y-axis. For a perfect ensemble

1 forecast the average agreement scale between ensemble members should match the agree-
2 ment scale between the ensemble forecast and observed flood map, i.e. they should align
3 along the 1:1 line. The SSS map plots the difference between the ensemble/ensemble and
4 the ensemble/observed average agreement scales at each grid cell (Fig. 4.10) and indicates
5 where the spatial spread-skill is over-, under-, or well-spread. Three numbered areas (Fig
6 4.9(a)) identify three different ensemble spread-skill relationships. Area 1 shows that the
7 agreement between ensemble members is close, but that they disagree with the observed
8 flood extent. This is displayed in orange shades as an under-spread or miss region on the
9 SSS map, Figure 4.10. This is the region close to the confluence area described in Section
10 4.5.2. Recall that in this region, most ensemble members did not predict the flooding that
11 occurred with the exception of one ensemble member (ens_1). In area 2 on Figure 4.9, both
12 (a) and (b) are in agreement at grid level, which indicates the ensemble is well-spread;
13 these are shown in white on Figure 4.10. Away from the miss and well-spread regions in
14 Figure 4.9, the overall visual impression is that the ensemble spread-skill lies below the 1:1
15 line and is over-spread, indicated by area 3. This corresponds to purple shading on the SSS
16 map (Fig. 4.10). Overall Figure 4.9 tells us that the spread-skill relationship for this exam-
17 ple case study is not uniform across the domain but is in fact location specific. The areas
18 identified (1, 2 and 3) lie within different sub-catchments, which are linked to different
19 GloFAS grid cells, driving the ensemble flood map selection for each sub-catchment. Infer-
20 ences can be made about the spread-skill of the driving discharge data at sub-catchment
21 level across the domain. Using the spatial spread-skill relationship shown on the ensemble
22 SSS map we can infer how well the ensemble forecasting system encompasses the multiple
23 sources of uncertainty and how meaningful the probabilistic ensemble forecast of flood in-
24 undation actually is. An ensemble flood map forecast that is well-spread suggests that the
25 probabilistic forecast is meaningful. The SSS map is a useful evaluation tool for validating
26 flood forecasts in un-gauged or partially gauged rivers. A simulation library approach, like
27 the Flood Foresight maps used here, relies on the accuracy of the return period thresholds

1 set, the (ensemble) forecast streamflow and the accuracy of the flood inundation map for a
 2 given streamflow. The forecast evaluation approaches presented here enable these system
 3 attributes to be evaluated even where observed streamflow is limited or erroneous. The
 4 SSS map summarises the whole ensemble, which makes it useful for forecasters attempting
 5 to convey uncertainty information to decision makers, highlighting regions where there is
 6 high/low confidence in the forecast.

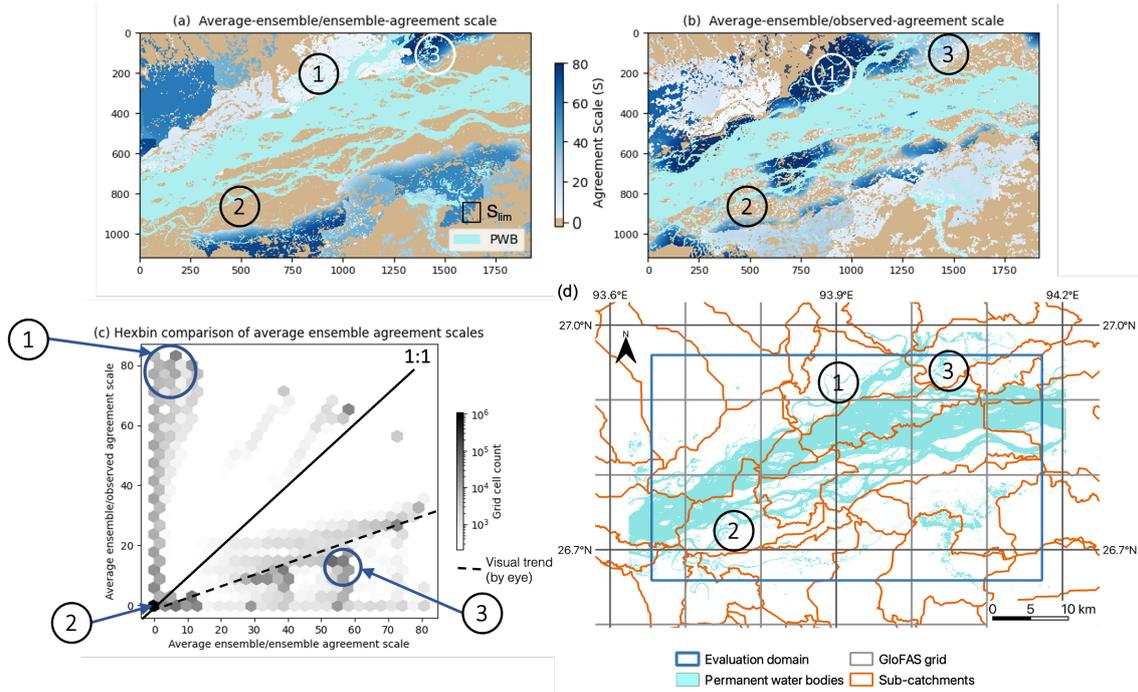


Figure 4.9: Brahmputra River, Assam region, August 2017. (a) The average agreement scale map of each unique pair of forecast ensemble flood maps and (b) between each ensemble member compared against the observed SAR-derived flood map. (c) A binned histogram scatter plot compares (a) and (b) to indicate the spatial spread-skill of the forecast ensemble. (d) indicates the corresponding sub-catchment locations. Areas labelled (1, 2 and 3) are discussed in Section 4.5.3. A fixed maximum scale S_{lim} (Eq. (6)) is drawn to scale (a). Note PWB means permanent water bodies.

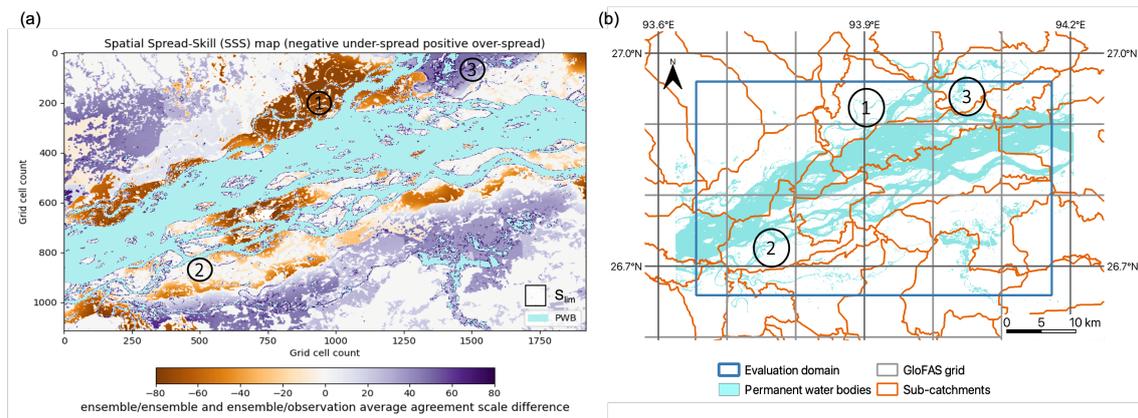


Figure 4.10: Brahmaputra River, Assam region, August 2017. (a) The Spatial Spread-Skill (SSS) map shows the difference between the ensemble/ensemble and the ensemble/observed average agreement scales at each grid cell. Negative values (orange) indicate where the ensemble is under-spread and positive values (purple) indicate where the ensemble is over-spread. White areas indicate where the average agreement scales match and indicate good spatial spread-skill. (d) Indicates the corresponding sub-catchment locations. Areas labelled (1, 2 and 3) are discussed in Section 4.5.3. A fixed maximum scale S_{lim} (Eq. (6)) is drawn to scale (a). Note PWB means permanent water bodies.

1 4.6 Conclusions

2 Differences between ensemble members in ensemble forecast flood map systems are mostly
3 driven by initial condition perturbations at the top of the hydro-meteorological forecast
4 chain within the numerical weather prediction system. Presently, there is limited under-
5 standing or evaluation of how these meteorological uncertainties link to mapped flooding
6 predictability, which involves additional sources of uncertainty. An evaluation of the spa-
7 tial predictability and the spread-skill relationship of the ensemble flood map forecast pro-
8 vides an improved understanding of the performance of the forecast system. Uncertainties
9 in other parts of the forecast chain are not truly represented by the ensemble flood maps
10 and evaluating the spatial spread-skill of the flood maps is important for understanding
11 the likelihood of flooding that the ensemble flood maps capture. In this paper, we present
12 a new scale-selective approach to assess the spatial predictability and spread-skill of an
13 ensemble flood map forecast by comparing against a satellite SAR-derived observation of
14 flood extent. By calculating a skilful scale at each grid cell for every unique ensemble
15 member pair we can determine the ensemble *spatial spread*, and between every ensemble
16 member and the SAR-derived flood map we can determine the ensemble *spatial skill*. The
17 hexbin scatter plot summarises the spread-skill relationship so that a trend across the
18 whole domain can be assessed. The difference between these skilful scales can be mapped
19 onto the Spatial Spread Skill (SSS) map which shows for each specific location in the
20 domain whether the ensemble is over-, under- or well-spread. The methods are applied
21 to an example flooding event of the Brahmaputra River in the Assam region of India in
22 August 2017.

23

24 In operational practice there are multiple options of ensemble flood map presentation
25 type such as presenting the ensemble median or other exceedance probability for delivery
26 to end-users and decision makers. An important aspect of developing an inundation flood

1 forecasting system is to determine the most useful way to present a spatial ensemble fore-
2 cast. Using a scale-selective approach we have evaluated the performance of individual
3 ensemble members, a combined total ensemble and the spatial ensemble median compared
4 to a SAR-derived observation of flood extent. Other options could be to exclude ensem-
5 ble member outliers, to spatially cluster similar ensemble members into groups of flood
6 extent or to present a most likely, best and worst case ensemble flood map. Whichever
7 presentation method is chosen, this should be fully explored using the spatial spread-
8 skill methods described here to evaluate the ensemble performance of historical flooding
9 events. We found for this example flooding event that one ensemble member significantly
10 outperformed the combined and median flood maps and that potentially in some flood
11 forecasting scenarios this member would have been excluded as an outlier. The categor-
12 ical scale maps show the ensemble spatial median could miss vital flooding information
13 and that all members should be considered as potential future flooding scenarios.

14

15 Through mapping the spatial-spread skill relationship, which varies with location, links
16 can be made between the spatial variations in spread-skill and the physical characteristics
17 of the flooding event. We found that one ensemble member outperformed all others in a
18 region close to a confluence zone and nearby observed heavy rainfall. The region correlates
19 to an area of under-spread ensemble members indicating that not enough members were
20 predicting flooding here. Future studies could investigate the physical processes further
21 using the methods presented here. The ensemble flood map spatial spread-skill could be in-
22 vestigated in the context of a particular physical process (such as rainfall intensity/location
23 or an improved aspect of the hydrological model such as antecedent soil moisture) and how
24 these uncertainties translate to the probabilistic flood map forecast. An understanding of
25 the spatial predictability is particularly important for un-gauged catchments where the
26 calibration of both forecast streamflow and return period thresholds (used to select the
27 simulation library flood map) are rarely practiced routinely. Ideally, in operational prac-

1 tice, these spatial verification approaches including the categorical scale and SSS maps
2 could be calculated and stored routinely as flooding events coincide with SAR-derived or
3 other remotely observed flood maps to build up a verification catalogue/database. This
4 database could then be used to investigate the spatial spread-skill model performance un-
5 der different scenarios such as forecast lead time, month or season, or flood type. More
6 locally, the impact of an improved DTM or the inclusion of a Digital Surface Model (DSM)
7 or other surface features in the hydraulic model such as embankments could be considered.
8 Over time, such a database would improve our understanding of the spatial predictability
9 of an ensemble flood map system and how well the uncertainties present are represented
10 by the ensemble forecast.

11 **4.7 Chapter summary**

12 In this chapter, we extend the scale-selective verification methods described in Chapter
13 3 with a new application to an ensemble flood map forecast. By considering the spatial
14 spread between ensemble member flood maps, and between each ensemble member and a
15 SAR-derived flood map, we evaluate the spread-skill of an ensemble flood map forecast.
16 The spatial uncertainty held within multiple ensemble flood maps can be summarised onto
17 a single map indicating the forecast performance at specific locations. For the case study
18 evaluated, we found that individual ensemble member's forecasts can be more accurate
19 than an aggregated forecast such as the ensemble median. This is an important con-
20 sideration when choosing how to present ensemble flood map information in operational
21 practice.

1 Chapter 5

2 A multi-system comparison of 3 forecast flood extent using a 4 scale-selective approach

5 In this chapter we address the third research question outlined in Chapter 1; How can a
6 scale-selective approach be applied to evaluate multiple flood forecasting systems?:

- 7 • How can we evaluate the performance of flood forecasting systems predicting flood
8 inundation extent at different spatial scales?
- 9 • What can we learn about the flood forecasting system performance and how does
10 each compare?

11 The remainder of this chapter (except for the chapter summary, Section 5.9), has been
12 published and is reproduced from (Hooker et al., 2023b).

1 **5.1 Abstract**

2 Fluvial flood forecasting systems increasingly couple river discharge to a flood map library
3 or a real-time hydrodynamic model to provide forecast flood maps to disaster management
4 teams and humanitarian agencies. The forecast flood maps can be linked to potential im-
5 pacts to inform disaster risk reduction schemes, such as forecast-based financing. The
6 success of forecast-based financing is dependent upon the accuracy of the forecast flood
7 maps. We investigate a new application of scale-selective verification by evaluating the
8 performance of three flood forecasting systems. Two simulation library systems, Flood
9 Foresight (30 m) and GloFAS Rapid Flood Mapping (1000 m) and one hydrodynamically
10 modelled system, the Bangladesh FFWC Super Model (300 m), all made predictions of
11 flood extent at different spatial scales (grid lengths, shown in brackets) for the Jamuna
12 River flood, Bangladesh, July 2020. These forecast flood maps are validated against
13 Synthetic Aperture Radar-derived observations of flooding across four districts using a
14 scale-selective approach that can compare directly across different spatial scales. Our re-
15 sults show that the simulation library system accuracy critically depends on the discharge
16 return period threshold set to trigger a flood map selection and the number of hydrological
17 model ensemble members that must exceed it. At short forecast lead times, the Super
18 Model outperforms the other systems in three districts. Near to the Bangladesh border,
19 the trans-boundary benefits of the two global systems are evident, with both outperform-
20 ing the local model. We conclude that a scale-selective verification approach can quantify
21 the skill of systems operating at different spatial scales so that the benefits and limitations
22 can be evaluated. Multi-system comparison of flood maps is important for improving
23 impact-based forecasts and ensuring funds and response activities are appropriately tar-
24 geted.

25

1 **HIGHLIGHTS:**

- 2 • Through a new application of our scale-selective validation method we compare
3 flood maps of different spatial scales (grid lengths) with SAR-derived observations
4 of flooding.
- 5 • Three flood forecasting systems (two global ensemble and one local deterministic)
6 operating during flooding in Bangladesh, July 2020 are evaluated.
- 7 • The return period threshold set and the number of ensemble members used to trigger
8 a flood map from a simulation library proves crucial to the system performance.
- 9 • The results show the importance of accounting for spatial scale when interpreting
10 skill scores in multi-system studies.

11 **5.2 Introduction**

12 Flood forecasting systems are increasingly used to improve preparedness ahead of a ma-
13 jor flooding event (Stephens & Cloke, 2014a; Wu et al., 2020). One of the main action
14 points from the recent Global Assessment Report (GAR2022) on Disaster Risk Reduction
15 (DRR) is to ‘design systems to factor in how human minds make decisions about risk’
16 (UNDRR, 2022b). Whilst flood forecasting systems have improved significantly and con-
17 tinue to improve both globally and locally, the reliance on government departments and
18 disaster managers to make the right decisions when faced with a potential crisis can result
19 in inappropriate actions and unpreparedness (e.g., Fekete & Sandholz, 2021; Coughlan de
20 Perez et al., 2022). The GAR2022 report shows that just 5.8% (\$5.5 billion USD) of official
21 development assistance contributes to disaster prevention and preparedness compared to
22 90.1% (\$119.8 billion USD) for emergency response. Yet it has been demonstrated (for
23 Europe) that financing for mitigation purposes such as flood forecasting systems can lead

1 to overall cost savings (Pappenberger et al., 2015).

2

3 Forecast-based financing (FbF) schemes can form a major element of DRR strategies
4 and aim to directly link forecasts of extreme events to humanitarian actions (Coughlan
5 de Perez et al., 2015, 2016). FbF schemes work by quantifying risks in advance of crises
6 or disasters, prepositioning funds, and agreeing in advance how funds will be released
7 based on forecasts, ahead of an event (OCHA, 2020). Anticipation and risk financing
8 allows humanitarians to be better prepared by making important decisions before disaster
9 strikes. These proactive decisions, directly linked to action (rapid fund release) remove
10 the potential for reactive, incorrect decisions in the midst of a disaster. The success of the
11 FbF system largely depends on the threshold triggers set and on the performance of the
12 flood forecasting system at mapping the flood hazard.

13

14 Advances in flood forecasting systems both at global and local levels link together
15 meteorological and hydrological forecasts to hydrodynamic models, simulating flood-wave
16 propagation (Emerton et al., 2016; Wu et al., 2020; Apel et al., 2022). The resulting
17 flood maps when directly linked to impacts can be used to inform DRR schemes. Multiple
18 trade-offs exist in the development of such systems that inherently depend on observation
19 data availability and computing power. These determine whether the maps can be mod-
20 elled in real-time or are pre-calculated and form part of a simulation library; the spatial
21 scale (grid size) of the forecast flood maps and whether the maps are deterministic or
22 probabilistic (Savage et al., 2016). A recent review of flood inundation prediction (Bates,
23 2022) states that a key task to drive forward the development of better global hydraulic
24 models will be more rigorous and comprehensive validation. Hoch and Trigg (2019) out-
25 line a Global Flood Model Validation Framework which includes a recommendation to
26 routinely validate flood extent. Quantitative performance evaluation forms an important
27 part of fitness-for-purpose assessment and continual system improvement. Currently, there

1 is limited quantitative validation of operational flood forecasting systems producing flood
2 maps and operating in the same area. A recent advancement in flood map validation
3 (Hooker et al., 2022) means that quantitative comparisons can be made across flood maps
4 at different spatial scales, which makes a multi-system evaluation possible.

5
6 The accuracy of ensemble forecasts of flood extent can be verified by comparing with
7 observations of flooding from unmanned aerial vehicles or satellite-based sensors. Satellite-
8 based synthetic aperture radar (SAR) sensors are active, which means they can operate at
9 night, through cloud and weather, and are well known for their flood detection capability
10 (e.g., Horritt et al., 2001; Mason, Davenport, et al., 2012; Schumann et al., 2022). The
11 SAR backscatter intensity depends on the smoothness of the surface, with unobstructed
12 flooded areas returning low backscatter values. Recent techniques used to extract flood
13 extent from SAR images have led to improved flood detection in urban areas (Mason
14 et al., 2018, 2021a, 2021b). Since late 2021, SAR-derived flood maps are produced for
15 every Sentinel-1 image detecting flooding around the world by the Global Flood Monitor-
16 ing (GFM) service (EU Science Hub, 2021; GFM, 2021; Hostache et al., 2021), part of
17 the Copernicus Emergency Management Service (CEMS) (Copernicus Programme, 2021).
18 Within eight hours of the Sentinel-1 image acquisition, three flood detection algorithms
19 are combined to give the flood class and uncertainty estimation per grid cell.

20
21 In this paper, through application of a new approach, we evaluate the inundation ac-
22 curacy of three fluvial flood forecasting systems operating during severe flooding of the
23 Jamuna River in Bangladesh, July 2020. The flood maps are compared against satel-
24 lite SAR-derived flood extent. The systems are: Flood Foresight, a FbF system run by
25 JBA Consulting working in partnership with the Start Network (Revilla-Romero et al.,
26 2017); the Global Flood Awareness System (GloFAS) Rapid Flood Mapping (RFM) service
27 (Alfieri et al., 2013; GloFAS, 2021) and the Bangladesh Flood Forecasting and Warning

1 Centre (FFWC) Super Model (BWDB, 2020), each producing forecast flood maps at dif-
2 ferent spatial scales (30 m, 1000 m and 300 m respectively). We investigate how a novel
3 scale-selective spatial verification approach (Hooker et al., 2022) can be applied to multi-
4 system studies where the forecast flood maps are presented at different spatial scales.
5 Through applying this approach, we determine a skilful scale of each flood map that can
6 be directly compared and discuss the benefits and limitations of each forecast system for
7 flood mapping purposes that underpin the triggering of FbF schemes.

8

9 We describe the characteristics of the Jamuna River flood, July 2020 in Section 5.3.
10 The three flood forecasting systems are described in Section 5.4 along with details of the
11 SAR-derived observation of flooding. The scale-selective methods used to evaluate the
12 forecast flood maps are outlined in Section 5.5. The performance of Flood Foresight with
13 forecast lead time is presented in Section 5.6.1 and the multi-system flood map comparisons
14 are presented in Section 5.6.2. In Section 5.6.3 we discuss the benefits and limitations of
15 each system and conclude with recommendations in Section 5.7.

16 **5.3 Flood event on the Jamuna River, Bangladesh July 2020**

17 Bangladesh lies on the world’s largest delta that drains the Tibetan plateau and the Hi-
18 malayas to the north via the Brahmaputra and Ganges River systems. The river system
19 capacity in Bangladesh is overwhelmed each monsoon season by the volume of water trav-
20 elling through. Flooding is exacerbated in coastal regions where tidal surges from tropical
21 cyclones impede the river drainage (Bernard et al., 2022). Due to its geographical location
22 and the low lying, low slope nature of the land, Bangladesh is susceptible and vulnerable
23 to flooding year after year and faces a worsening situation as a result of climate change
24 and sea level rise (Hossain et al., 2021).

25

1 This study focuses on the Brahmaputra River (locally named the Jamuna River), which
2 is characterised by braided, meandering channels that migrate continually due to frequent
3 silting and erosion, particularly during flood events. The total length of the Brahmaputra
4 is 2,900 km with a catchment area of around 583,000 km². Several flashy tributaries such
5 as the Teesta join the main channel in the north from steep catchments in the southern
6 Himalayas. The main distributary of the Jamuna River is the Old Brahmaputra, described
7 by the Bangladesh Flood Forecasting and Warning Centre (FFWC) as a high flow spill
8 river contributing largely to flooding, depending on the variations of siltation at the entry
9 point (BWDB, 2020).

10

11 Bangladesh is divided locally into administrative districts, locally named zilas. Four
12 of these that align along the Jamuna River have been chosen for the comparison of flood
13 mapping skill, these are Kurigram, Gaibandha, Jamalpur and Sirajganj (Fig. 5.1).

14

15 Bangladesh experienced an active monsoon season during the summer of 2020 which
16 brought severe and prolonged flooding in multiple spells. An unusually wet May following
17 cyclone Amphan meant that water levels were already raised ahead of the monsoon season
18 (Hossain, 2020). According to the Bangladesh Water Development Board (BWDB) the
19 flooding had some remarkable characteristics. It began earlier than usual in late June
20 and had a triple peak that had never been seen before. The flooding affected 40% of the
21 country, inundating over 34,000 km². In 2020 this resulted in the second highest level of
22 flooding since 1989 and the second longest flood duration since 1998. An estimated 5.5
23 million people were affected with 1 million houses waterlogged. Around 1.1 million people
24 were displaced with almost 100,000 evacuated to over 1,500 shelters. Almost 1 million
25 tube-wells and more than 100,000 latrines were damaged, 83,000 hectares of paddy fields
26 were affected, and 257 people lost their lives.

27

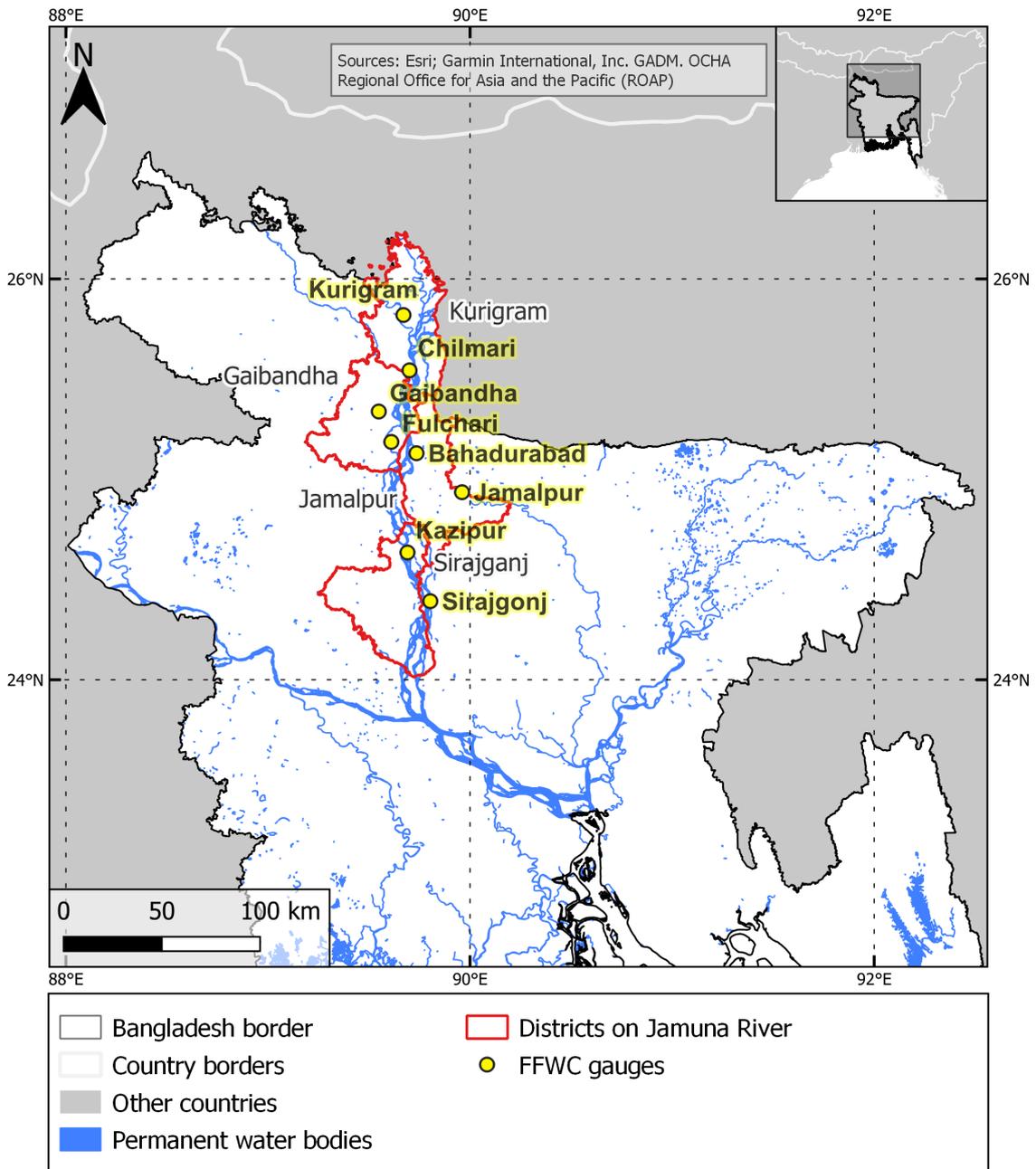


Figure 5.1: Four districts (zilas) of interest in the Jamuna catchment in northern Bangladesh.

1 The FFWC estimate that the Jamuna basin received 20% more rainfall in July than
 2 normal (BWDB, 2020). Gaibandha recorded the highest 1-day maximum rainfall across

1 the basin at 250 mm, with a 10-day consecutive maximum rainfall of 549.5 mm. The heavy
2 monsoon rainfall in July caused two flood peaks in one month, the first peak around 15
3 July and the second around 25 July.

4 **5.4 Flood forecasting systems and data**

5 The focus of this multi-system comparison is to evaluate the performance of three sys-
6 tems at forecasting the flood inundation extent for the second flood peak on 25 July. In
7 this section we briefly outline these three flood forecasting systems and summarise their
8 similarities and differences in simulating flood inundation extent (Sections 5.4.1, 5.4.2 and
9 5.4.3). Table 5.1 details the main system attributes. The flood map data used for com-
10 parison from each system is described in Section 5.4.4.

11

Table 5.1: Flood forecasting system comparison. Any ensemble member forecast discharge value is defined as ens_{any} . The mean forecast discharge value of all ensemble members is defined as ens_{mean} .

| Attribute | Flood Foresight | GloFAS Rapid Flood Mapping | FFWC |
|---|--|--|--|
| System application | Forecast-based financing for humanitarian early action | Global, broad scale medium range flooding prediction for large river basins | National flood forecasting and warning |
| System type | Ensemble simulation library (globally scalable) | Global ensemble simulation library (upstream drainage area >5000 km ² , river width >100 m) | Deterministic |
| Forecast type | Daily, 10-day lead time | Daily, maximum flood extent next 30 days | Daily, 5-day lead time |
| Meteorological model | ECMWF IFS | ECMWF IFS | BMD (WRF) |
| Hydrological model | LISFLOOD | LISFLOOD | MIKE II FF rainfall-runoff |
| Hydraulic model | RFLOW/JFLOW | CA2D | MIKE II GIS |
| Observed driving/input data | None | None | Rainfall and river water level |
| Grid length (m) | 30 | 1000 | 300 |
| DSM/DEM | NEXTMAP World30 DSM | SRTM (adjusted) | Survey of Bangladesh (>10 yrs old) |
| DSM/DEM grid resolution (m) | 30 | 90 (re-scaled to 1000) | 300 |
| Modelled flood map return period thresholds (yrs) | 20, 50, 100, 200, 500 and 1500 plus 30 interpolated flood maps | 10, 25, 50, 100, 250, 500 and 1000 | N/A |
| Flood map selection | $ens_{any} > 5yr$ return period threshold | $ens_{mean} > 10yr$ return period threshold | N/A |
| Defences included? | No | No | Yes |

1 **5.4.1 GloFAS Rapid Flood Mapping**

2 GloFAS couples state-of-the-art numerical weather prediction (NWP) forecasts with a
3 distributed hydrological model. With its continental scale set-up, it provides downstream
4 countries with forecasts of upstream river conditions up to one month ahead as well as
5 continental and global overviews for large river basins. As of version 3.1 (released in
6 May 2021), this modelling chain is based on the full configuration of the LISFLOOD hy-
7 drological model, forced by an ensemble of meteorological inputs (GloFAS, 2021). The
8 meteorological forecast data are provided to LISFLOOD by the ECMWF Integrated Fore-
9 casting System (IFS), the operational 51 ensemble member NWP system from ECMWF
10 (Alfieri et al., 2013; GloFAS, 2021). The NWP data (precipitation, temperature, potential
11 evapotranspiration, and evaporation rates for open water and bare soil surfaces) are taken
12 as inputs into the hydrological model, LISFLOOD. LISFLOOD is a distributed hydro-
13 logical rainfall-runoff model, simulating surface, groundwater and subsurface water flow
14 and then routing the water to river channels and simulating the routing of the channel
15 flow (LISFLOOD, 2022). LISFLOOD includes consideration of snow melt, infiltration,
16 vegetation interactions (interception, evapotranspiration, water uptake) and exchange of
17 soil moisture between a 3-layer soil water balance sub-model. The runoff data produced
18 is routed through a representation of the river network using a double kinematic wave ap-
19 proach. The river network used is taken from the HydroSHEDS dataset (Lehner, 2014b).
20 GloFAS is calibrated using historical streamflow records from selected stations worldwide
21 (Hirpa et al., 2018). For Bangladesh, four river gauges have been used to calibrate GloFAS
22 as reported on the GloFAS web viewer (GloFAS, 2022a). Two of these gauges are on river
23 reaches that would impact the four Jamuna River districts: one at Bahadurabad on the
24 main Jamuna River, and another upstream at Kaunia on the Teesta River, a tributary
25 of the Jamuna River. The observation record at Kaunia is very short at around 7 years
26 (1985-1992), while the record at Bahadurabad is over 35 years (1981-2015). Modified

1 Kling-Gupta Efficiency (KGE) is calculated for each station; a performance measure that
2 indicates how well the model reanalysis (at day 0) replicates the flows observed, greater
3 than 0.8 is very good and less than 0.2 is very poor. Both stations have KGE values above
4 0.7, indicating relatively good hydrological model performance.

5

6 GloFAS Rapid Flood Mapping (RFM) (GloFAS, 2022b) displays a maximum flood
7 extent over the next 30 days by matching the return periods from the GloFAS streamflow
8 forecast to a catalogue of modelled inundation extents. Flood maps are triggered for
9 basins greater than 5000 km² and where the 10-year RP threshold is exceeded by the
10 ensemble mean. The RP flood maps available are listed in Table 5.1. The flood maps
11 were developed using the semi-inertial formulation of the CA2D hydraulic model which is a
12 reduced complexity model based on the cellular automata approach and the diffusive wave
13 equations, specifically designed to simulate flood inundation events involving wide areas
14 (Dottori & Todini, 2011). Dottori et al. (2016) describe the methods used to derive the
15 flood maps at specified return periods on a global scale using a vegetation corrected version
16 of the global DEM SRTM (Farr et al., 2007). The hydraulic modelling was performed at
17 1 km grid resolution.

18 **5.4.2 Flood Foresight**

19 Start Network (Start Network, 2022) is a charity and network of over 80 humanitarian
20 agencies and aims to develop locally led, early action by moving to a model of proac-
21 tive funding to alleviate crises before they happen. JBA Consulting, in partnership with
22 Start Network, have developed a Disaster Risk Financing (DRF) system for the Jamuna
23 River that links a fluvial probabilistic flood inundation forecasting system, Flood Fore-
24 sight (Revilla-Romero et al., 2017), to populations impacted by flooding (Fig. 5.2). The
25 DRF system quantifies the flood risk to the population for the purposes of setting trig-
26 ger threshold levels through a probabilistic global catastrophe risk model, FLY (Dunning,

1 2019).

2

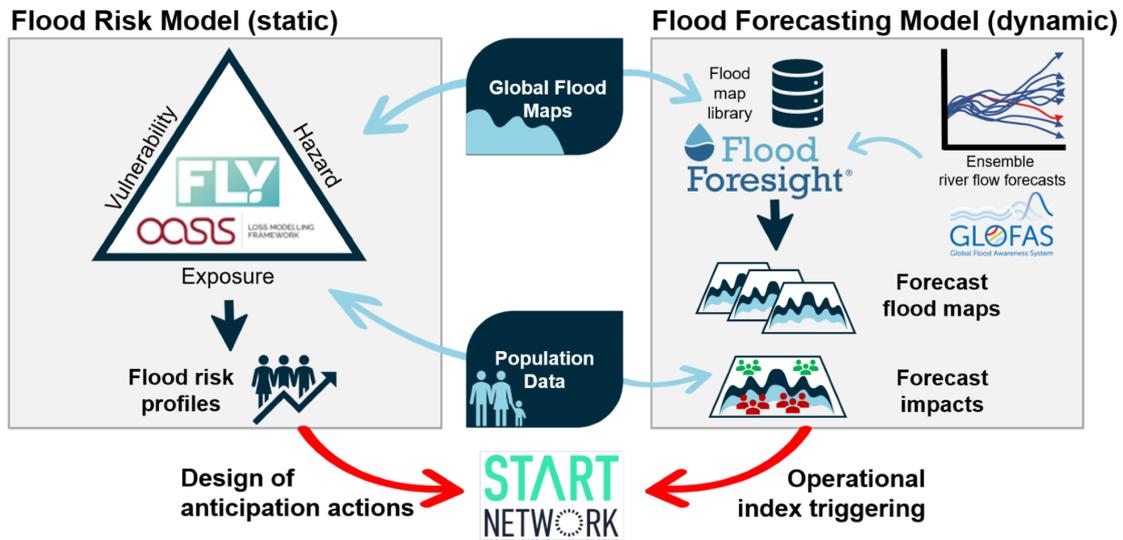


Figure 5.2: Flood Foresight/Start Network ensemble flood inundation forecast and population impacts work flow.

3 A domain of interest is divided into ‘Impact Zones’ (IZ) or sub-catchments using the
 4 HydroBASINS data-set (Lehner, 2014b). The Flood Foresight system links each IZ to
 5 GloFAS grid cells providing 51 ensemble member forecasts of river discharge (Section
 6 5.4.1). Based on the forecast discharge for each IZ, a precomputed flood map is selected
 7 from a simulation library. The flood map library was hydrodynamically modelled using
 8 JFlow® (Bradbrook, 2006) and RFlow using a detailed DSM at 30 m spatial scale at
 9 specified return period (RP) thresholds (detailed in Table 5.1). These were subsequently
 10 linearly interpolated at 5 intermediate intervals between each RP threshold and extrap-
 11 olated between zero and the 20 year RP flood map (totalling 36 flood maps). The flood
 12 map selected is determined by the RP threshold exceeded within each IZ. An example
 13 forecast domain in Figure 5.3 shows neighbouring IZ trigger flood maps at different RP
 14 thresholds and the RP threshold is not exceeded in every IZ.

15 The simulation library approach enables rapid flood map selection so that the system

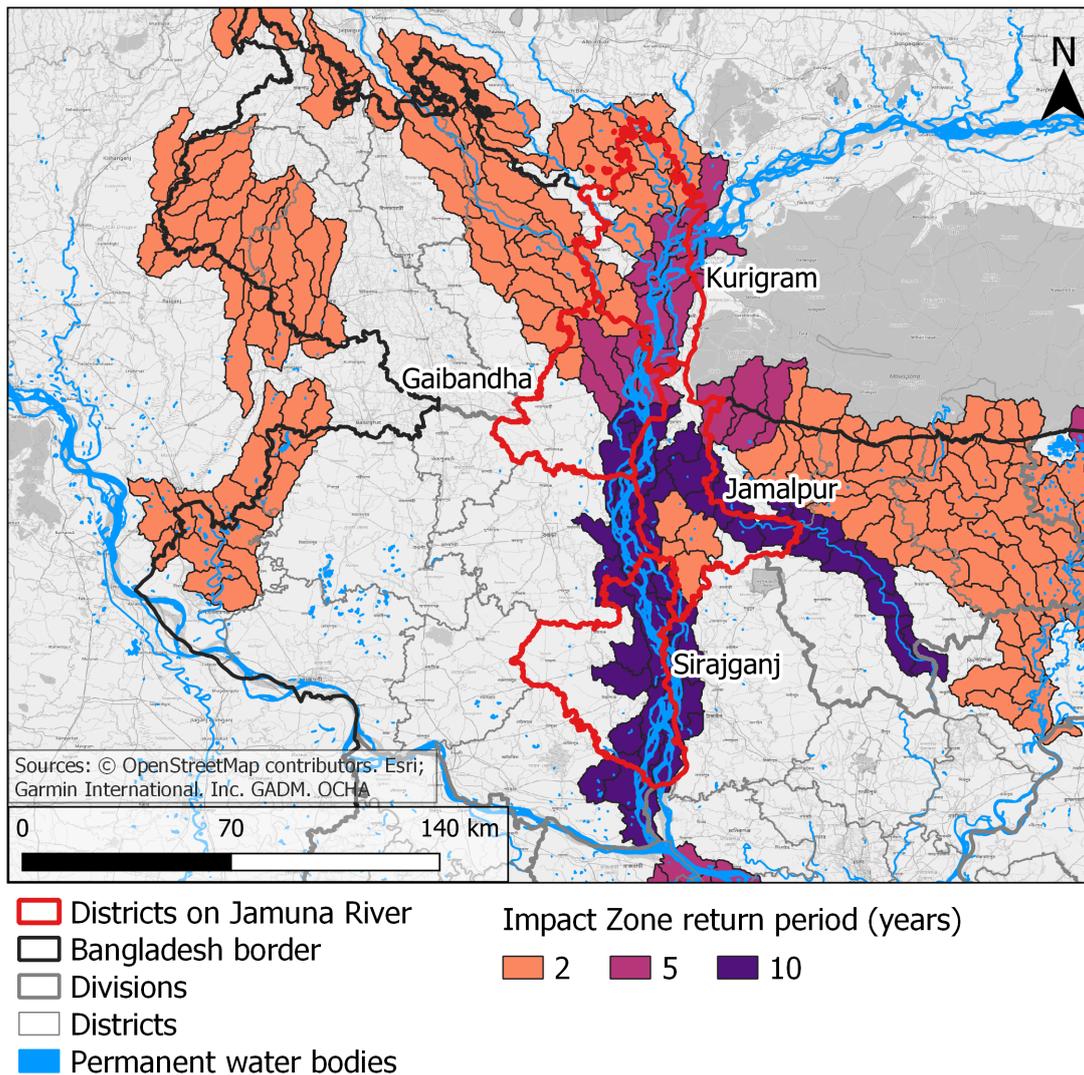


Figure 5.3: Example Flood Foresight forecast domain divided into Impact Zones (shown in colour with black outline where the 2-year return period threshold is exceeded), each linked to a GloFAS grid cell. The Impact Zone colour shows the corresponding return period threshold exceeded, determined by the GloFAS forecast discharge.

- 1 can be run in near real-time. Where the forecast discharge exceeds a 5-year RP threshold
- 2 the probabilistic flood maps are triggered and linked to populations impacted. IZ linked
- 3 to a 2-year RP threshold (Fig. 5.3) will not trigger a flood map due to low confidence in
- 4 the RP threshold levels and the uncertain flood map interpolation process at low discharge

1 values. The system runs daily and produces 51 ensemble member flood extent and depth
2 maps for forecast lead-times up to 10 days ahead. Hooker et al. (2023a) recommended
3 consideration of all ensemble members as indicators of potential flooding so we evaluate
4 each grid cell where any ensemble member indicates flooding (ens_{any}).

5 5.4.3 Bangladesh Flood Forecasting and Warning Centre

6 Historical catastrophic flooding in Bangladesh has led to well developed and forward-
7 thinking flood forecasting and warning services. Under the Ministry of Water Resources,
8 flood forecasting in Bangladesh is the responsibility of the Bangladesh Water Development
9 Board (BWDB) following the BWDB Act-2000. The Flood Forecasting and Warning Cen-
10 tre (FFWC), established in 1972, is the lead organisation for flood forecasting and warning
11 services. The FFWC act as coordinators between other Bangladesh agencies and ministries
12 involved in flood disaster management. During the event in July 2020, FFWC provided a
13 5-day deterministic and a 10-day probabilistic flood forecast. FFWC have identified two
14 main priority areas of improvement: to increase warning lead time and to make location
15 specific flood forecasts (BWDB, 2020). Operationally, FFWC use real time hydrologi-
16 cal data from water level and rainfall stations at 3-hourly intervals. Rainfall estimates
17 are based on the preceding three days of rainfall along with analysis derived from NWP
18 forecasts from the Bangladesh Meteorology Department (BMD, NCAR (2022)). Forecast
19 flood bulletins are prepared daily and disseminated through various modes to multiple
20 recipients.

21

22 During the monsoon flood season, the FFWC generate a daily hydrodynamically mod-
23 elled flood inundation map for the Jamuna, Ganges and Meghna river basins. The flood
24 maps are generated using output files from MIKE 11 FF Rainfall-Runoff hydrological
25 model and hydrodynamic modelling simulations using a customized MIKE 11 GIS model
26 (Havnø et al., 1995; Gourbesville, 1998). The Digital Elevation Model (DEM) used for the

1 hydrodynamic modelling has a 300 m spatial resolution that was collected by the Survey
2 of Bangladesh (SoB) more than 10 years previously.

3 **5.4.4 Observation data**

4 The data described here are used to evaluate the three flood forecasting systems. Obser-
5 vations of river water level from eight river gauges across the four districts were provided
6 by the FFWC for validation purposes. The gauge locations are shown on Figure 5.1 in
7 yellow, five are located on the main Jamuna River channel and three are located on trib-
8 utaries/distributaries of the Jamuna River. Three satellite SAR images from Sentinel-1
9 (S1A) acquired on 25 July 2020 and three pre-flood images from the same track from 7
10 June were used to derive a remotely observed flood map. The HASARD flood mapping
11 algorithm (Chini et al., 2017) hosted on WASDI (WASDI, 2022) uses a statistical, hierar-
12 chical split-based approach to separate the two classes (flooded and unflooded) using the
13 pre-flood and flood images. The HASARD mapping algorithm removes permanent water
14 bodies, such as the river channel, reservoirs and lakes. Flooded areas beneath vegetation,
15 near to buildings and under bridges will not be detected using this method. To smooth
16 the HASARD flood maps and allow a fairer comparison we apply a morphological closing
17 operation, without impacting the location of the flood extent, to flood fill buildings and
18 vegetation. So that the flood prediction accuracy alone can be evaluated, the pre-flood
19 occurrence of surface water using the JRC Global Surface Water database (Pekel et al.,
20 2016) has been removed from the forecast inundation maps. The observed flood extent
21 mosaic derived from the three SAR images at 20 m grid size was re-scaled to match the
22 relevant forecast flood map grid lengths using majority (mode) aggregation.

5.5 Scale-selective evaluation methods

The flood forecasting systems detailed in Table 5.1 produce flood maps at 30 m, 300 m and 1000 m spatial scales (grid lengths). Validation of forecast flood maps against remotely observed flood extent is typically carried out by labelling each grid cell using a contingency table with categories: correctly predicted flooded, under-prediction (miss), over-prediction (false alarm) and correctly predicted unflooded. After labelling, conventional binary performance measures such as Critical Success Index (CSI) and Pierce Skill Score (PSS) are calculated and give a domain average skill score (Stephens et al., 2014). Comparing the binary performance measures of these flood maps at different spatial scales would not be meaningful as the higher resolution maps would be overly penalised due to the double penalty impact (Roberts & Lean, 2008; Hooker et al., 2022). Several commonly applied binary performance measures will be included here for demonstration and comparison purposes only and the details of these can be found in Appendix 5.8, Table 5.2.

To tackle the issue of validating across differing spatial scales, we apply a scale-selective approach to flood map evaluation. The scale-selective evaluation approach includes calculation of the Fraction Skill Score (FSS, Roberts and Lean (2008)) and location specific agreement scales (Dey, Roberts, et al., 2016), which are plotted on a Categorical Scale Map (CSM, Hooker et al. (2022, 2023a)). A brief summary of the method is given here. For full methodology please see Roberts and Lean (2008); Dey, Roberts, et al. (2016); Hooker et al. (2022, 2023a). The FSS calculates the accuracy of a forecast flood map by comparing against an observed flood map across a range of neighbourhood lengths (n , see Figure 5.4). First, every grid cell is compared ($n = 1$). Then, the next largest neighbourhood size, $n = 3$, surrounding each grid cell is compared and the process continues to $n = 5$, $n = 7$ and so on. The fraction flooded (number of flooded grid cells in the neighbourhood divided by the total number of grid cells in the neighbourhood) in each of the forecast and

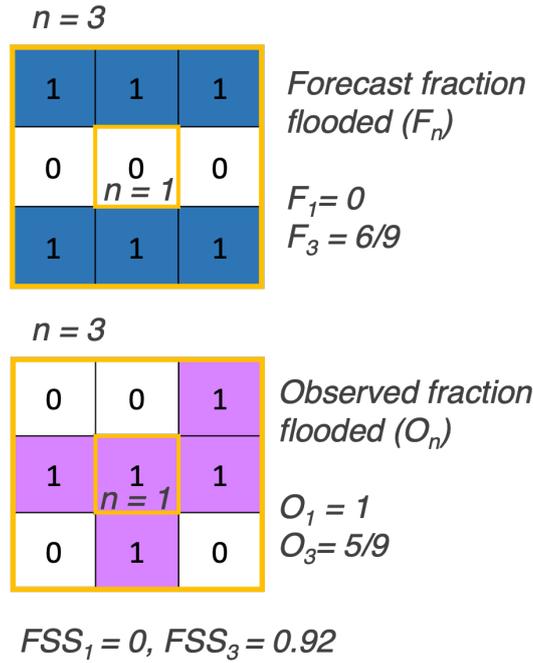


Figure 5.4: An example FSS calculation applied to forecast flood extent, 1 = flooded (in blue), 0 = unflooded (in white) compared to a remotely observed flood extent in pink. The FSS is calculated for two neighbourhood sizes, $n = 1$ (small gold box) and $n = 3$ (large gold box).

1 observed flood maps are used to calculate the FSS at each n . The FSS calculation is based
 2 on the Brier Skill Score. The FSS for each n is derived by calculating the mean squared
 3 error (MSE) between the forecast and observed fractions and dividing this by a reference
 4 MSE. This value is subtracted from 1 to give the FSS score. Increasingly larger neigh-
 5 bourhoods are compared until a target FSS score has been reached and exceeded at which
 6 point the skilful scale has been met (e.g. see Fig. 5.9). The target FSS score, given by
 7 $FSS_T \geq 0.5 + \frac{f_o}{2}$ depends on the fraction of observed flooding in the domain of interest, f_o .

8
 9 The FSS gives a domain averaged skill score. We also calculate a local agreement
 10 scale (S) for each grid cell that can be mapped onto a Categorical Scale Map (CSM).
 11 The relationship between S and n is given by $S = (n - 1)/2$. An acceptable level of

1 background bias between the forecast and observed flood maps can be pre-set. This is
2 used to determine an agreement criterion. Like the FSS method, the comparison begins at
3 each grid cell, if the agreement criterion is met, the grid cell is labelled with an agreement
4 scale $S = 0$. Where the criterion is not met, a larger neighbourhood size is compared
5 (e.g. $n = 3$). The fraction flooded in each of the forecast and observed flood maps are
6 compared and if the criterion is met the agreement scale assigned would be $S = 1$. The
7 process continues to larger neighbourhoods (e.g. $n = 5$, $S = 2$) until either the criterion is
8 met or a predetermined limit is reached (S_{lim} , set to 9 for this application). The agreement
9 scale at this limit would indicate a miss or false alarm for the grid cell. Combining a map
10 of agreement scales with a conventional contingency map creates a CSM which shows
11 the level of agreement (S) and whether the forecast is over or under-predicting the flood
12 extent (e.g. see Fig. 5.8). The CSM shows a location specific skill score that can be linked
13 to different aspects of the forecast system such as IZ and their associated river discharge
14 forecast or return period thresholds, river channel bathymetry, the DSM or flood defences.

15 **5.6 Results and discussion**

16 Forecast flood maps for the Jamuna River flooding, July 2020 from each of the three
17 systems described in Sections 5.4.2 - 5.4.3 are compared against SAR-derived flood maps
18 (described in Section 5.4.4) using scale selective evaluation methods (as discussed in Sec-
19 tion 5.5). First, in Section 5.6.1 the Flood Foresight system performance is evaluated
20 against forecast lead time, where flood maps out to 10-days lead time are available. Per-
21 formance evaluation of the other two systems with forecast lead time was not possible due
22 to the availability of flood maps. In Section 5.6.2 we compare the forecast flood maps from
23 each of the systems described in Section 5.4 and discuss their benefits and limitations in
24 Section 5.6.3.

1 **5.6.1 Flood Foresight Jamuna River case study**

2 The performance of the Flood Foresight system with forecast lead time is evaluated here.
 3 Flood Foresight predicts the flood extent for second flood peak on 25 July at all lead times
 4 (out to 10-days). To evaluate the spatial accuracy within each district with forecast lead
 5 time, the absolute agreement scale score has been averaged across each district (solid lines,
 Fig. 5.5). An average agreement scale of zero would indicate the forecast and observed

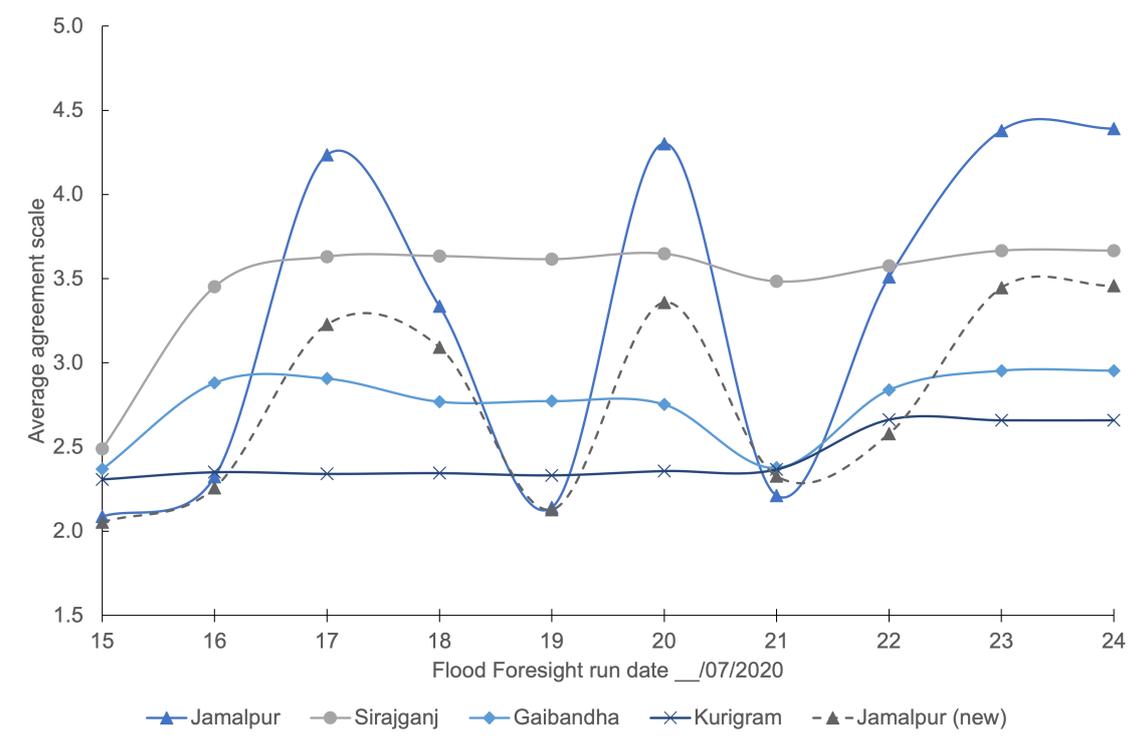


Figure 5.5: Flood Foresight average agreement scale against forecast lead time for each district and updated Jamalpur following reassociation of IZ with GloFAS grid cells, forecast valid for 25 July.

6
 7 fields are in agreement at grid level, $S = 2$ means agreement is reached within a 5 by 5
 8 neighbourhood. Across the 10-day forecast Flood Foresight performs best in Kurigram
 9 district in the north with consistently the worst performance in Sirajganj. Three of the
 10 districts (except Jamalpur) show a similar trend with forecast lead time with the best
 11 performance for all districts (smallest average agreement scale) occurring at a 10-day lead

1 time (2.31) with a second peak of performance at a 4-day lead time (2.61). The perfor-
2 mance worsens from a 4-day lead time to a 1-day lead time across all districts, with an
3 average agreement scale at a 1-day lead time of 3.42. The unusual variation in skill with
4 forecast lead time for Jamalpur prompted further investigation into the driving data.

5
6 Section 5.4.2 describes how the domain is divided into IZ with each of these linked
7 to the driving data, GloFAS river discharge (Fig. 5.3). The flood map selection within
8 each IZ depends on the forecast discharge exceeding a particular RP threshold. A direct
9 comparison of observed and modelled river conditions is not possible as observed data are
10 river water levels, GloFAS provides forecast discharge, and the stage-discharge relation-
11 ships are not available. The observed river water levels have been aligned using human
12 judgment to the nearest GloFAS grid cell forecast discharge (1-day lead time control mem-
13 ber) and compared for all river gauges. We aimed to match the trend in the two series
14 whilst keeping the local station risk level close to the GloFAS 2 and 5-year RP threshold
15 levels. Two of the eight river gauges are located in Jamalpur (Fig. 5.1). Bahadurabad
16 is located on the main Jamuna channel and Jamalpur is on a distributary of the Jamuna
17 River crossing the Jamalpur district (Fig. 5.6). (Hydrographs for gauges located outside
18 of Jamalpur are plotted in Appendix 5.8, Fig. 5.10.) Bahadurabad (Fig. 5.6 (a)) is well
19 calibrated in GloFAS (Section 5.4.1) and the overall forecast aligns well with the observed
20 river water level; the second flood peak magnitude is slightly underestimated. The dis-
21 charge forecast skill for Jamalpur station (Fig. 5.6 (b)) is very different to the observed
22 river water level in terms of variation in water levels/discharge, timing and magnitude of
23 flood peaks. A closer look at the river routing network in GloFAS shows that the Old
24 Brahmaputra distributary in Jamalpur is not connected to the main river channel; this
25 connection should occur at flood flows. To overcome the disconnection of the distributary
26 in the river network, the IZ association to GloFAS grid cells in the Flood Foresight system
27 design was manually updated. IZ aligning the Old Brahmaputra, crossing Jamalpur, were

1 reassociated to GloFAS grid cells located upstream on the main Jamuna channel. The new
2 forecast hydrograph for Jamalpur (Fig. 5.6 (c)) shows an improved forecast compared to
3 observed river water levels. However, the first flood peak arrives quicker than observed
4 and the second flood peak magnitude is underestimated.

5

6 Following the gauge reassociation update to the Flood Foresight system, updated CSMs
7 were reanalysed and this led to an improvement in the overall skill for the Jamalpur district
8 (Fig. 5.5) but with variation in skill with forecast lead time. The variation is investigated
9 by examining the CSMs for Jamalpur. Fig. 5.7(a) maps the original CSM for Jamalpur
10 for run date 20 July. CSMs show a location specific (at each grid cell) agreement scale
11 between the forecast flood map and the observed SAR-derived flood map. There is a
12 large area of under-prediction around Jamalpur caused by the disconnection of the Old
13 Brahmaputra distributary (as discussed above).

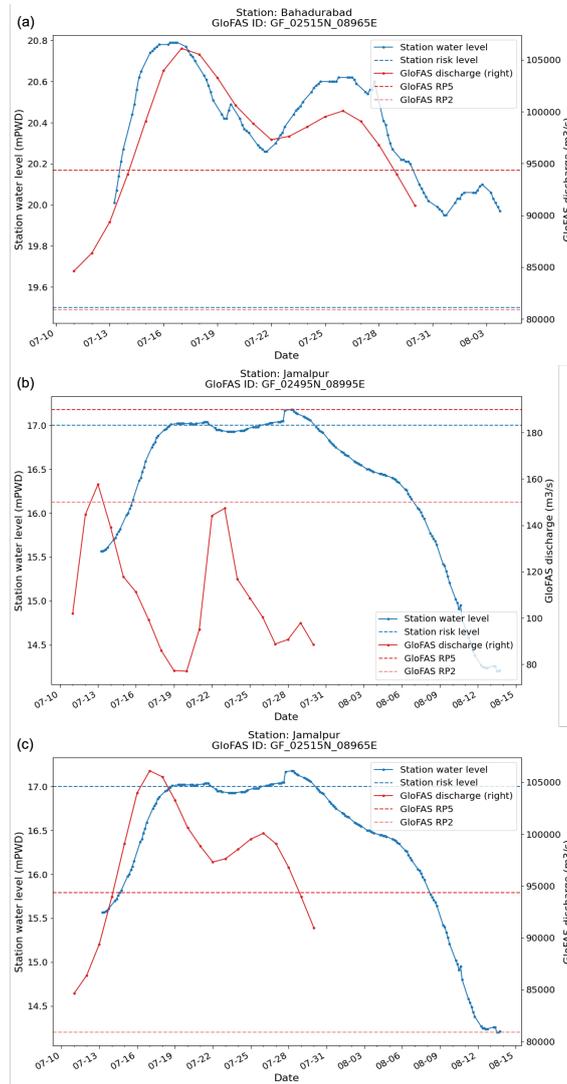


Figure 5.6: (a) GloFAS forecast discharge (control member, 1-day lead time) compared to FFWC observed river water level for the main Jamuna channel at Bahadurabad and (b) the old Brahmaputra distributary in the Jamalpur district. (c) The old Brahmaputra distributary forecast discharge following reassociation of IZ with GloFAS grid cells. The GloFAS RP threshold levels are taken from the nearest GloFAS grid cell to the gauge station location. Station risk levels are provided by FFWC.

1 The CSM change map (Fig. 5.7(b)) is calculated by taking the difference between the
 2 absolute CSM values, $|updated\ CSM| - |original\ CSM|$. The CSM change map for Ja-
 3 malpur shows where the reassociation has impacted the flood map. A negative agreement
 4 scale change indicates an increase in skill (purple, smaller grid size), whereas a positive
 5 agreement scale change indicates a decrease in skill (orange, larger grid size). The CSM
 6 change values can highlight areas where the agreement scale has increased/decreased but
 7 have not reached $S = 0$ (agreement at grid level). Areas surrounding the Old Brahmapu-
 8 tra (in purple) show most of the increased skill following the reassociation. However, there
 9 are regions of Jamalpur not impacted by the reassociation that remain under-predicted
 10 where smaller distributaries run, that would not be captured/calibrated in GloFAS due
 11 to their small size (Fig. 5.7(b)). This results in some variation in skill with forecast
 12 lead time remaining after the reassociation (Fig. 5.5). The visualisation of incremental
 13 improvements at specific locations will benefit future flood map development work. The
 14 CSM change map shows more sensitivity to changes in skill that would not be visible on
 15 a conventional contingency map.

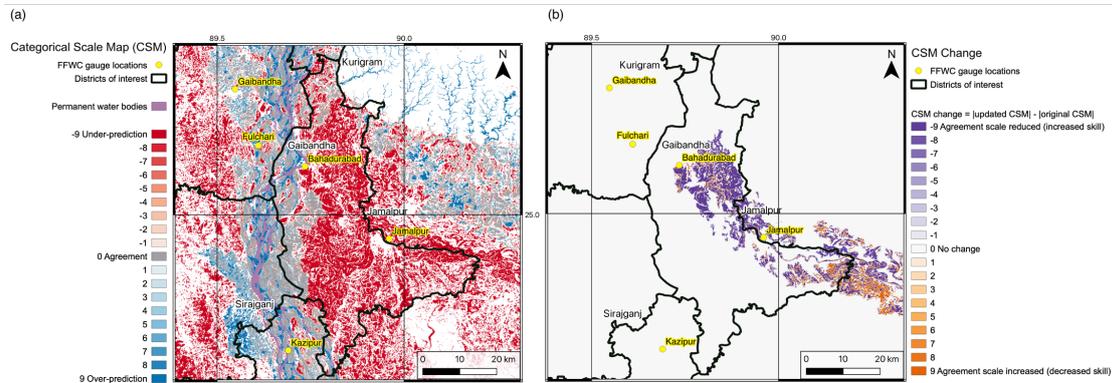


Figure 5.7: (a) Original CSM for Jamalpur. (b) CSM change (updated CSM - original CSM) for Jamalpur following reassociation of IZ with GloFAS grid cells. Run date 20 July, forecast valid for 25 July for (a) and (b).

1 **5.6.2 Multi-system flood map comparison**

2 GloFAS RFMs are displayed when the mean ensemble member exceeds the 10-year RP
3 threshold within the next 30 days. By inspecting the reporting point hydrograph at Ba-
4 hadurabad on each of the 10 days preceding the flood peak (25 July), we found that the
5 10-year RP threshold was exceeded just once by the ensemble mean on the 18 July (run
6 date), which means only one forecast flood map was available from GloFAS RFM for eval-
7 uation. FFWC provided daily flood maps, based on most recent observations, valid for
8 the same day. For the multi-system comparison, based on the availability of forecast flood
9 maps, the following have been selected to compare in detail: Flood Foresight's flood map
10 for the same run date as the available RFM forecast (18 July forecast date, valid date 25
11 July), the RFM (run date 18 July) and the FFWC flood map for 25 July (run date and
12 forecast valid date) have each been compared to the SAR-derived flood map (25 July) and
13 CSMs calculated (Fig. 5.8). CSMs have also been calculated for all available forecast lead
14 times for the Flood Foresight system (not shown).

15

16 The large area of flood under-prediction (in red) to the northwest of Sirajganj on each
17 CSM (Fig. 5.8) can be linked to observed very heavy rainfall and possible surface water
18 flooding (SWF) detected by the SAR data. The recorded rainfall described earlier (Section
19 5.3) is confirmed by rainfall derived from satellite data, (Climate Hazards Group InfraRed
20 Precipitation with Station (CHIRPS) data (Funk et al., 2015)), which shows that the July
21 rainfall anomaly in this area was 125 - 175 mm. The accumulated rainfall in GloFAS
22 for the ensemble mean for the 20 days preceding 25 July amounts to a maximum of 300
23 mm, with up to 400 mm north of Bahadurabad. This is an under-prediction compared to
24 observed rainfall (Gaibandha 1-day maximum rainfall 250 mm, 10-day consecutive max-
25 imum rainfall 549.5 mm), which explains the under-prediction seen in Fig. 5.8 (a) and
26 (b). Note that each system forecasts fluvial flooding and does not account for SWF, which

1 was likely given the extreme rainfall accumulations recorded. This is a limitation of both
 2 Flood Foresight and RFM systems and demonstrates the need for combining fluvial and
 3 SWF forecasting systems or combining the forecast flood maps with SAR data using data
 4 assimilation (e.g., García-Pintado et al., 2015; Cooper et al., 2019).

5

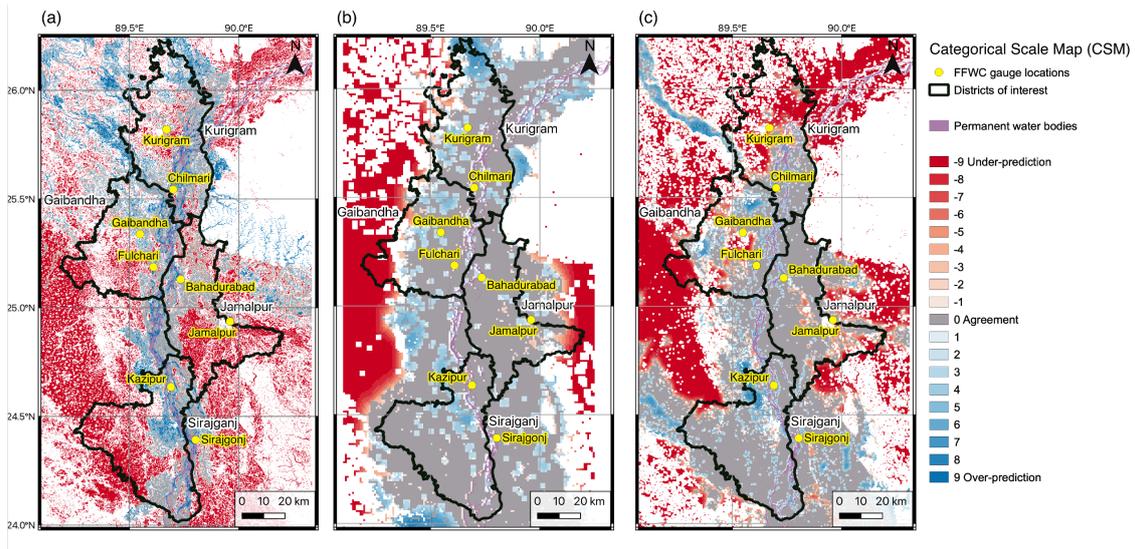


Figure 5.8: CSM for flood inundation forecasts from three forecast systems for flood peak 25 July compared to SAR-derived flooding. (a) Flood Foresight run date 18 July, (b) GloFAS RFM run date 18 July and (c) FFWC flood map run date 25 July.

6 The CSM for Flood Foresight (Fig. 5.8 (a)) shows areas of over-prediction (in blue)
 7 next to the Jamuna River. It is likely that more of the river channel has been removed from
 8 the SAR image during the flood mapping process compared to the 12-month occurrence
 9 of water in the Global Surface Water database. Also, the Jamuna River channel migration
 10 will also contribute to errors in this area as the DSM was acquired in 2016. To the west of
 11 Sirajganj, the FFWC correctly maps flooding associated with the Atrai River, a tributary
 12 of the Jamuna River that is not forecast by the other two systems. Multiple tributaries
 13 flow across Gaibandha and Sirajganj that are not currently resolved by GloFAS. This also
 14 impacts the performance of Flood Foresight when flood maps were not triggered above

1 the 5-year RP threshold.

2

3 Flood waters are (Fig. 5.8 (b)) spread further from the main river channel in the
4 RFM, compared to the Flood Foresight flood map. This is due to the GloFAS configu-
5 ration where clusters of cells are linked to the main channel reporting points. The RFM
6 extent is also due to a smoother DEM created by re-scaling from 90 m to 1000 m. This
7 will effectively remove flood barriers such as embankments and roads. An element of
8 smoothing (albeit to a lesser extent compared to the RFM) would also occur in the Flood
9 Foresight flood maps at 30 m grid length. The FFWC CSM shows good accuracy in
10 Sirajganj and Jamalpur (Fig. 5.8 (c)). In Sirajganj there is a region of over-prediction
11 on both the Flood Foresight and the RFM CSM not present on the FFWC CSM. The
12 FFWC model includes flood defences and water level observations from Kazipur, both
13 could contribute to the better performance seen here. FFWC maps perform less well in
14 Kurigram where the flood extent is under-estimated. Kurigram is next to the northern
15 border of Bangladesh where upstream water level data are unavailable for hydrological
16 model calibration and validation. The benefits of the trans-boundary systems of GloFAS
17 and Flood Foresight are evident here.

18

19 To quantify the district performance of each system the FSS has been calculated for
20 neighbourhood sizes up to $n = 19$ or larger (not plotted) where the FSS target has not
21 been reached (Fig. 5.9). The FSS gives a measure of spatial accuracy for each system
22 flood map, however the score is not directly comparable across different spatial scales since
23 the scores are calculated in terms of a neighbourhood size (Section 5.5). We can calculate
24 a skilful scale, which is half of the neighbourhood size at which the FSS exceeds FSS_T ,
25 this accounts for the size of the grid cell and can be directly compared across each of the
26 forecast systems. For example, in Figure 5.9(a) for Kurigram, at grid level ($n = 1$) the
27 FFWC FSS (0.52) exceeds the Flood Foresight FSS (0.48). We saw on the CSM map (Fig.

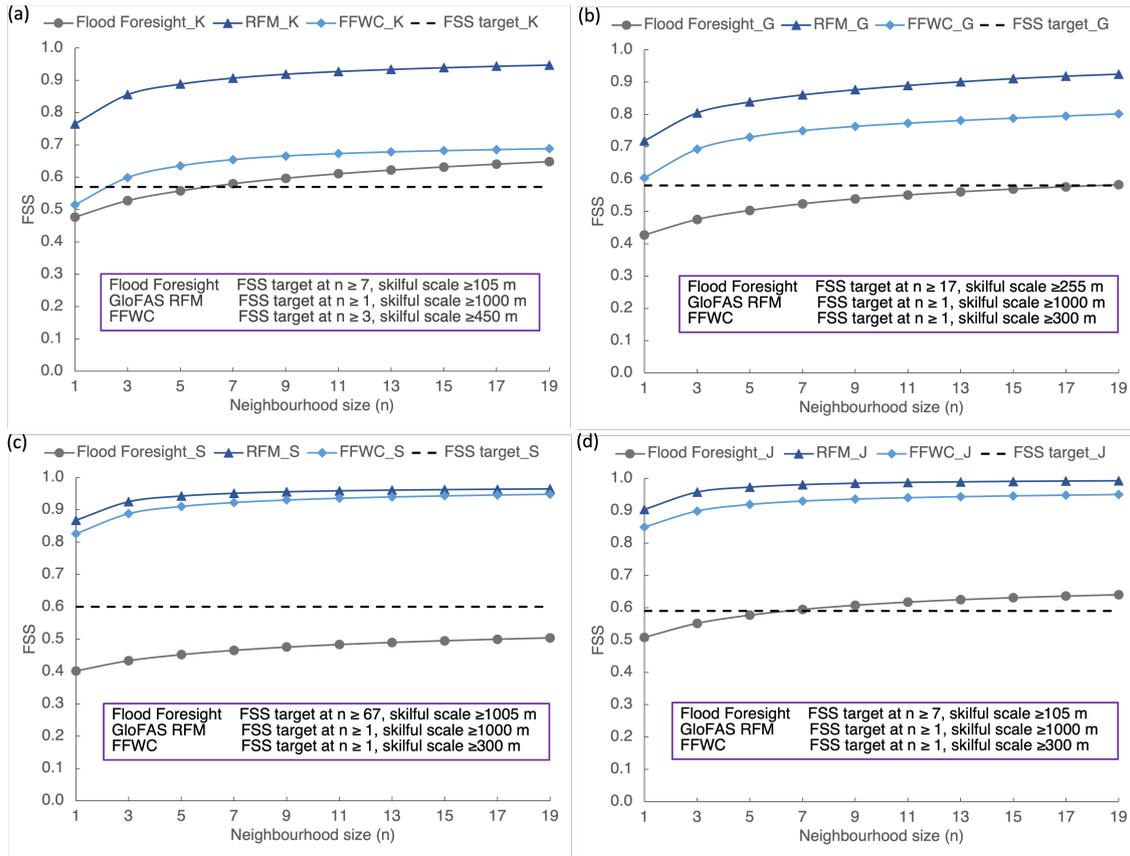


Figure 5.9: Average FSS plotted against neighbourhood size (n) for each forecast system (run dates as described in Figure 5.8) in Kurigram (a), Gaibandha (b), Sirajganj (c) and Jamalpur (d) and the target skill score for each district.

5.8) that the Flood Foresight map appeared to capture the observed flooding at Kurigram more accurately than the FFWC map. Despite this observation, the FFWC map has a higher CSI score (0.35, Fig. 5.11(c)) compared to the Flood Foresight CSI (0.31, Fig. 5.11(a)) because the grid size is not accounted for by the CSI. However, neither of the two systems flood maps have reached the target FSS_T at grid level (Fig. 5.9(a)). The FFWC map exceeds FSS_T at $n = 3$ and the Flood Foresight map exceeds FSS_T at $n = 7$. By accounting for the impact of the grid size, the skilful scale for Flood Foresight for Kurigram is 105 m ($\frac{1}{2}(7 \times 30)$) compared to 450 m ($\frac{1}{2}(3 \times 300)$) for the FFWC flood map indicating the Flood Foresight system is more accurate in Kurigram. In Gaibandha, the skilful scale

1 is similar for Flood Foresight and the FFWC maps (255 m verses 300 m, Fig. 5.9(b)).
2 In Sirajganj, Flood Foresight requires a neighbourhood size similar to the grid size of the
3 GloFAS RFM to exceed FSS_T at 1005 m (Fig. 5.9(c)). Following the reassociation of
4 the IZ in Jamalpur, Flood Foresight has a skilful scale of 105 m here compared to a high
5 FSS score for the FFWC model at 300 m grid level (Fig. 5.9(d)). Across all districts,
6 GloFAS RFM exceeds the FSS_T at grid level with the best performance in Jamalpur and
7 the worst performance in Gaibandha which can be linked to the under-estimation of flood
8 extent seen on the CSM (Fig. 5.8(b)).

9
10 The scale selective approach is also useful for comparing performance above the FSS_T
11 line where two systems exceed FSS_T at grid level ($n = 1$). For example, in Jamalpur (Fig.
12 5.9(d)) at $n = 1$ the RFM FSS is 0.90 compared to the FFWC FSS of 0.85. However,
13 FFWC reaches the same score as the RFM (at $n = 1$, 0.90), at $n = 3$. In terms of spatial
14 scale, the same FSS is reached by FFWC at $n = 3$ (450 m) as RFM at $n = 1$ (1000 m),
15 indicating that the FFWC flood map is more accurate than RFM in Jamalpur. The same
16 result is true in Sirajganj (Fig. 5.9(c)) where the FFWC FSS at $n = 3$ (0.89) exceeds the
17 RFM FSS at $n = 1$ (0.87). In Gaibandha (Fig. 5.9(b)) the FFWC FSS at $n = 5$ (0.73,
18 750 m) exceeds the RFM FSS at $n = 1$ (0.72). Overall, by accounting for grid length,
19 the FFWC provides the most accurate forecast flood map (albeit at a 0-day lead time)
20 in Jamalpur, Sirajganj and Gaibandha. RFM is most skilful in Kurigram with Flood
21 Foresight outperforming the FFWC model here.

22 5.6.3 Discussion

23 GloFAS RFMs are designed to give an early indication, up to a month in advance that se-
24 vere flooding is possible for large rivers across the globe. The RFMs are triggered when the
25 ensemble mean discharge exceeds the 10-year RP threshold over the next 30 days. The
26 RP threshold levels are calculated using ERA5 reanalysis data (Harrigan et al., 2020).

1 GloFAS RFM is triggered for the Jamuna River only once at an 8-day lead time on 18
2 July (mapping maximum extent over the next 30 days) where the flood map shows a high
3 level of skill across the districts of interest (Fig. 5.8), comparable to the local FFWC flood
4 maps, which use local observations. Unfortunately, the RFM skill reduces closer to the
5 event with no flood maps triggered after 18 July. We also see this impacting the Flood
6 Foresight skill, which reduces closer to the flood peak (Fig. 5.5). The Flood Foresight
7 system indicates flooding surrounding the Jamuna River at all forecast lead times. The
8 flood extent is generally under-predicted. This is partly due to the discharge forecast not
9 exceeding the 5-year RP threshold, which also links to unresolved/uncalibrated smaller
10 tributaries/distributaries in the GloFAS river network.

11

12 Boelee (2022) finds (for Africa, based on GloFAS ensemble reforecast flows) that the
13 forecast discharge exceedance above return period thresholds depends on both the forecast
14 lead time and the number of ensemble members considered (the probability trigger set).
15 Boelee found more flood occurrences were predicted at medium-range lead times, compared
16 to short-range lead times. We infer that these results could also apply to Bangladesh and
17 that they partly explain the RFMs best performance at 8 days lead time. The ERA5
18 reanalysis data is used to initialise the GloFAS forecast and determine the RP thresholds.
19 Currently, the RP thresholds do not account for either forecast lead time or ensemble vari-
20 ability. The GloFAS hydrographs at Bahadurabad indicate a reanalysis discharge value of
21 around the 2-year RP threshold for all lead times within 5 days of the flood peak. This is
22 a significant underestimation compared to observed river levels (both in historical context
23 and compared to FFWC danger levels and severe flood thresholds). High confidence is
24 assigned to the initial conditions in the streamflow forecast and also in the short-term
25 forecast before the ensembles show more variation at longer lead times. This leads to an
26 overall under-estimation of the flood magnitude at shorter lead times. Zsoter et al. (2020)
27 found that ensemble-reforecast-based thresholds would lead to an improved forecast at

1 lead times beyond a few days as they can account for variations in forecast skill with
2 lead-time and ensemble variability. Ensemble-forecast-based thresholds could improve
3 the flood map selection in both the RFM and Flood Foresight systems, which presently
4 are the main limitation of both systems.

5

6 The RFM uses the ensemble mean or 50% of ensemble members must exceed the RP
7 threshold to trigger a flood map. However, Boelee (2022) found that the percentage of
8 flood events exceeding the threshold for any return period dropped to less than 50% as
9 soon as the required ensemble size was increased to two ensemble members or more, for all
10 the lead times. Extreme flood event prediction, can lie in the ensemble member outliers
11 (Hooker et al., 2023a). The Flood Foresight system uses information from all ensemble
12 members for impact forecasts, which is a major benefit of this system. The automated
13 probabilistic flood maps can be produced quickly in near real-time indicating a spread of
14 possible conditions that could support the decision-making process. The RFM could see
15 an improved forecast at more lead times if any ensemble member exceeding the RP thresh-
16 old (rather than the mean) triggered the flood map selection. This would require further
17 investigation in flood prone areas so that the number of ensemble members chosen can be
18 optimised to avoid increasing false alarms. The RFM could also provide probabilistic infor-
19 mation by combining flood maps from all ensemble members that exceed the RP threshold.

20

21 The FFWC model performs significantly better at shorter lead times compared to the
22 other two systems, which is not surprising as locally observed river level and rainfall data
23 are used as input data. The benefits of the FFWC model are that no RP thresholds need
24 to be determined or exceeded to produce a flood map and there are no flood map inter-
25 polation uncertainties. The deterministic FFWC models forecast skill drops significantly
26 with forecast lead time (BWDB, 2020), which reduces the usefulness for flood mitigation
27 and FbF purposes at longer lead times. The FFWC model performs less well compared

1 to RFM and Flood Foresight near to the country's border where observations upstream
2 are unavailable.

3 **5.7 Conclusions**

4 Forecast flood maps are increasingly sought to accurately link flooding hazard to popula-
5 tions impacted to inform FbF schemes. Humanitarian agencies in Bangladesh would like
6 the impacts mapped in detail at Union level (4,571 Unions in Bangladesh) at long forecast
7 lead times (out to 10 days) so that insurance funds can be locally targeted in good time.
8 This creates a conflict between the detail or spatial scale (grid size) of the flood maps and
9 the forecast skill of the flood forecasting system. Spatial validation of forecast flood maps
10 from multiple systems has received little attention, partly due to the problem of comparing
11 skill scores from maps at different spatial scales. Here, we applied a validation approach
12 using scale-selective methods (Hooker et al., 2022) that determines a skilful scale, which
13 can be directly compared across forecast systems.

14

15 We evaluated three flood forecasting systems, Flood Foresight (30 m), GloFAS RFM
16 (1000 m) and the FFWC Super Model (300 m) each predicting flood extent at differ-
17 ent spatial scales (shown in brackets) for the Jamuna River in July 2020. Each of the
18 maps were compared against SAR-derived observations of flood extent. Evaluating Flood
19 Foresight skill at all lead times out to ten days for four districts revealed issues with unre-
20 solved/uncalibrated tributaries/distributaries in GloFAS (used to input forecast discharge
21 to Flood Foresight). Reconfiguring the Flood Foresight system led to an improved flood
22 map forecast in one district (Jamalpur), but similar issues remained in other areas. This
23 highlights one problem with trying to combine a gridded global hydrometeorological model
24 with a detailed sub-catchment network and linking this to detailed flood maps. The flood
25 mapping skill in the Jamuna basin is linked to the detail of the river network and whether

1 flood maps are triggered, which depends on exceeding the RP threshold set. Where Flood
2 Foresight maps were triggered, such as in Kurigram, the flood map accuracy outperformed
3 the local FFWC model. This is due to its location next to the border of Bangladesh and
4 the lack of upstream observations. In other areas the FFWC model captures more de-
5 tail in the river network and shows less under-prediction compared to the other systems.
6 For FbF applications and humanitarian response in Bangladesh, a combination of Flood
7 Foresight and the local FFWC model could produce flood inundation maps at a useful
8 scale that can be linked to flooding impacts. Flood Foresight has the benefit of forecast
9 skill at a longer lead times (up to 10 days) with probabilistic maps accounting for some
10 of the forecast uncertainty. Flood Foresight could be supplemented or linked to driving
11 data from the FFWC model at shorter lead times to incorporate local observations and a
12 higher resolution river network. This would avoid regions of non-trigger where no flood
13 map is selected from the Flood Foresight library due to the forecast discharge not ex-
14 ceeding the RP threshold. GloFAS RFM is designed as a deterministic ‘heads-up’ tool at
15 a coarse resolution that would be difficult to link to impacts for FbF applications in its
16 current configuration. All systems miss flooding across a wide area captured by the SAR
17 data that is possibly due to surface water flooding (SWF). These fluvial flood forecasting
18 systems are not designed to map SWF and we recommend combining the forecast flood
19 maps with SAR-derived flood maps through data assimilation so that SWF can be ac-
20 counted for in post event impact calculations for FbF schemes. Alternatively, a combined
21 fluvial/pluvial flood forecasting system would be ideal, however pluvial flood forecasting
22 practice is currently less developed in part due to the difficulties in observing SWF and
23 accurately predicting convective rainfall (Speight et al., 2021).

24

25 For future spatial validation of flood maps using SAR data, we recommend making
26 use of the Copernicus GFM product (GFM, 2021) which maps flooding detected from all
27 Sentinel-1 images since October 2021. Importantly, an exclusion mask layer is available

1 which can be used to exclude areas where SAR is unable to detect flooding such as in dense
2 urban areas, under vegetation and near steep topography. The forecast flood maps would
3 no longer be penalised for over-prediction in regions where the SAR cannot reliably detect
4 flooding. Another opportunity for improvement lies with the flood map library. Both
5 Flood Foresight and RFM maps are currently undefended. However, both would likely be
6 improved if flood defence information from FFWC could be incorporated, allowing areas
7 benefiting from those defences to be discounted when flood conditions are at return periods
8 lower than the standard of protection of the defence. Alternatively, a new high resolu-
9 tion DTM including local defence features, ideally through acquiring LiDAR data (where
10 locally obtained), would improve the local accuracy of the flood maps. The maps held
11 within the simulation library could be hydrodynamically precomputed at lower discharge
12 return periods, which would avoid inaccuracies caused by flood map interpolation beneath
13 the lowest RP level. This would increase the confidence in these flood maps so that a lower
14 RP threshold could be used to trigger FbF allowing more people access to insurance funds.

15

16 Fortunately, some of the issues discussed here such as the river network detail will be
17 partly resolved by the major upgrade to GloFAS with version 4.0 due in 2023 (Grimaldi,
18 2022). Significantly, the spatial resolution of GloFAS will increase 4 fold to approximately 5
19 km grid size. The river network will increase similarly, and more distributaries/tributaries
20 will be included in Bangladesh. This upgrade along with the use of ensemble-reforecast-
21 based RP thresholds should improve the flood map selection process used by both GloFAS
22 RFM and Flood Foresight. Scale-selective validation methods will enable future system
23 changes to be evaluated and compared meaningfully. Ideally, this would be automated
24 and integrated into the flood forecasting system.

1 5.8 Appendix

Table 5.2: Binary performance measures and formulas. Contingency categories: correctly predicted flooded (A), over-prediction (B), under-prediction (C), correctly predicted unflooded (D)

| 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 |

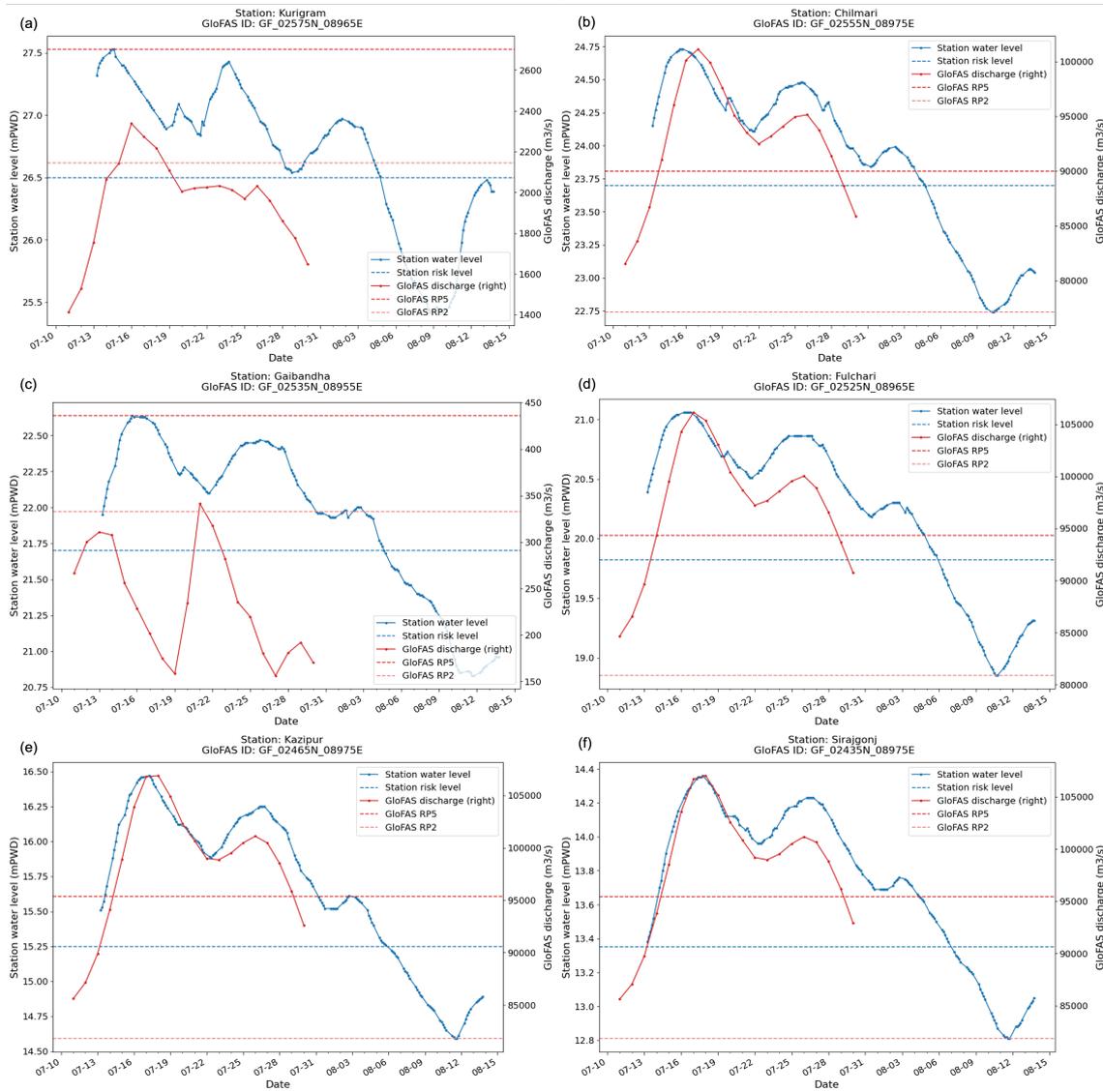


Figure 5.10: Six FFWC station water levels during July and August 2020 across the four districts compared with closest GloFAS grid cell discharge.

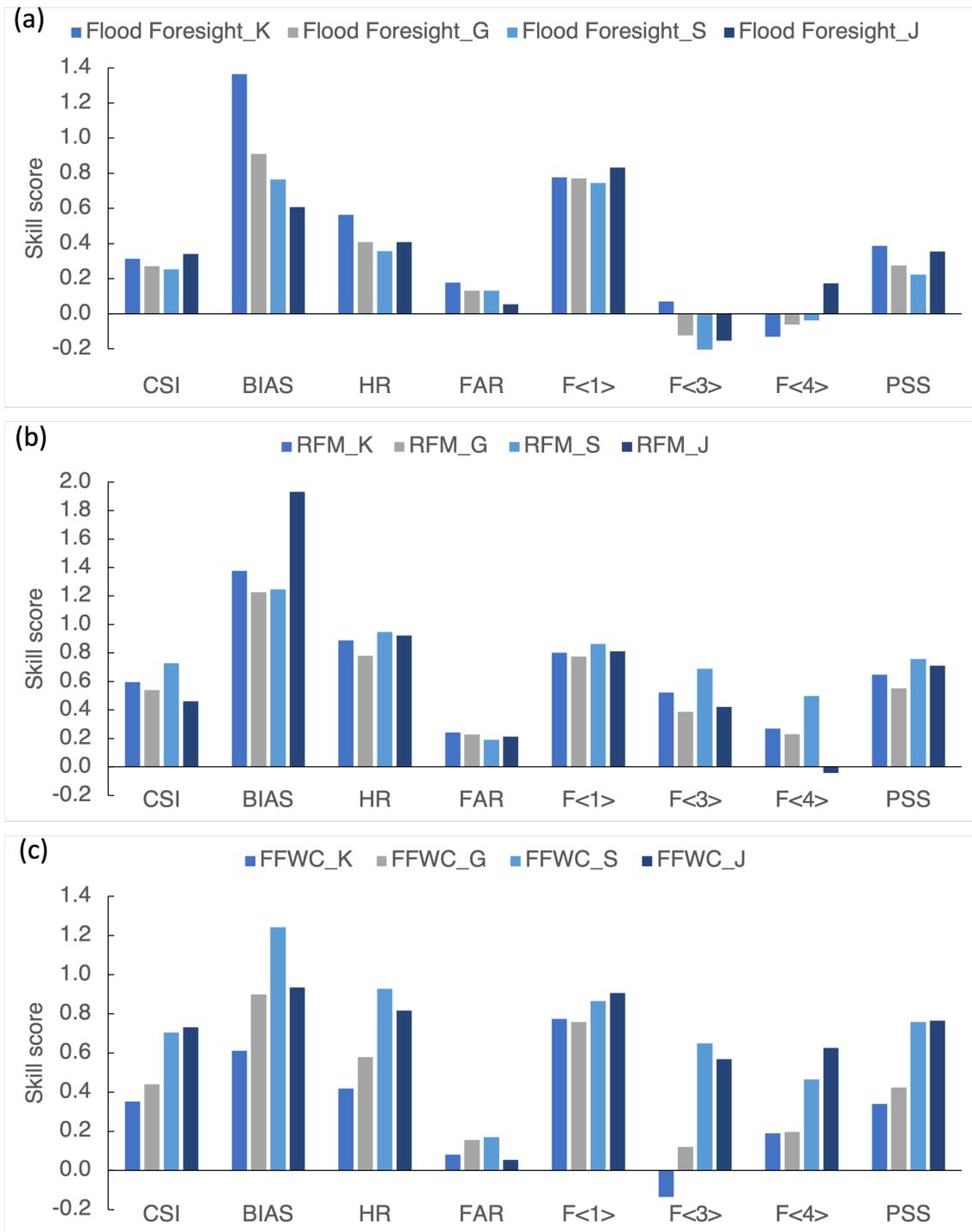


Figure 5.11: Binary performance measures for each district. Flood Foresight run date 18 July (a), GloFAS RFM 18 July (b) and FFWC flood map 25 July (c).

1 **5.9 Chapter summary**

2 In this chapter we apply the scale-selective methods presented in Chapter 3 to evaluate
3 three flood forecasting systems, each presenting flood inundation forecasts at different
4 spatial scales. Conventional binary performance measures would result in biases linked
5 to spatial scale. We show how the scale-selective methods can be used to meaningfully
6 compare across different flood forecasting systems. The benefits and limitations of the
7 systems could be discussed as a result of the evaluation process. The issues identified
8 here, that can lead to a non-trigger of flood maps in some sub-catchments, are addressed
9 in Chapter 6. In Chapter 6 we will use the SAR data directly to improve the flood map
10 selection for previously non-triggered flood maps in the Flood Foresight system.

1 Chapter 6

2 Updating simulation library flood 3 map selection through assimilation 4 of probabilistic SAR-derived flood 5 extent

6 In this chapter we address the fourth research question outlined in Chapter 1; Does a data
7 assimilation framework improve the analysis of flood inundation from a simulation library
8 system?:

- 9 • Can we incorporate probabilistic information from remotely observed flood inunda-
10 tion into a data assimilation framework to improve the flood map selection within a
11 simulation library flood forecasting system?
- 12 • How does the analysis flood map compare to independent validation data?

The Art of DA



1 **6.1 Abstract**

2 Mitigating against the impacts of catastrophic flooding requires funding for the commu-
3 nities at risk, ahead of an event. Simulation library flood forecasting systems are being
4 deployed for forecast-based financing (FbF) applications. FbF schemes use hydromete-
5 orological predictions, linked to flood inundation maps that overlay local population (or
6 other) impact maps to trigger a financial payment in advance of flooding. The FbF trigger
7 is usually automated and relies on the accuracy of the flood inundation forecast. This re-
8 liance can lead to missed events that were forecast below the trigger threshold required or
9 were not predicted at all. However, earth observation data from satellite-based synthetic
10 aperture radar (SAR) sensors can reliably detect most large flooding events, although
11 their use is limited in urban areas. Data assimilation (DA) combines forecast information
12 with observation information to improve the current system state (the analysis). A new
13 DA framework is presented to update the flood map selection from a simulation library
14 system using SAR data, taking account of observation uncertainties. By utilising flood
15 extent likelihood data derived from Sentinel-1 SAR images, we derive a new cost function
16 that must be minimised. By iteration through the flood map library we optimise the flood
17 map selection per sub-catchment. We have applied our DA method to the Pakistan 2022
18 flood. During this flood, the Indus River in the Sindh province downstream of the Sukkur
19 barrage was not forecast to reach flood levels, which resulted in a non-trigger of the FbF
20 scheme for this region. Our experiments have focused on three different scenarios: a large
21 city area with limited SAR flood detection (due to a dense urban area); sub-catchments
22 with mixed rural and urban areas; and flood edge sub-catchments. We found that the
23 flood map selection could be triggered in four out of five sub-catchments tested, with
24 the exception occurring in the dense urban area due to the simulation library flood map
25 accuracy here. Thus, the analysis flood map, created by assimilating observations from
26 SAR flood likelihood data, has potential to be used to trigger a secondary finance scheme

1 during a flood event and avoid missed financing opportunities for humanitarian action.

2 **6.2 Introduction**

3 Our warmer climate is increasing the frequency and intensity of extreme weather events
4 and the exposure and vulnerability of communities and individuals (Pörtner et al., 2022).
5 Large-scale flood forecasting systems predicting flood inundation extent are increasingly
6 used for disaster risk reduction to improve preparedness ahead of a major flooding event
7 (Stephens & Cloke, 2014a; Hooker et al., 2023b; Wu et al., 2020). An ensemble flood
8 forecasting system creates probabilistic flood maps indicating the likelihood of flooding
9 across a region or country. Flood impact risk factors such as population density, land-use
10 types or vulnerable infrastructure can also be mapped for the same area. The forecast-
11 flood-likelihood maps can be overlaid with impact maps and depending on the severity of
12 the hazard and the level of impact, a risk profile can be determined. The flood risk profile
13 can be used to inform forecast-based financing (FbF) schemes that enable the pre-release
14 of funds based on the flood forecast, ahead of the flood event (Coughlan de Perez et al.,
15 2015, 2016). Automation of FbF schemes is important for rapid action to take place to
16 mitigate against flooding impacts. The skill of the flood forecasting system is key to *trig-*
17 *gering* the FbF scheme. A *non-trigger* of FbF ahead of or during a flood event might
18 prove catastrophic for those impacted.

19

20 Advances in flood forecasting both at global and local levels link together meteorolog-
21 ical and hydrological forecasts of river discharge that drive the selection of pre-computed
22 flood maps from a simulation library (Speight et al., 2021; Hooker et al., 2023a). The
23 use of a simulation library obviates the need to run a hydrodynamic model as part of the
24 forecast process, reducing computation time and allowing near real-time updating for large
25 areas, which otherwise presents a significant challenge. The flood maps within the library

1 are at a relatively higher spatial resolution (e.g. 30 m) compared to the resolution of the
2 driving global hydrological model (e.g. approximately 5 km). This mismatch in scales can
3 lead to problems with flood map selection and can cause gaps where the minimum return
4 period threshold has not been exceeded (a non-trigger) by the forecast discharge (Hooker
5 et al., 2023b). The three main issues that cause this in the global scale model are the rep-
6 resentation of river networks, the return period thresholds determined and the exclusion
7 of dam operations. Rivers that are narrower than a particular width, or catchment areas
8 smaller than a pre-determined size are not resolved by global scale models. In addition,
9 the return period thresholds set may be poorly calibrated due to a lack of ground truth
10 observational data such as river discharge or river water level (Boelee, 2022). These two
11 limitations can lead to a non-trigger, i.e. no flood map is selected from the simulation
12 library for a particular sub-catchment. Also, local dam operations such as diversions of
13 river water for irrigation purposes or rapid releases of flood waters downstream, are not
14 generally included in global scale models. This can lead to over- or under-prediction of
15 forecast discharge, resulting in inaccurate or non-trigger of flood map selection in the fore-
16 cast.

17
18 Satellite-derived observations of flooding have the potential to bring additional spatial
19 information into flood inundation forecasts compared to in situ point gauging stations.
20 These observations could be used to update and improve the FbF scheme either as part of
21 a secondary finance payment following the acquisition of the satellite data or to improve
22 the flood inundation forecasts going forwards as the flood event evolves. Synthetic aper-
23 ture radar (SAR) sensors are particularly useful for remote flood detection, since they can
24 see through cloud, most weather and are active both day and night (e.g. Mason, Daven-
25 port, et al., 2012; Schumann et al., 2022). Previously, SAR data have been used in several
26 different ways to improve hydraulic models and flood prediction through data assimila-
27 tion (DA). Data assimilation finds an optimal state (such as water level) and/or model

1 parameter values by accounting for the previous forecast, the observations available, and
2 both of their associated uncertainties. The updated state (analysis) and/or parameter set
3 are used to initiate the next forecast in a feedback loop or cycle. A review of approaches
4 used to assimilate satellite-derived data into hydraulic models (from 2007 until 2015) can
5 be found in Table 7 of Grimaldi et al. (2016) and Table 1 of Revilla-Romero et al. (2016).

6
7 Here, we summarise the different ways that observations of flood extent and water
8 levels are used in various DA approaches and discuss some limitations. In order to extract
9 flood extent from SAR data an image classification technique must be applied (see Section
10 3, Grimaldi et al. (2016)). Binary maps of flood extent were assimilated by Lai et al. (2014)
11 using a 4D variational approach to successfully estimate the roughness parameter over the
12 flood plain. More recent variational approaches combine in situ observations with high-
13 resolution hydrometeorology and satellite altimetry data into a hydraulic–hydrological
14 numerical model (Pujol et al., 2022). Other DA methods applied to flooding include fil-
15 tering methods such as particle filtering (PF) and ensemble Kalman filter (EnKF) methods
16 (van Leeuwen, 2009; Evensen, 1994). These methods have been used to assimilate SAR-
17 derived water levels (WL), found by intersecting the edge of the binary flood map with
18 a digital elevation model (DEM) (Mason et al., 2007; Mason, Davenport, et al., 2012).
19 SAR-derived WL only provide information at the flood edge and rely on the spatial res-
20 olution and the vertical accuracy of the underlying DEM (Dasgupta et al., 2021a), which
21 makes them difficult to obtain.

22
23 A probabilistic flood mapping procedure for SAR data was first introduced by Giustarini
24 et al. (2016). This created the potential for flooding probabilities to be assimilated directly.
25 More recently, observation uncertainty associated with classifying flood extent from SAR
26 data is openly available through the Copernicus Emergency Management Service (CEMS)
27 (Copernicus Programme, 2021). Probabilistic flood maps from SAR were used to quan-

1 tify observation uncertainty for assimilating flood extents by Hostache et al. (2018) using
2 a PF approach. Hostache et al. found an improvement in forecast performance of WL
3 compared to those found by García-Pintado et al. (2013); García-Pintado et al. (2015)
4 who assimilated WL using a local ensemble transform Kalman filter (ETKF). Cooper et
5 al. (2019) considered the direct assimilation of SAR backscatter values, avoiding the step
6 of deriving an intermediate flood extent first. Cooper et al. assimilated pseudo-observed
7 SAR backscatter values directly into a 2D hydrodynamic model using an ETKF approach
8 and compared three different observation operators. The authors found a new backscatter
9 observation operator performed well compared to more conventional options.

10

11 In general, filtering methods suffer from degeneracy and limited persistence of improve-
12 ment following DA. Several attempts have been made to address this. By using a mutual
13 information-based likelihood function, Dasgupta et al. (2021b) increased the persistence of
14 the DA impact on improvement of both water depth and discharge. Di Mauro et al. (2021,
15 2022) aimed to overcome degeneracy issues by applying a tempered particle filter (TPF)
16 and particle mutation so that variation across the ensemble is maintained. Comparisons
17 against sequential importance sampling methods concluded that the TPF approach im-
18 proved the persistence of the DA. More recent studies using EnKF methods include the
19 assimilation of distributed WL derived from optical satellite data combined with WL from
20 river gauges (Annis et al., 2022). Annis et al. found that improvements from the DA by
21 including the satellite-derived WL data spread further across the domain compared to
22 those without. However, these improvements also suffered from limited persistence.

23

24 One major assumption made when assimilating SAR data is that the observation er-
25 rors are from a Gaussian distribution. Nguyen, Ricci, Piacentini, Fatras, et al. (2023)
26 used a novel approach by assimilating flood extent data expressed as a wet surface ratio
27 (Nguyen, Ricci, Piacentini, Simon, & Rodriguez-suquet, 2023) and Gaussian anamorpho-

1 sis (GA) to transform the non-Gaussian observation errors into Gaussian. They found
2 this gave slightly better results compared to DA without GA but concluded that the sub-
3 optimal assumption of Gaussian errors in EnKF may still give reliable and valid results.

4

5 All of these studies involved the assimilation of satellite-derived data to update hy-
6 drodynamic model states and/or parameter values by taking a data assimilation cycling
7 approach. The assimilation provides updated initial conditions ahead of the next forecast
8 cycle. Our new approach differs significantly as we aim to develop a DA framework to
9 assimilate probabilistic flood maps into a simulation library flood inundation forecasting
10 system. We aim to improve the flood map selection creating a new analysis flood map
11 without a feedback loop by utilising spatially distributed flood likelihood information de-
12 rived from SAR data. We test the DA framework using forecast and optical satellite
13 observation data from a major flood event in Pakistan, August 2022 (Floodlist, 2022).

14

15 In this chapter, the flood forecasting system and derivation of the static simulation
16 library along with satellite-derived observations of flood likelihood are outlined in Section
17 6.3. The development of the DA framework and verification methods are described in
18 Section 6.4. Section 6.5 presents an overview of the 2022 Pakistan flood and details of
19 the data used. The DA framework successfully triggered flood maps in 4 out of 5 sub-
20 catchments tested as shown in our results, discussed in Section 6.6. We conclude in Section
21 6.7 with recommendations for future work to improve the use of SAR flood likelihood data
22 to update a simulation library flood forecasting system through data assimilation.

1 **6.3 Simulation library forecasting system and observation** 2 **data**

3 Flood Foresight, a simulation library flood inundation forecasting system and its appli-
4 cation for disaster risk reduction through FbF is outlined in Section 6.3.1. Section 6.3.2
5 details the derivation of the flood likelihood data from SAR that will be used as obser-
6 vation data in the assimilation process. The extraction of flood extent information from
7 optical images that will be used for validation is explained in Section 6.3.3.

8 **6.3.1 Flood Foresight and Forecast-based-Financing**

9 The Global Flood Awareness System (GloFAS) couples global ensemble weather forecasts
10 with a hydrological model and provides daily ensemble forecast river discharge at approx-
11 imately 10 km grid size (v3.2, GloFAS (2021)). The Flood Foresight system (Revilla-
12 Romero et al., 2017; Hooker et al., 2023a) is a fluvial, probabilistic flood inundation
13 forecasting system. Flood Foresight is set up by dividing the catchment into ‘Impact
14 Zones’ (IZ) or sub-catchments using the HydroBASINS data set (Lehner, 2014b). Each
15 IZ in Flood Foresight is linked to a GloFAS grid cell that provides a 51 ensemble member
16 forecast of river discharge. Flood Foresight contains a simulation library of precomputed
17 flood depth and extent maps. The flood map library was hydrodynamically modelled us-
18 ing JFlow[®], (Bradbrook, 2006) and RFlow using a detailed 30 m digital surface model.
19 The maps were modelled at specific return period (RP) thresholds (20, 50, 100, 200, 500
20 and 1500 years). Subsequently, these were linearly interpolated at 5 intermediate intervals
21 between each RP threshold and extrapolated between zero and the 20 year RP flood map
22 (totalling 36 flood maps). Depending on the forecast discharge from GloFAS for each IZ,
23 a flood map is selected from the simulation library. The flood map selected is determined
24 by the RP threshold exceeded within each IZ. The resultant forecast flood map is created
25 by stitching together individual IZ flood maps (at various RP’s) and is produced daily out

1 to 10 days ahead.

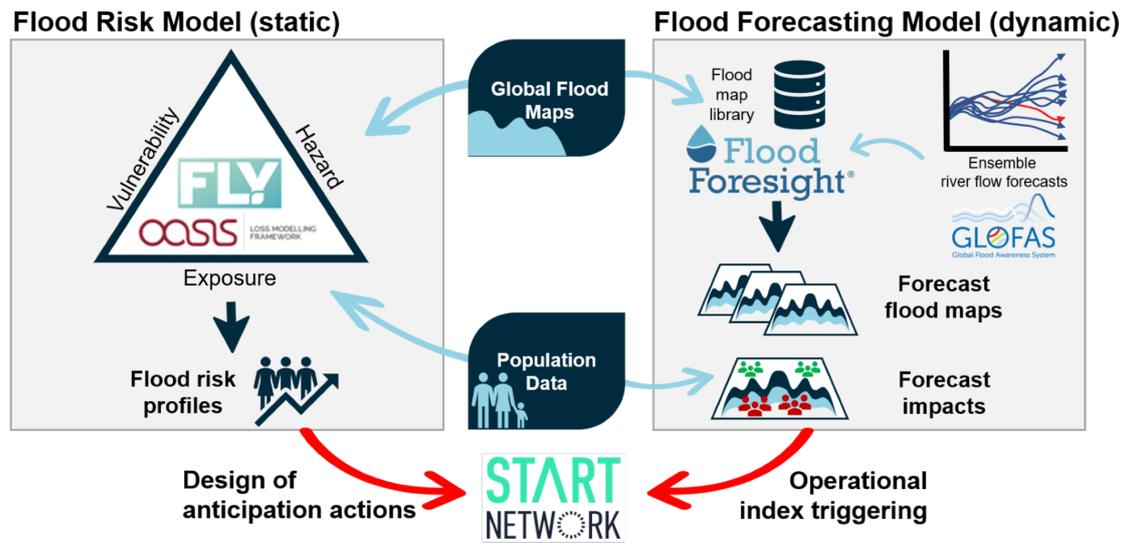


Figure 6.1: Flood Foresight/Start Network ensemble flood inundation forecast and population impacts work flow.

2 The charity Start Network (Start Network, 2022) brings together a group of over 80
 3 humanitarian agencies and aims to develop local community-led, early action through a
 4 model of proactive funding to mitigate against the impacts of crises. JBA Consulting, in
 5 partnership with Start Network, have developed a Disaster Risk Financing (DRF) system
 6 for the Indus River basin in Pakistan that links Flood Foresight forecast flood maps to
 7 populations impacted by flooding (Fig. 6.1). For the purposes of setting FbF trigger
 8 threshold levels, the DRF system quantifies the flood risk to the population through a
 9 probabilistic global catastrophe risk model, FLY (Dunning, 2019). The analysis flood map
 10 produced here as a result of the DA relates to the dynamic operational index triggering.

11 6.3.2 Satellite-derived flood likelihood

12 The GFM service (GFM, 2021) combines the outputs of three different algorithms to ex-
 13 tract flood extent and uncertainty information from Sentinel-1 SAR data. The process is
 14 automatic and runs continuously in near real-time (within 8 hours after image acquisition)

1 for every SAR image detecting flooding on a global-scale. The resulting flood informa-
2 tion layers are openly available through open access. The mini-ensemble approach used
3 for flood detection increases the confidence in the flood detecting capabilities of SAR.
4 The first flood mapping algorithm developed by the Luxembourg Institute of Science and
5 Technology (LIST) use a pair of SAR images (pre-flood and flood) and a hierarchical
6 split-based change detection approach to classify permanent and flood waters (Chini et
7 al., 2017). The classification uncertainty depends on the Bayesian inference classification
8 probabilities. The second flood mapping algorithm by the German Aerospace Research
9 Centre (DLR) uses a hierarchical tile-based thresholding approach and the optimization
10 of the classification by combining various information sources using fuzzy-logic theory and
11 region growing (Martinis et al., 2015; Twele et al., 2016). The uncertainty information
12 depends on fuzzy memberships. The final algorithm from TU Wien uses the historical
13 time series of the SAR backscatter values per pixel and classifies flood extent from the
14 backscatter probability distribution (Wagner et al., 2020; Bauer-Marschallinger et al.,
15 2022). The classification uncertainty is based on the Bayesian posterior probability. The
16 output flood layer is derived using the mini-ensemble with a pixel classified as flooded
17 where two out of three of the algorithms determines a flood class. The flood likelihood
18 values are aggregated, first by converting each to lie in the same range $[0, 100]$ before
19 averaging the likelihood values. Regions where SAR is unable to detect flooding due to
20 shadow or layover effects are removed from the classification process. This usually in-
21 cludes dense urban areas, densely vegetated areas, regions with steep slopes and regions
22 that might appear flooded such as dry, sandy desert-like surfaces. The exclusion mask is
23 available to download as an additional layer. Permanent and seasonal water bodies are
24 classified separately as a reference water mask layer.

25

26 Krullikowski et al. (2023) applied and assessed the usefulness of GFM ensemble like-
27 lihood on two test sites in Myanmar and Somalia, both situated in challenging areas

1 for flood detection using SAR data. Krullikowski et al. found that the GFM ensemble
2 likelihood layer resulted in a simplified appraisal of trust in the ensemble flood extent
3 detection approach and provides more reliable and robust uncertainty information for
4 detecting flooding compared to using a single algorithm only.

5 **6.3.3 Optical Normalized Difference Water Index (NDWI)**

Occasionally, optical satellite data can be useful for observations of flood extent. Flood
detection from optical satellites depends on a near cloud free sky where the satellite ac-
quisition coincides with the flood event. The Normalized Difference Water Index (NDWI)
for flood and surface water detection is calculated with Sentinel-2 optical data using the
green band (B03) and NIR band (B08) (Albertini et al., 2022). The NDWI is given by:

$$NDWI = \frac{B03 - B08}{B03 + B08}, \quad (6.1)$$

6 where positive values indicate water. Albertini et al. (2022) reviewed the performance of
7 surface water and flood detection metrics using multispectral satellite data. They found
8 that the average overall accuracy from previous flood studies applying NDWI to be 87.85%
9 and for permanent surface water studies scored 94.41%. This included studies using data
10 from different satellite sensors with spatial resolutions ranging from 10 m for Sentinel-2
11 to 500 m for Terra-Aqua MODIS.

12 **6.4 Methods**

13 **6.4.1 Data assimilation framework**

14 The aim of the DA framework is to update a previous forecast of flood inundation extent
15 and depth from Flood Foresight (the background) where a non-trigger has occurred in the
16 forecast system but where flooding was derived from concurrent satellite-based SAR data.

1 Using observation uncertainty information from the GFM flood likelihood layer (Section
 2 6.3.2), we aim to improve the flood map selection for non-triggered IZ by minimising a
 3 cost function per IZ.

4

The data assimilation framework aims to update the state vector, $\mathbf{x} \in \mathbb{R}^n$ representing flood depths at each grid cell location. The total number of grid cells across an IZ is n , the total of observed unmasked grid cells is m . To find the optimum state accounting for observation uncertainty we define the observation *likelihood* term $P(\mathbf{y}|\mathbf{x})$ where observations, $\mathbf{y} \in \mathbb{R}^m$ have two possible binary outcomes, $y = 1$ flooded and $y = 0$ unflooded, following classification from SAR data. The likelihood term can be represented by a Bernoulli distribution (Lauritzen, 2023), defined as

$$P(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^m L_i^{H(x_i)} (1 - L_i)^{1-H(x_i)}, \quad (6.2)$$

where L_i is the GFM flood likelihood value (see Section 6.3.2). The observation operator $\mathbf{H}(\mathbf{x})$ defined as

$$\mathbf{H}(\mathbf{x}) = \begin{cases} 1 & \text{flooded} & x_i > 0.2m \\ 0 & \text{unflooded} & \text{otherwise} \end{cases}, \quad (6.3)$$

acts to convert flood depths (state space) to a binary flood class (observation space) at unmasked grid cells. Thus, we exclude observation likelihood information for masked grid cells where SAR data cannot reliably detect flooding. To find the maximum posterior likelihood of the state variable, we take the negative log likelihood of $P(\mathbf{y}|\mathbf{x})$ and divide by the number of unmasked grid cells m to derive the cost function (averaged across unmasked grid cells per IZ)

$$J(\mathbf{x}) = -\frac{1}{m} \sum_{i=1}^m \{H(x_i) \ln(L_i) + (1 - H(x_i)) \ln(1 - (L_i))\}. \quad (6.4)$$

1 The value of $J(\mathbf{x})$ is calculated per IZ by iterating through the flood map library (36 flood
2 maps) and finding the flood map return period by minimising $J(\mathbf{x})$ and maximising the
3 posterior likelihood, given the observation uncertainty data. To ensure that the minimum
4 is reached across the flood map library, $J(\mathbf{x})$ is calculated for all 36 flood maps. Note
5 that this is different to the standard data assimilation approach, minimisation would be
6 accomplished via a gradient descent algorithm (Bannister, 2017). Following the assimila-
7 tion process, replacing the non-triggered IZ with updated flood maps results in an analysis
8 flood depth and extent map (the flood depth information is contained within the simula-
9 tion library). This means that the analysis flood map remains consistent with the Flood
10 Foresight system where the flood maps have been hydrodynamically modelled, i.e. they
11 are physically realistic. Retaining the flood depth information is important for FbF appli-
12 cations for quantifying the risk of flood impacts. Since the observations have binary values
13 (flooded/unflooded), we cannot distinguish between floods that have the same extent but
14 different depths from the observation data. This property is inherited in the cost function.

15 **6.4.2 Validation methods**

The resulting analysis flood map, following assimilation of SAR-derived flood likelihood data, is validated by comparing against independent flood extent observation data derived from optical satellite data, the NDWI (Section 6.3.3). The results will be validated by calculating the Fraction Skill Score (FSS, Roberts and Lean (2008); Hooker et al. (2022)) and by mapping the performance on a Categorical Scale Map (CSM, Dey et al. (2014); Hooker et al. (2022, 2023a)). Both the FSS and the CSM avoid issues with the double penalty impact of conventional binary performance measures as well as the impact of flood magnitude on the skill score (Hooker et al., 2022). The FSS is based on the Brier Skill Score and uses a neighbourhood approach to determine the skillful spatial scale of the analysis flood map. The fraction of flooding within a given square neighbourhood size of length n is compared by calculating the mean-squared-error (MSE) between the analysis

and the validation flood maps to give

$$FSS_n = 1 - \frac{MSE_n}{MSE_{n(ref)}}, \quad (6.5)$$

where $MSE_{n(ref)}$ is a potential maximum MSE that depends on the fraction of flooding across the IZ on the analysis and the validation flood maps. A skilful scale is determined when $FSS \geq FSS_T$, the target FSS score, where $FSS_T \geq 0.5 + \frac{f_o}{2}$ depends on the fraction of observed flooding across the IZ, f_o . When the analysis and validation flood extents are equal in area across an IZ there is said to be no background bias and the maximum FSS is 1. Otherwise, the maximum asymptotic FSS (AFSS) is given by

$$AFSS = \frac{2f_o f_a}{f_o^2 + f_a^2}, \quad (6.6)$$

1 where f_a is the is the fraction of flooding on the analysis flood map per IZ.

2

3 The CSM plots a local agreement scale (S) at every grid cell. An overview of the
 4 method is presented here. Please see Chapter 3 or Dey et al. (2014); Hooker et al. (2022,
 5 2023a) for full details of the methodology. A background bias between the analysis and
 6 verification flood maps that is deemed acceptable is predetermined. The pre-set bias is
 7 used to calculate an agreement criterion that must be reached by the flood map comparison
 8 calculation. The comparison begins at each grid cell $n = 1$, if the agreement criterion is
 9 met at grid level, the grid cell is labelled with an agreement scale $S = 0$. Where the
 10 criterion is not met, a larger neighbourhood size is compared (e.g. $n = 3$). The fraction
 11 flooded in each of the analysis and validation flood maps are compared and if the criterion
 12 is met, the agreement scale assigned would be $S = 1$. The process continues to larger
 13 neighbourhoods (e.g. $n = 7, S = 3$) until either the criterion is met or a predetermined
 14 limit is reached (S_{lim} , set to 9 for this application). The agreement scale at this limit
 15 would indicate a false alarm or miss for the grid cell. Note that the relationship between

1 n , the neighbourhood size used for the FSS, and S is given by $S = (n - 1)/2$. The
2 agreement scales are combined with data from a conventional contingency map (Stephens
3 et al., 2014) and are mapped across an IZ. The CSM indicates a location-specific level of
4 agreement and shows where the flood maps are over- or under-estimating flooding.

5 **6.5 Pakistan flood 2022**

6 **6.5.1 Event overview**

7 In Spring 2022, Pakistan saw a record-breaking heatwave with temperatures exceeding
8 50°C. The heat exacerbated upstream glacial snow melt feeding the Indus River basin,
9 which runs over 3000 km across the length of Pakistan, draining the Himalayas to the
10 Arabian Sea. An intense monsoon season followed in July and August, driving multi-
11 ple flood-producing mechanisms including multi-day extreme precipitation that was the
12 primary driver of floods (World Weather Attribution, 2022; Nanditha et al., 2023). At-
13 tribution studies indicate that the 5-day maximum rainfall over the provinces Sindh and
14 Balochistan, which led to catastrophic flooding, was made 75% more intense by 1.2°C of
15 global warming (World Weather Attribution, 2022). The northern Sindh province received
16 an estimated 442.5 mm of rainfall in August, 784% more than usually recorded, causing
17 inundation of 55,000 km² across the region (Floodlist, 2022). Despite early warnings of
18 the potential for significant flooding from GloFAS, the unimaginable scale and magnitude
19 of the flood impacted over 33 million people with over 1700 lives lost and costing more
20 than \$40 billion in economic damages (Floodlist, 2022; World Resources Institute, 2023).

21 **6.5.2 Data**

22 The DA framework (Section 6.4.1) was tested in the northern Sindh province where
23 widespread flooding occurred during August 2022. Figure 6.2 maps the NDWI derived
24 from Sentinel-2 optical data (Section 6.3.3) that was used to verify the resultant analysis

1 flood map following DA. Local reports and photographs of flooding were made in the cities
of Sukkur and Larkana (DAWN, 2022; The Guardian, 2022; Sky News, 2022). The DA

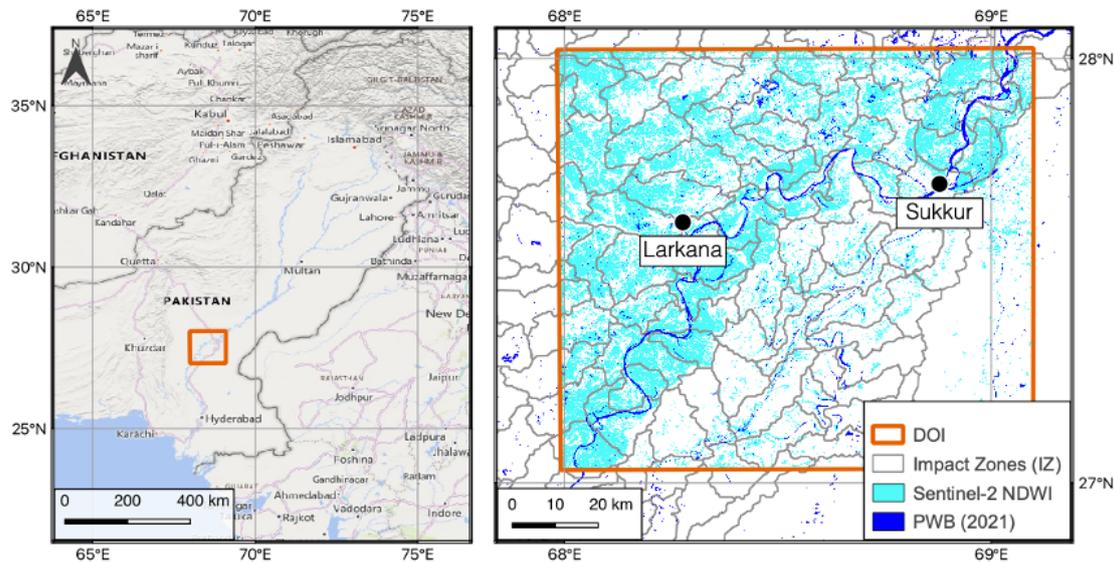


Figure 6.2: The domain of interest (DOI) is located on the Indus Basin, Sindh province, Pakistan (left). The region is divided into sub-catchments or Impact Zones (IZ) in Flood Foresight (right). Satellite-derived flooding (NDWI) from Sentinel-2 data (Section 6.3.3) from 31 August 2022 is highlighted along with permanent water bodies (PWB).

2

3 was applied to 5 IZ covering 3 different scenarios (Fig. 6.3): (1) One IZ where a large pro-
4 portion of the IZ is a dense urban area and is masked (where the GFM product is currently
5 unable to detect flooding), Sukkur (S), see Figure 6.3(a); (2) Two IZ with mixed urban
6 and rural areas, Larkana north and south (LN, LS); and (3) two flood edge locations (FE1,
7 FE2). The GFM flood likelihood data used to represent observation uncertainty in the
8 DA is mapped in Figure 6.3(a) where darker shades of orange indicate a higher likelihood
9 of flooding (Section 6.3.2). The forecast data from Flood Foresight is mapped in Figure
10 6.3(b) where the purple shades indicate the maximum return period flood map triggered
11 by the system from 10 to 31 August, 2022. Each of the IZ selected were non-triggered
12 IZ during this period. The driving forecast river discharge data from GloFAS did not
13 reach the required threshold to trigger a flood map along the central Indus channel. This

1 is likely due to poor calibration of GloFAS due to a lack of observation of river stage or
 2 discharge. The Sukkur barrage operations for diversion and altering of river water flows
 3 are not currently included in the GloFAS system, which makes forecasts unreliable along
 this stretch of the Indus River.

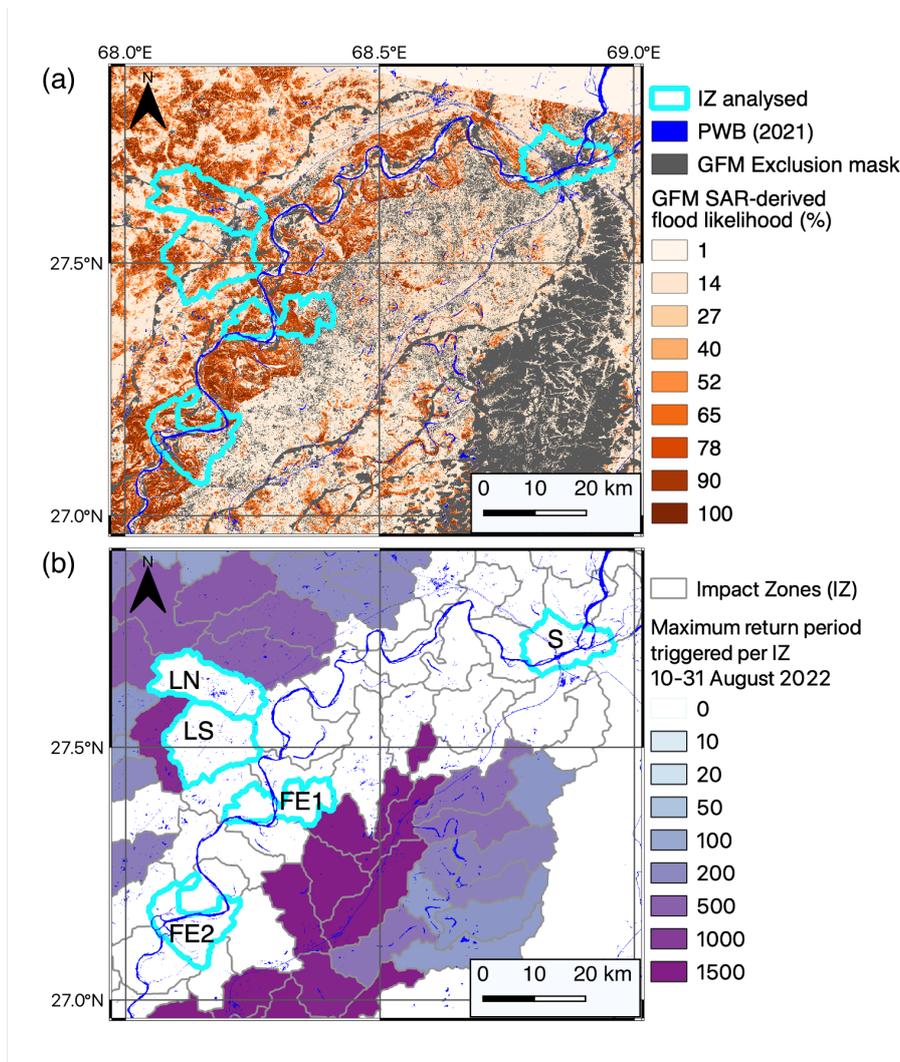


Figure 6.3: (a) GFM flood likelihood derived from Sentinel-1 SAR data, masked areas (grey) indicate where flooding cannot be reliably detected from SAR data. (b) The maximum return period threshold triggered per IZ by the Flood Foresight system during peak flooding 10-31 August, 2022. Five non-triggered IZ of interest labelled: S - Sukkur, LN - Larkana North, LS - Larkana South, FE1 - flood edge 1, FE2 - flood edge 2.

1 **6.6 Results and discussion**

2 Results are presented for the 3 scenarios tested in the following section. We discuss the
3 benefits and limitations of the approach and how it could potentially be modified for
4 improved performance.

5 **6.6.1 Scenario 1**

6 The DA framework was applied to 3 different scenarios totalling 5 IZ. The first scenario
7 tested was an IZ centred on the city of Sukkur. Sukkur is located just south of a large
8 barrage, used to control flood waters. Significant flooding was observed locally in the
9 Sukkur region (Sky News, 2022), however the dense city centre means that flooding is
10 difficult to detect using Sentinel-1 imagery at 20 m spatial resolution. Around one third
11 of the IZ is masked by the GFM process (Fig. 6.4(c)) but high flood likelihood values
12 are visible across some areas of the IZ (Fig. 6.3(a)). The aim is to test whether the DA
13 framework can select a flood map from the simulation library based on limited usable SAR
14 data.

15
16 The value of the cost function $J(\mathbf{x})$ from eqn. (4), Section 6.4.1 is plotted against
17 the RP value of each flood map from the simulation library in Figure 6.4(a). The cost
18 function was minimised at the lowest RP flood map (3 years) and we found that a ‘no
19 flood’ map gave a slightly lower value of $J(\mathbf{x})$. In this instance, the lower RP flood maps
20 over-estimated flooding in areas where low flood likelihood values were derived from SAR.
21 The influence of the Sukkur barrage and river canals running across the IZ made the
22 hydrodynamic modelling more difficult. Also, the flood maps do not include local defence
23 information and the flood map interpolation process is highly uncertain at RP less than 20
24 years. The results mean that no flood map was triggered following the DA (Fig. 6.4(b and
25 c)) with the CSM map (Fig. 6.4(d)) indicating where the flooding was underestimated,

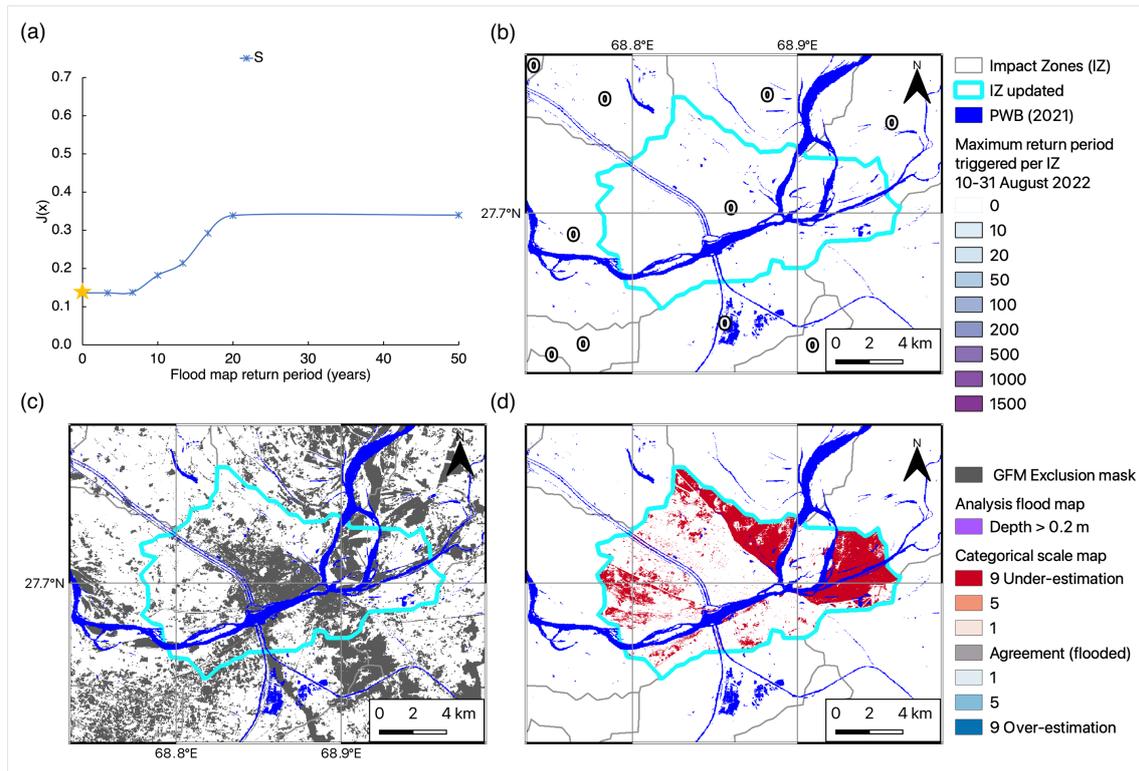


Figure 6.4: Scenario 1: Sukkur, a dense urban area. (a) The cost function (eqn. (4)) plotted against the RP value, the yellow star highlights the minimum value of $J(x)$. (b) The analysis RP triggered following DA. (c) The analysis flood extent map (note that no map was triggered for Sukkur) and (d) the CSM comparing the analysis flood extent map against Sentinel-2 NDWI.

1 particularly upstream of the Sukkur Barrage, with no flood map selected.

2 6.6.2 Scenario 2

3 The second scenario focused on a mixed urban and rural area with 2 IZ chosen around
 4 Larkana city. The dense urban area is again masked and is split across the 2 IZ (Fig.
 5 6.3(a)), but there are large unmasked areas with high and low flood likelihood values.
 6 The assimilation results for scenario 2 are plotted in Figure 6.5(a) where the cost function
 7 minimum value is similar for both LN and LS with a 7 year RP flood map triggered for LN
 8 and a 13 year RP flood map triggered for LS (Fig. 6.5(b)). The resultant analysis flood
 9 maps selected (Fig. 6.5(c)) also indicate flooding within Larkana city, overlapping the

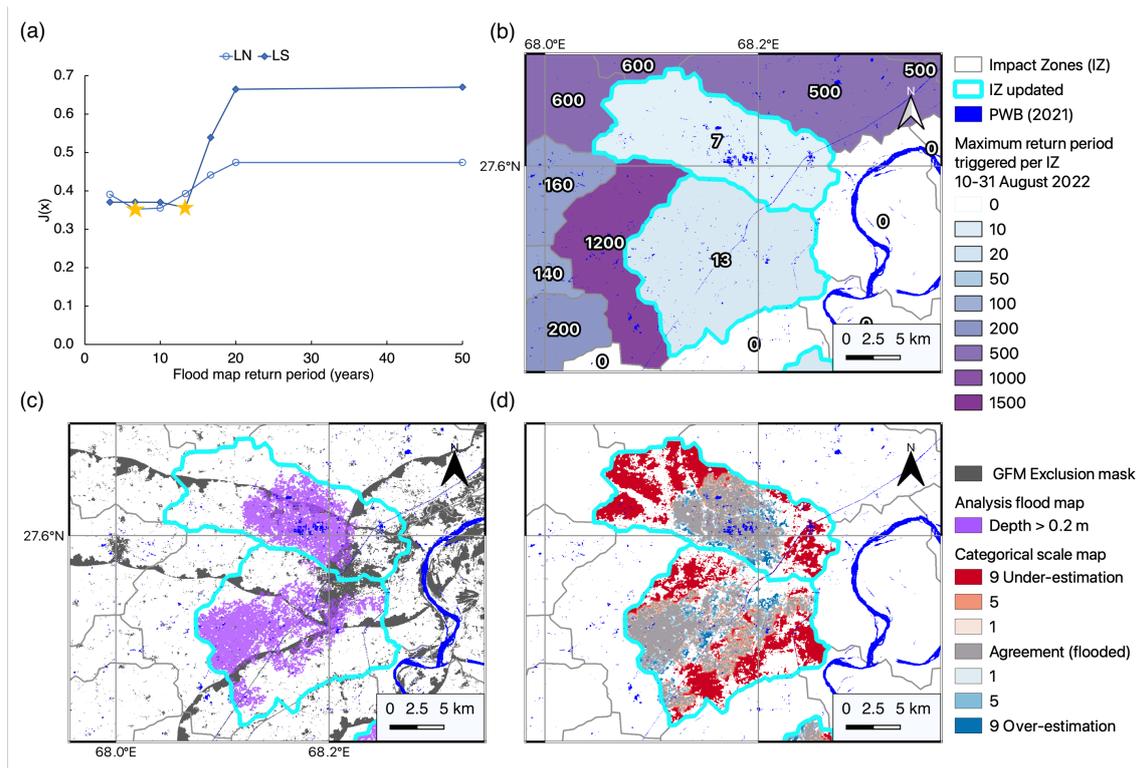


Figure 6.5: Scenario 2: Larkana, a mixed urban and rural area. (a) The cost function (eqn. (4)) plotted against the RP value, the yellow star highlights the minimum value of $J(\mathbf{x})$. (b) The analysis RP triggered following DA. (c) The analysis flood extent map and (d) the CSM comparing the analysis flood extent map against Sentinel-2 NDWI.

1 masked area. This is consistent with local observations and is important for population
 2 impact calculations for FbF schemes. The CSM indicates that whilst a large area is now
 3 correctly indicating flooding there are also large areas that are under-estimated by the
 4 analysis flood map (Fig. 6.5(d)). The neighbouring IZ that were triggered by the forecast
 5 system (Fig. 6.5(b)) are at much higher RP thresholds than the ones selected following
 6 the DA. By inspecting higher RP flood maps than those selected by the DA (for LN and
 7 LS) it became clear that these were over-estimating flooding in locations where low flood
 8 likelihood values were located causing $J(\mathbf{x})$ (Fig. 6.5(a)) to increase.

9

10 One potential solution to the inconsistency seen across the domain could be overcome

1 by including information from the forecast system. The assimilation process could be
2 carried out across a region including multiple IZ at the same time, rather than consider-
3 ing individual IZ. Conditions could be imposed, such as a consistent flood depth across
4 IZ boundaries away from the flood edge. One way to impose some smoothness would
5 be through the use of a background error term. Note that the background error is the
6 prior or forecast error. Information from neighbouring IZ could be spread across a do-
7 main by the background error covariance (\mathbf{B}) matrix used in variational DA (Bannister,
8 2008). The \mathbf{B} matrix could be calculated offline using the content of the simulation library.
9

10 Once the entire IZ becomes inundated at a 20-year RP, $J(\mathbf{x})$ remains constant with
11 increasing RP (Fig. 6.5(a)). Although the depth values are increasing, there are no
12 significant changes in flood extent possible across the IZ, meaning the cost function cannot
13 distinguish between flood maps over a 20-year RP. For the IZ tested here, the minimum
14 has already occurred at lower RP, but it is possible that the minimum could occur where
15 $J(\mathbf{x})$ is constant, meaning a range of potential RP are possible solutions. In order to
16 distinguish between equally plausible flood maps, additional observation data would be
17 required to measure flood depth. Flood depth data for sufficiently large floods could be
18 derived from satellite altimetry data such as the Surface Water and Ocean Topography
19 (SWOT) mission (Frasson et al., 2019; de Moraes Frasson et al., 2023).

20 **6.6.3 Scenario 3**

21 The final scenario investigates the impact of assimilating SAR-derived flood likelihood
22 data on flood map selection where the flood edge lies within the IZ. The cost function
23 value of $J(\mathbf{x})$ drops relatively sharply for FE1 in the north (Fig. 6.6(a)) from 0.56 at 3
24 years RP to a minimum of 0.34 at 17 years RP. Further south, $J(\mathbf{x})$ for FE2 is initially
25 lower at 0.21 at 3 years RP, gradually decreasing to a minimum of 0.17 at 20 years RP. The
26 shape of the cost function shows a smoother descent to a minimum compared to scenarios

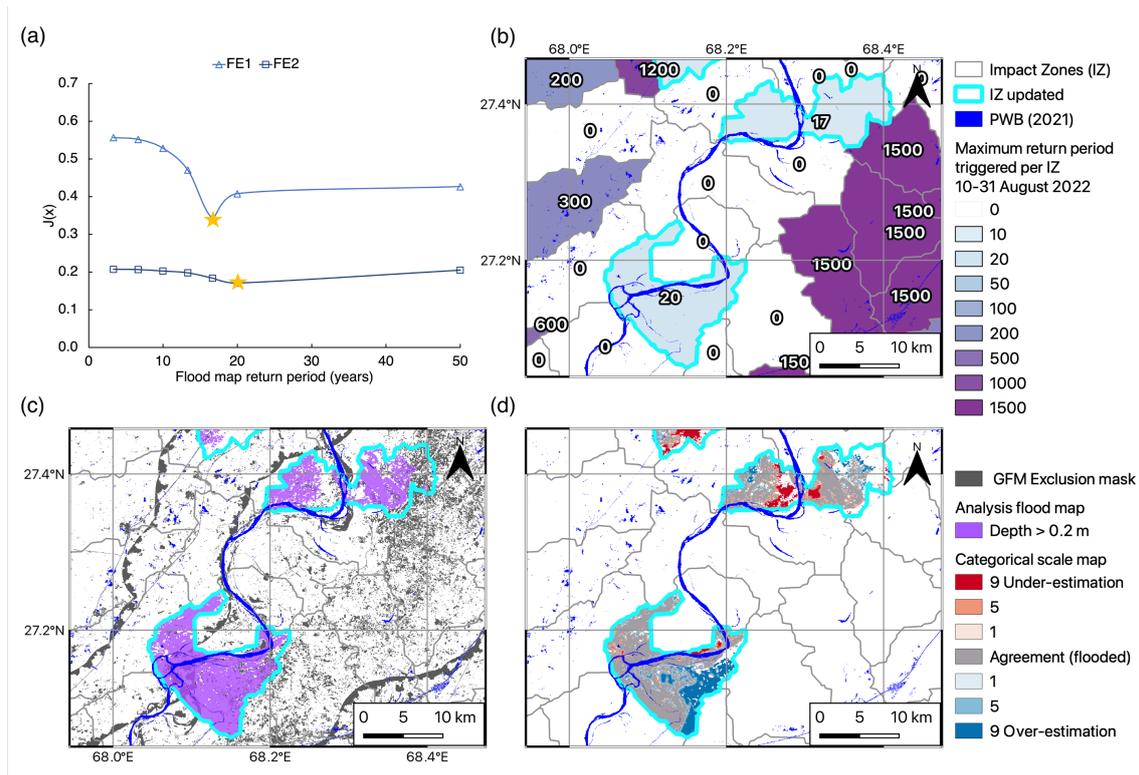


Figure 6.6: Scenario 3: Flood edge location. (a) The cost function (eqn. (4)) plotted against the RP value, the yellow star highlights the minimum value of $J(\mathbf{x})$. (b) The analysis RP triggered following DA. (c) The analysis flood extent map and (d) the CSM comparing the analysis flood extent map against Sentinel-2 NDWI.

1 1 and 2 as there is more variation in flood extent between flood maps at different RP at
 2 the flood edge location. In similarity to scenario 2 (Larkana), neighbouring IZ are again at
 3 very high RP levels (Fig. 6.6(b)). The analysis flood map for FE1 does not reach the edge
 4 of the IZ, whereas FE2 virtually flood fills the IZ (Fig. 6.6(c)). The CSM (Fig. 6.6(d))
 5 shows that some flooding close to the main Indus River has been under-estimated in FE1
 6 but with overall good accuracy and limited over-estimation. FE2 shows over-estimation in
 7 the west but again a large area in agreement with the flood extent derived from Sentinel-2
 8 NDWI. Using flood extent observation likelihood data would be more useful where the
 9 flood edge stays within the IZ tested for the maximum RP flood map. There would be
 10 more chance of variation in the cost function value across the full range of RP flood maps.

1 In contrast, for a (near) flood filled IZ there would be less variation seen in $J(\mathbf{x})$ due to
 2 limited changes in flood extent. Future work could focus on assimilating flood edge IZ
 3 and sharing information with neighbouring IZ by including a background term to update
 4 regions closer to the river channel where the IZ are more likely to be flood filled (Section
 5 6.6.2).

6 6.6.4 Analysis flood map validation

Table 6.1: Validation skill scores for each IZ analysis flood map compared against independent Sentinel-2 NDWI

| IZ code | Analysis RP | FSS at ($n = 1$) | FSS_T | $AFSS$ | n at FSS_T |
|---------|-------------|-------------------------|---------|--------|----------------|
| S | 0 | 0 | n/a | n/a | n/a |
| LN | 7 | 0.32 | 0.64 | 0.40 | $AFSS < FSS_T$ |
| LS | 13 | 0.38 | 0.65 | 0.51 | $AFSS < FSS_T$ |
| FE1 | 17 | 0.55 | 0.66 | 0.71 | n = 35 (525 m) |
| FE2 | 20 | 0.49 | 0.63 | 0.77 | n = 15 (225 m) |

7 In Table 6.1 the FSS (Section 6.4.2) has been calculated by comparing the analysis
 8 flood map selected per IZ with the corresponding Sentinel-2 NDWI representing observed
 9 flooding, with permanent water bodies excluded from the validation. The target skill
 10 score FSS_T and the asymptotic FSS $AFSS$ are also calculated. The FSS for scenario 1
 11 (Sukkur) was 0 as no flood map was triggered. For scenario 2 (LN and LS) the FSS score
 12 at grid level $n = 1$ (0.32 and 0.38) is around half of FSS_T . Usually, by increasing the
 13 neighborhood size the value of FSS also increases, eventually exceeding FSS_T . In this
 14 case $AFSS$ (which is calculated using the fraction flooded across the IZ from both the
 15 analysis and observed flood maps) is less than FSS_T meaning the total differences in flood
 16 extent are too large for the FSS to reach or exceed FSS_T . The result of this is that there

1 is not a meaningful or skilful scale of the analysis flood maps for scenario 2. This is due
2 to the under-estimation of flood extent seen on the CSM (Fig. 6.5(d)). For scenario 3,
3 $FSS_T < AFSS$ which makes it possible to calculate a skilful scale where $FSS \geq FSS_T$.
4 For FE1 this occurs when $n = 35$ or 525 m and FE2 at $n = 15$ or 225 m confirming that
5 FE2 was the most accurate analysis. These results confirm that future work should focus
6 on flood edge IZ first during the assimilation process.

7 **6.7 Conclusions**

8 In this chapter we introduced a new DA framework to update and improve the flood map
9 selection within a flood forecasting system designed for FbF applications. Open access
10 flood likelihood data derived from satellite-based SAR is used to update the flood map
11 selection for previously non-triggered sub-catchments or IZ during a flood event. The
12 framework is tested on the catastrophic flooding in Pakistan, August 2022 for 3 scenarios.

13

14 The first scenario tested an IZ where limited useful SAR data was available due to a
15 dense urban area. This resulted in no flood map selection following the DA. The second
16 scenario, where two IZ contained a mix of urban and rural areas did trigger flood maps
17 but at low RP levels, relative to neighbouring IZ that were previously triggered. This re-
18 sulted in under-estimation of the flood extent. However, the analysis flood map included
19 flooding across parts of the city of Larkana. Information from the flood likelihood data
20 from other areas of the IZ could select a flood map that included urban flooding. This is
21 useful for FbF applications where population impacts are considered. The final scenario
22 examined flood edge locations and these gave the best results as the variation in flood
23 extent selected higher RP flood maps that were more closely matched to the validation
24 data. The skilful scale of the analysis flood maps in the flood edge IZ was 225 to 525 m.
25 Out of the 5 IZ tested, 4 resulted in a flood map selection with the dense urban area and

1 limited SAR coverage the exception. Each of these 4 IZ were non-triggered, and now they
2 are, which is beneficial, even if the extents are not perfect.

3

4 The non-triggered flood maps could be updated quickly following the production of
5 the GFM flood likelihood layer (approximately 8-hours after SAR acquisition). Although
6 observed flood extent is used in the assimilation, the flood maps selected contain depth
7 values that are already linked to a catastrophe risk model and population impact maps.
8 Therefore the analysis is suitable to inform secondary financing schemes for flood response
9 and recovery, during an event. Improvements could be made by the inclusion of prior in-
10 formation from the simulation library system. An additional background term in the data
11 assimilation framework could improve the consistency of the flood maps selected across a
12 region.

13

14 An additional benefit of our approach is that the analysis flood map could be used in a
15 feedback loop to update the river discharge (e.g. in the associated GloFAS grid cell). This
16 could be useful for hydrological model calibration or in updating the initial conditions for
17 the next forecast.

18

19 Future options for optimising the simulation library flood maps (where remote flood
20 depth observations are available) could use a conventional iterative approach (such as gra-
21 dient descent methods) to step through depth values within each individual grid cell. The
22 resultant analysis depth map would represent the best estimate of the *true* flood extent and
23 could be used to inform secondary insurance payments. However, the analysis flood map
24 would no longer be consistent with the simulation library system and it would be difficult
25 to use this to update the hydrological model (i.e. in a feedback loop). The analysis will
26 possibly include surface water flooding (SWF), which is likely for large relatively flat river
27 basins where monsoon rainfall contributes to flooding, such as the Indus basin in Pakistan

1 and the Ganges-Brahmaputra-Meghna catchments of India and Bangladesh. The analysis
2 is more likely to represent flooding as observed ‘from the ground’, which makes it more
3 consistent with locally observed flooding. This combined SWF-fluvial analysis flood map
4 would mean that secondary insurance payments are more fairly distributed as they do not
5 depend on the type of flooding mapped within the fluvial simulation library. The inclusion
6 of SWF in the analysis flood map causes additional inconsistencies with the fluvial flood
7 forecasting system, which makes the analysis map less useful for updating the system in
8 a feedback loop ahead of the next forecast.

9
10 In this application of the DA framework to selected IZ we were able to analyse the
11 entire flood map library. For operational applications across a wider area, optimal iteration
12 methods could be used to save computation time and storage.

13 **6.8 Chapter summary**

14 In this chapter we address some of the limitations of simulation library flood forecasting
15 systems described in Chapter 5 by creating an assimilation framework using satellite flood
16 likelihood data to improve the flood map selection for non-triggered flood maps. We ap-
17 plied the framework to three different scenarios where the flood maps were not previously
18 triggered following flooding in Pakistan in 2022: a large city area with limited SAR flood
19 detection (due to a dense urban area); sub-catchments with mixed rural and urban areas;
20 and flood edge sub-catchments. The updated flood maps were evaluated against inde-
21 pendent satellite data using the scale-selective verification methods introduced in Chapter
22 3.

1 Chapter 7

2 Conclusions

3 7.1 Main conclusions

4 In this thesis we have applied scale-selective verification metrics in a new application to
5 evaluate several aspects of flood map forecasts from a simulation library system such
6 as the spatial accuracy and the ensemble spatial spread-skill. The forecast maps were
7 validated against SAR-derived observations of flooding for case studies both in the UK
8 and internationally. We applied a new scale-selective verification approach to a multi-
9 system comparison and addressed some of the limitations found by developing a new data
10 assimilation (DA) framework to improve the flood map analysis. The main conclusions,
11 which answer the four questions posed in Chapter 1, are:

12 **1. What are the skilful spatial scales in flood inundation forecasts made**
13 **using a simulation library approach?**

- 14 • In Chapter 3 we describe a scale-selective approach to evaluate forecast flood
15 inundation maps. A skilful spatial scale for forecast flood maps was determined
16 by calculating the Fraction Skill Score, a validation metric, found by comparing
17 the forecast flood maps against a satellite SAR-derived observation of flooding

1 across a range of neighbourhood sizes. A target skill score was calculated and
2 this depends on the magnitude of the observed flood. We found that the skilful
3 spatial scale determined for the flood edge was more sensitive to changes in
4 spatial accuracy and spatial scale compared to the skilful scale found by evalu-
5 ating the entire flood extent. Categorical scale maps indicated that the skilful
6 scale varies with location within a domain. Based on the results from evaluat-
7 ing a forecast of severe flooding on the River Wye and the River Lugg (UK)
8 in February 2020, we found that the skilful scales depend on forecast inputs
9 and can be linked to catchment characteristics such as flood plain topography,
10 land use type and urban infrastructure. There are multiple operational uses of
11 scale-selective evaluation of forecast flood maps such as model development and
12 improvement and determining a meaningful spatial scale at which to present
13 the forecast flood maps to end users.

14 **2. How skilfully does an ensemble of forecast flood maps represent the spa-**
15 **tial uncertainty within the flood forecast?**

- 16 • In Chapter 4 we presented a new scale-selective approach to assess the spatial
17 predictability and spread-skill of an ensemble flood map forecast that accounts
18 for the individual spatial prediction of flood extent held within each ensemble
19 member flood map. The method determines, at specific locations within the
20 domain, whether the ensemble forecast is over-, under- or well-spread. The
21 spatial spread-skill relationship was mapped onto our new spatial spread-skill
22 map. Results following the application of the method to an ensemble forecast
23 of flooding on the Brahmaputra (Assam, India) in August 2017 show that
24 the spatial spread-skill relationship is highly spatially variable and depends on
25 multiple factors throughout the flood forecasting system chain. We found that
26 one ensemble member flood map outperformed all others including summary

1 flood maps such as the ensemble median and a total combined ensemble.

2 **3. How useful are scale-selective evaluation approaches when applied to mul-**
3 **ti-ple flood forecasting systems?**

4 • In Chapter 5 we investigated a new application of scale-selective verification by
5 evaluating the performance of three flood forecasting systems. Two simulation
6 library systems, Flood Foresight (30 m) and GloFAS Rapid Flood Mapping
7 (1000 m) and one hydrodynamically modelled system, the Bangladesh FFWC
8 Super Model (300 m), all made predictions of flood extent at different spatial
9 scales (grid lengths, shown in brackets) for the Jamuna River flood, Bangladesh,
10 July 2020. Our results show that the simulation library system accuracy crit-
11 ically depends on the discharge return period threshold set to trigger a flood
12 map selection and the number of hydrological model ensemble members that
13 must exceed it. At short forecast lead times, the Super Model outperforms the
14 other systems in three out of four districts. Near to the Bangladesh border,
15 the trans-boundary benefits of the two global systems are evident, with both
16 outperforming the local model. We conclude that a scale-selective verification
17 approach can quantify the skill of systems operating at different spatial scales
18 so that their benefits and limitations can be evaluated.

19 **4. Does a data assimilation framework improve the analysis of flood inun-**
20 **duction from a simulation library system?**

21 • In Chapter 6 a DA framework was developed to integrate probabilistic flood
22 extent maps from satellite-based SAR sensors into the simulation library flood
23 map selection process. The method was tested on the severe flood event in
24 Pakistan, 2022, where several sub-catchments resulted in a non-trigger of the
25 forecast-based financing system deployed here, despite significant flooding evi-

1 dent from earth observation data. The DA successfully triggered flood maps in
2 4 out of 5 sub-catchments tested and we found that evaluating sub-catchments
3 over the flood edge gave the best results. The analysis flood map remains con-
4 sistent with the forecast-based financing system so could be easily linked to
5 impacts to inform secondary finance payments during an event.

6 **7.2 Thesis synthesis**

7 The overall aim of this thesis was to improve flood inundation forecasts using satellite
8 derived observations of flooding. It is important to understand meaningful length scales
9 when comparing spatial observations with forecasts. Assimilating observations at grid
10 level could lead to over-fitting or indicating over-confidence in the forecast. The forecast
11 improvement may suffer from limited persistence through time. Finding meaningful or
12 skilful spatial scales was the focus of Chapter 1. Visualising the skilful scale on an agree-
13 ment scale map, brought additional benefits to the model evaluation, since improvements
14 could be targeted to specific locations. The scale-selective evaluation approach overcomes
15 issues such as the double penalty impact of high resolution verification and the impact of
16 flood magnitude on skill scores. The method developed was used as an evaluation tool in
17 Chapters 3 to 6.

18
19 In Chapter 4 the scale-selective verification approach developed in Chapter 3 was ex-
20 tended to an ensemble forecast system. The ensemble verification approach adds an extra
21 dimension, the ensemble spread. The spread-skill relationship was evaluated for a simula-
22 tion library system predicting flood inundation. The results showed that the most skilful
23 member may lie in an ensemble outlier and the ensemble mean flood map may miss local
24 detail. These are important considerations for interpreting ensemble flood maps in fore-
25 cast applications.

1

2 The scale-selective verification approach has been applied to a simulation library sys-
3 tem in Chapters 3 and 4. In Chapter 5 we also verify a local flood forecasting system
4 along with two simulation library systems to show how the verification approach can be
5 useful in multi-system evaluations where the flood maps are presented at different spatial
6 scales. The skilful scale can be interpreted as a displacement distance, which can be di-
7 rectly compared across the systems. The simulation library system limitations found in
8 Chapter 5 gave motivation to develop the data assimilation framework in Chapter 6.

9

10 In some forecast situations the forecast discharge driving the flood map selection may
11 not exceed the return period threshold set to trigger the flood map selection. This can
12 result in missed opportunities for forecast-based financing where flooding did occur but
13 was not well forecast. The data assimilation approach developed in Chapter 6 aimed to
14 overcome this by using satellite-derived flood likelihood data to improve the flood map
15 selection from the simulation library, creating a new analysis flood map. The resultant
16 analysis flood map was evaluated using the approaches developed earlier in Chapter 1.

17 **7.3 Limitations**

18 The satellite-derived observations of flooding from SAR data used for verification in Chap-
19 ters 3, 4, and 5 make the assumption that the observed flooding extent is accurate. The
20 limitations of SAR-derived flooding are discussed in each Chapter. However, new meth-
21 ods such as masking areas where SAR is unable to reliably detect flooding, described in
22 Chapter 6 would have benefited the results in previous chapters. The inclusion of flood de-
23 fences in the forecast flood maps would also give a better evaluation of model performance.

24

25 The flood inundation evaluation in Chapter 1 was developed on a conventional grid

1 across a region of interest. The neighbourhood approach may not be meaningful where
2 the neighbourhood extends beyond a sub-catchment boundary, where flood waters are
3 contained within each sub-catchment. The evaluation and data assimilation could be car-
4 ried out on a sub-catchment level, rather than across a domain. We evaluated the spatial
5 extent of flooding. Flood depth evaluation would add an extra dimension to the analysis,
6 however this is complicated by the need for a detailed digital terrain model. Extrapolating
7 water levels across a region from observations at the flood edge creates large uncertainties
8 in the flood depth estimation. More recent techniques may reduce these uncertainties
9 (Amitrano et al., 2024; Betterle & Salamon, 2024) so that agreement scales could be cal-
10 culated for flood depths.

11

12 In Chapter 4 the spread-skill relationship was evaluated in detail for one forecast lead
13 time. At the time, this was the only data available. It would be beneficial to evaluate
14 the spread-skill at longer forecast lead times and for different flooding events. The multi-
15 system comparison in Chapter 5 was limited by data availability, the study would benefit
16 from additional data from the local model at more lead times to compare against the
17 simulation library systems. The data assimilation framework in Chapter 6 uncovered
18 occasions where the sub-catchments became flood filled and distinguishing between flood
19 maps proved difficult. Recommendations were made in Chapter 7 and the next section to
20 address this limitation.

21 **7.4 Recommendations for future work**

22 The application of scale-selective verification methods (Chapter 3) and the ensemble
23 spatial-spread skill methods (Chapter 4) could be applied to other flood events. The
24 case studies might include additional emphasis on the physical characteristics of the flood-
25 ing or the driving meteorological situation, e.g., do the antecedent soil moisture conditions

1 impact the skilful scale of the forecast? This relies on the availability of a historical cata-
2 logue of events with flood inundation predictions. It would be interesting to investigate the
3 impact of forecast lead time or an increased spatial resolution of the driving model, on the
4 ensemble spatial-spread skill and the probabilistic skilful spatial scale (Necker et al., 2023).

5
6 An important aspect of developing an inundation flood forecasting system is to deter-
7 mine the most useful way to present both deterministic and ensemble flood map forecasts.
8 Future work might include investigating the use of presenting forecasts to end-users using
9 variable spatial scales, so that the scale reflects the forecast uncertainty. There is also
10 the potential to smooth and contour the flood edge, particularly in rural areas, so that it
11 is more representative of a flooding situation. Presentation options of ensemble forecasts
12 could be investigated such as spatially clustering similar ensemble members into groups
13 of flood extent or to present a most likely, best and worst case scenario ensemble flood map.

14
15 Several recommendations were made at the end of Chapter 6 for future data assimila-
16 tion approaches applied to simulation library systems. An additional background term in
17 the data assimilation framework could improve the consistency of the flood maps selected
18 across a region. A recent paper assimilates precipitation features based on the FSS. Re-
19 sults based on synthetic forecast experiments showed that the proposed method improved
20 the forecast accuracy (Otsuka et al., 2023). Future flood inundation assimilation could
21 incorporate the skilful scale found using the scale-selective methods presented in Chapter
22 3. Assimilating at a coarser scale would save on computation time and reduce the risk of
23 over-fitting. For example, if the flood maps were found to be skilful at $n = 5$, the fraction
24 flooded at $n = 3$ could be used to improve the flood map selection from the simulation
25 library.

1 References

- 2 Albertini, C., Gioia, A., Iacobellis, V., & Manfreda, S. (2022). Detection of Surface Water
3 and Floods with Multispectral Satellites. Remote Sensing, 14(23). doi: 10.3390/
4 rs14236005
- 5 Alfieri, L., Burek, P., Dutra, E., Krzeminski, B., Muraro, D., Thielen, J., & Pappenberger,
6 F. (2013). GloFAS-global ensemble streamflow forecasting and flood early warning.
7 Hydrology and Earth System Sciences, 17(3), 1161–1175. doi: 10.5194/hess-17-1161
8 -2013
- 9 Alfonso, L., Mukolwe, M. M., & Di Baldassarre, G. (2016, feb). Probabilistic Flood Maps
10 to support decision-making: Mapping the Value of Information. Water Resources
11 Research, 52(2), 1026–1043. doi: 10.1002/2015WR017378
- 12 Amitrano, D., Di Martino, G., Di Simone, A., & Imperatore, P. (2024). Flood Detection
13 with SAR: A Review of Techniques and Datasets. Remote Sensing, 16(4). Retrieved
14 from <https://www.mdpi.com/2072-4292/16/4/656> doi: 10.3390/rs16040656
- 15 Anderson, S. R., Csima, G., Moore, R. J., Mittermaier, M., & Cole, S. J. (2019). Towards
16 operational joint river flow and precipitation ensemble verification: considerations
17 and strategies given limited ensemble records. Journal of Hydrology, 577(March),
18 123966. Retrieved from <https://doi.org/10.1016/j.jhydrol.2019.123966> doi:
19 10.1016/j.jhydrol.2019.123966
- 20 Annis, A., Nardi, F., & Castelli, F. (2022). Simultaneous assimilation of water levels from

- 1 river gauges and satellite flood maps for near-real-time flood mapping. Hydrology
2 and Earth System Sciences, 26(4), 1019–1041. doi: 10.5194/hess-26-1019-2022
- 3 Apel, H., Vorogushyn, S., & Merz, B. (2022). Brief communication: Impact forecasting
4 could substantially improve the emergency management of deadly floods: case study
5 July 2021 floods in Germany. Natural Hazards and Earth System Sciences, 22(9),
6 3005–3014. doi: 10.5194/nhess-22-3005-2022
- 7 Arnal, L., Anspoks, L., Manson, S., Neumann, J., Norton, T., Stephens, E., ... Cloke,
8 H. L. (2020). "Are We talking just a bit of water out of bank? or is it Armageddon?"
9 Front line perspectives on transitioning to probabilistic fluvial flood forecasts in
10 England. Geoscience Communication, 3(2), 203–232. doi: 10.5194/gc-3-203-2020
- 11 ASDMA. (2017). Assam state disaster management authority flood alert. [https://](https://asdma.assam.gov.in/information-services/assam-flood-report)
12 asdma.assam.gov.in/information-services/assam-flood-report, last access
13 10th November 2021.
- 14 Bannister, R. N. (2008). Review A review of forecast error covariance statistics in atmo-
15 spheric variational data assimilation. I: Characteristics and measurements of fore-
16 cast error covariances. Q. J. R. Meteorol. Soc., 134, 1951–1970. Retrieved from
17 www.interscience.wiley.com doi: 10.1002/qj.339
- 18 Bannister, R. N. (2017). A review of operational methods of variational and ensemble-
19 variational data assimilation. Quarterly Journal of the Royal Meteorological Society,
20 143(703), 607–633. doi: 10.1002/qj.2982
- 21 Bates, P. D. (2022). Flood Inundation Prediction. Annual Review of Fluid
22 Mechanics, 54(1), 287–315. Retrieved from [https://doi.org/10.1146/annurev](https://doi.org/10.1146/annurev-fluid-030121-113138)
23 [-fluid-030121-113138](https://doi.org/10.1146/annurev-fluid-030121-113138) doi: 10.1146/annurev-fluid-030121-113138
- 24 Bauer-Marschallinger, B., Cao, S., Tupas, M. E., Roth, F., Navacchi, C., Melzer, T., ...
25 Wagner, W. (2022). Satellite-Based Flood Mapping through Bayesian Inference from
26 a Sentinel-1 SAR Datacube. Remote Sensing, 14(15), 1–28. doi: 10.3390/rs14153673
- 27 Ben Bouallègue, Z., & Theis, S. E. (2014). Spatial techniques applied to precipitation en-

- 1 semble forecasts: From verification results to probabilistic products. Meteorological
2 Applications, 21(4), 922–929. doi: 10.1002/met.1435
- 3 Bernard, A., Long, N., Becker, M., Khan, J., & Fanchette, S. (2022). Bangladesh’s vul-
4 nerability to cyclonic coastal flooding. Natural Hazards and Earth System Sciences,
5 22(3), 729–751. Retrieved from [https://nhess.copernicus.org/articles/22/](https://nhess.copernicus.org/articles/22/729/2022/)
6 729/2022/ doi: 10.5194/nhess-22-729-2022
- 7 Betterle, A., & Salamon, P. (2024). Water depth estimate and flood extent enhance-
8 ment for satellite-based inundation maps. Natural Hazards and Earth System
9 Sciences Discussions, 2024, 1–21. Retrieved from [https://nhess.copernicus.org/](https://nhess.copernicus.org/preprints/nhess-2024-22/)
10 doi: 10.5194/nhess-2024-22
- 11 Beven, K. (2016). Facets of uncertainty: Epistemic uncertainty, non-stationarity, likeli-
12 hood, hypothesis testing, and communication. Hydrological Sciences Journal, 61(9),
13 1652–1665. Retrieved from <http://dx.doi.org/10.1080/02626667.2015.1031761>
14 doi: 10.1080/02626667.2015.1031761
- 15 Boelee, L. (2022). Evaluation of global flood forecasts in ungauged catchments (Doctoral
16 dissertation, The Open University, UK). Retrieved from [http://oro.open.ac.uk/](http://oro.open.ac.uk/84856/)
17 84856/
- 18 Boelee, L., Lumbroso, D. M., Samuels, P. G., & Cloke, H. L. (2019). Estimation of
19 uncertainty in flood forecasts—A comparison of methods. Journal of Flood Risk
20 Management. doi: 10.1111/jfr3.12516
- 21 Bouttier, F., & Courtier, P. (2002). Data assimilation concepts and methods march 1999.
22 Training.
- 23 Bradbrook, K. (2006). JFLOW: A multiscale two-dimensional dynamic flood model.
24 Water and Environment Journal. doi: 10.1111/j.1747-6593.2005.00011.x
- 25 Breiman, L., Friedman, J., Stone, C. J., & Olshen, R. (1984). Classification and Regression
26 Trees. Chapman and Hall/CRC.
- 27 Briand, T., & Monasse, P. (2018). Theory and practice of image B-spline interpolation.

- 1 Image Processing On Line, 8, 99–141. doi: 10.5201/ipol.2018.221
- 2 Briggs, W. M., & Levine, R. A. (1997). Wavelets and field forecast verification. Monthly
3 Weather Review, 125(6), 1329–1373. doi: 10.1175/1520-0493(1997)125<1329:waffv>
4 2.0.co;2
- 5 Buizza, R. (1997). Potential forecast skill of ensemble prediction and spread and skill
6 distributions of the ECMWF ensemble prediction system. Monthly Weather Review,
7 125(1), 99–119. doi: 10.1175/1520-0493(1997)125<0099:PFSOEP>2.0.CO;2
- 8 BWDB. (2020). Bangladesh water development board annual flood report 2020. [http://](http://ffwc.gov.bd/)
9 ffwc.gov.bd/, last access 26th October 2022.
- 10 Casati, B., & Wilson, L. J. (2007). A new spatial-scale decomposition of the brier score:
11 Application to the verification of lightning probability forecasts. Monthly Weather
12 Review, 135(9), 3052–3069. doi: 10.1175/MWR3442.1
- 13 Central Water Commission. (2023). Flood forecast dashboard. [https://cwc.gov.in/](https://cwc.gov.in/ffm_dashboard)
14 [ffm_dashboard](https://cwc.gov.in/ffm_dashboard), last access 23rd January 2023.
- 15 Chen, X., Yuan, H., & Xue, M. (2018). Spatial spread-skill relationship in terms of agree-
16 ment scales for precipitation forecasts in a convection-allowing ensemble. Quarterly
17 Journal of the Royal Meteorological Society, 144(710), 85–98. doi: 10.1002/qj.3186
- 18 Chini, M., Hostache, R., Giustarini, L., & Matgen, P. (2017). A hierarchical split-based
19 approach for parametric thresholding of SAR images: Flood inundation as a test
20 case. IEEE Transactions on Geoscience and Remote Sensing, 55(12), 6975–6988.
21 doi: 10.1109/TGRS.2017.2737664
- 22 Cloke, H. L., & Pappenberger, F. (2008). Evaluating forecasts of extreme events for hy-
23 drological applications: an approach for screening unfamiliar performance measures.
24 Meteorological Applications, 15(1), 181–197. doi: <https://doi.org/10.1002/met.58>
- 25 Cloke, H. L., & Pappenberger, F. (2009). Ensemble flood forecasting: A review. Journal
26 of Hydrology, 375(3-4), 613–626. Retrieved from [http://dx.doi.org/10.1016/](http://dx.doi.org/10.1016/j.jhydrol.2009.06.005)
27 [j.jhydrol.2009.06.005](http://dx.doi.org/10.1016/j.jhydrol.2009.06.005) doi: 10.1016/j.jhydrol.2009.06.005

- 1 Cooper, E. S., Dance, S. L., Garcia-Pintado, J., Nichols, N. K., & Smith, P. J. (2018, jun).
2 Observation impact, domain length and parameter estimation in data assimilation
3 for flood forecasting. *Environmental Modelling & Software*, *104*, 199–214. doi:
4 10.1016/J.ENVSOFT.2018.03.013
- 5 Cooper, E. S., Dance, S. L., García-Pintado, J., Nichols, N. K., & Smith, P. J. (2019).
6 Observation operators for assimilation of satellite observations in fluvial inundation
7 forecasting. *Hydrology and Earth System Sciences*, *23*(6), 2541–2559. doi: 10.5194/
8 hess-23-2541-2019
- 9 Copernicus Programme. (2021). *Copernicus Emergency Management Service*. [https://
10 emergency.copernicus.eu/](https://emergency.copernicus.eu/), last access 14 September 2021.
- 11 Coughlan de Perez, E., Berse, K. B., Depante, L. A. C., Easton-Calabria, E., Evidente,
12 E. P. R., Ezike, T., ... Sant, C. V. (2022). Learning from the past in moving
13 to the future: Invest in communication and response to weather early warnings to
14 reduce death and damage. *Climate Risk Management*, *38*, 100461. Retrieved from
15 <https://doi.org/10.1016/j.crm.2022.100461> doi: 10.1016/j.crm.2022.100461
- 16 Coughlan de Perez, E., Van Den Hurk, B., Van Aalst, M. K., Amuron, I., Bamanya,
17 D., Hauser, T., ... Zsoter, E. (2016). Action-based flood forecasting for triggering
18 humanitarian action. *Hydrology and Earth System Sciences*, *20*(9), 3549–3560. doi:
19 10.5194/hess-20-3549-2016
- 20 Coughlan de Perez, E., Van Den Hurk, B., Van Aalst, M. K., Jongman, B., Klose, T., &
21 Suarez, P. (2015). Forecast-based financing: An approach for catalyzing humani-
22 tarian action based on extreme weather and climate forecasts. *Natural Hazards and
23 Earth System Sciences*, *15*(4), 895–904. doi: 10.5194/nhess-15-895-2015
- 24 Dance, S. L., Dalton, H., Carolin, C., Clark, J., Ferrando Jorge, N., Hooker, H., & Mason,
25 D. (2022). *Assimilated watercolours: Pop up art exhibitions in care homes, egu
26 general assembly 2022, vienna, austria, 23–27 may 2022, egu22-11694*. Retrieved
27 from <https://doi.org/10.5194/egusphere-egu22-11694>

- 1 Dasgupta, A., Grimaldi, S., Ramsankaran, R. A., Pauwels, V. R., & Walker, J. P.
2 (2018). Towards operational SAR-based flood mapping using neuro-fuzzy texture-
3 based approaches. Remote Sensing of Environment, 215(June), 313–329. doi:
4 10.1016/j.rse.2018.06.019
- 5 Dasgupta, A., Grimaldi, S., Ramsankaran, R. A., Pauwels, V. R. N., Walker, J. P., Chini,
6 M., ... Matgen, P. (2018). Flood mapping using synthetic aperture radar sensors
7 from local to global scales. In Global flood hazard (p. 55-77). American Geophysical
8 Union (AGU). Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/9781119217886.ch4)
9 [abs/10.1002/9781119217886.ch4](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/9781119217886.ch4) doi: <https://doi.org/10.1002/9781119217886>
10 [.ch4](https://doi.org/10.1002/9781119217886.ch4)
- 11 Dasgupta, A., Hostache, R., Ramasankaran, R., Schumann, G. J., Grimaldi, S., Pauwels,
12 V. R. N., & Walker, J. P. (2021a). On the impacts of observation location, timing
13 and frequency on flood extent assimilation performance. Water Resources Research.
14 doi: 10.1029/2020wr028238
- 15 Dasgupta, A., Hostache, R., Ramsankaran, R. A., Schumann, G. J., Grimaldi, S., Pauwels,
16 V. R., & Walker, J. P. (2021b). A Mutual Information-Based Likelihood Function for
17 Particle Filter Flood Extent Assimilation. Water Resources Research, 57(2), 1–28.
18 doi: 10.1029/2020WR027859
- 19 Davies, P. A., McCarthy, M., Christidis, N., Dunstone, N., Fereday, D., Kendon, M.,
20 ... Sexton, D. (2021). The wet and stormy uk winter of 2019/2020. Weather,
21 76(12), 396-402. Retrieved from [https://rmets.onlinelibrary.wiley.com/doi/](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/wea.3955)
22 [abs/10.1002/wea.3955](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/wea.3955) doi: <https://doi.org/10.1002/wea.3955>
- 23 DAWN. (2022). Larkana faces threat of outbreak due to inundation. [https://www.dawn](https://www.dawn.com/news/1707072)
24 [.com/news/1707072](https://www.dawn.com/news/1707072), last access 15th June 2023.
- 25 de Moraes Frasson, R. P., Turmon, M. J., Durand, M. T., & David, C. H. (2023). Es-
26 timating the Relative Impact of Measurement, Parameter, and Flow Law Errors
27 on Discharge from the Surface Water and Ocean Topography Mission. Journal of

- 1 Hydrometeorology, 24(3), 425–443. Retrieved from [https://journals.ametsoc](https://journals.ametsoc.org/view/journals/hydr/24/3/JHM-D-22-0078.1.xml)
2 [.org/view/journals/hydr/24/3/JHM-D-22-0078.1.xml](https://journals/hydr/24/3/JHM-D-22-0078.1.xml) doi: [https://doi.org/](https://doi.org/10.1175/JHM-D-22-0078.1)
3 [10.1175/JHM-D-22-0078.1](https://doi.org/10.1175/JHM-D-22-0078.1)
- 4 Dey, S. R., Leoncini, G., Roberts, N. M., Plant, R. S., & Migliorini, S. (2014). A spatial
5 view of ensemble spread in convection permitting ensembles. Monthly Weather
6 Review. doi: [10.1175/MWR-D-14-00172.1](https://doi.org/10.1175/MWR-D-14-00172.1)
- 7 Dey, S. R., Plant, R. S., Roberts, N. M., & Migliorini, S. (2016). Assessing spatial
8 precipitation uncertainties in a convective-scale ensemble. Quarterly Journal of the
9 Royal Meteorological Society. doi: [10.1002/qj.2893](https://doi.org/10.1002/qj.2893)
- 10 Dey, S. R., Roberts, N. M., Plant, R. S., & Migliorini, S. (2016). A new method for the
11 characterization and verification of local spatial predictability for convective-scale
12 ensembles. Quarterly Journal of the Royal Meteorological Society. doi: [10.1002/](https://doi.org/10.1002/qj.2792)
13 [qj.2792](https://doi.org/10.1002/qj.2792)
- 14 Dhar, O. N., & Nandargi, S. (2000). A study of floods in the Brahmaputra basin in
15 India. International Journal of Climatology, 20(7), 771–781. doi: [10.1002/1097](https://doi.org/10.1002/1097-0088(20000615)20:7<771::AID-JOC518>3.0.CO;2-Z)
16 [-0088\(20000615\)20:7<771::AID-JOC518>3.0.CO;2-Z](https://doi.org/10.1002/1097-0088(20000615)20:7<771::AID-JOC518>3.0.CO;2-Z)
- 17 Dhar, O. N., & Nandargi, S. (2003). Hydrometeorological Aspects of Floods in India. In
18 M. M. Q. Mirza, A. Dixit, & A. Nishat (Eds.), Flood problem and management in
19 south asia (pp. 1–33). Dordrecht: Springer Netherlands. Retrieved from [https://](https://doi.org/10.1007/978-94-017-0137-2_1)
20 doi.org/10.1007/978-94-017-0137-2_1 doi: [10.1007/978-94-017-0137-2.1](https://doi.org/10.1007/978-94-017-0137-2_1)
- 21 Di Mauro, C., Hostache, R., Matgen, P., Pelich, R., Chini, M., Van Leeuwen, P. J., ...
22 Blöschl, G. (2021). Assimilation of probabilistic flood maps from SAR data into a
23 coupled hydrologic-hydraulic forecasting model: A proof of concept. Hydrology and
24 Earth System Sciences, 25(7), 4081–4097. doi: [10.5194/hess-25-4081-2021](https://doi.org/10.5194/hess-25-4081-2021)
- 25 Di Mauro, C., Hostache, R., Matgen, P., Pelich, R., Chini, M., van Leeuwen, P. J., ...
26 Blöschl, G. (2020). Assimilation of probabilistic flood maps from SAR data into a
27 hydrologic-hydraulic forecasting model: a proof of concept. Hydrology and Earth

- 1 System Sciences Discussions(September), 1–24. doi: 10.5194/hess-2020-403
- 2 Di Mauro, C., Hostache, R., Matgen, P., Pelich, R., Chini, M., van Leeuwen, P. J., ...
- 3 Blöschl, G. (2022, aug). A Tempered Particle Filter to Enhance the Assimila-
- 4 tion of SAR-Derived Flood Extent Maps Into Flood Forecasting Models. Water
- 5 Resources Research, 58(8). Retrieved from [https://onlinelibrary.wiley.com/](https://onlinelibrary.wiley.com/doi/10.1029/2022WR031940)
- 6 doi/10.1029/2022WR031940 doi: 10.1029/2022WR031940
- 7 Dottori, F., Salamon, P., Bianchi, A., Alfieri, L., Hirpa, F. A., & Feyen, L. (2016). Devel-
- 8 opment and evaluation of a framework for global flood hazard mapping. Advances
- 9 in Water Resources, 94, 87–102. Retrieved from [http://dx.doi.org/10.1016/](http://dx.doi.org/10.1016/j.advwatres.2016.05.002)
- 10 j.advwatres.2016.05.002 doi: 10.1016/j.advwatres.2016.05.002
- 11 Dottori, F., & Todini, E. (2011). Developments of a flood inundation model based
- 12 on the cellular automata approach: Testing different methods to improve model
- 13 performance. Physics and Chemistry of the Earth, 36(7-8), 266–280. Retrieved
- 14 from <http://dx.doi.org/10.1016/j.pce.2011.02.004> doi: 10.1016/j.pce.2011
- 15 .02.004
- 16 Dunning, P. (2019). Fly. [https://www.jbarisk.com/news-blogs/fly-technology](https://www.jbarisk.com/news-blogs/fly-technology-revolutionising-the-world-of-catastrophe-modelling/)
- 17 -revolutionising-the-world-of-catastrophe-modelling/, last access 20th Jan-
- 18 uary 2023.
- 19 ECMWF. (2022). Skill scores of forecasts of weather parameters by tige centres. [https://](https://www.ecmwf.int/en/forecasts/charts)
- 20 www.ecmwf.int/en/forecasts/charts, last access 20 April 2022.
- 21 Emerton, R. E., Stephens, E. M., Pappenberger, F., Pagano, T. C., Weerts, A. H., Wood,
- 22 A. W., ... Cloke, H. L. (2016). Continental and global scale flood forecasting
- 23 systems. Wiley Interdisciplinary Reviews: Water, 3(3), 391–418. doi: 10.1002/
- 24 wat2.1137
- 25 Environment Agency. (2021). National LIDAR Programme. [https://](https://data.gov.uk/dataset/f0db0249-f17b-4036-9e65-309148c97ce4/national-lidar-programme)
- 26 [data.gov.uk/dataset/f0db0249-f17b-4036-9e65-309148c97ce4/national](https://data.gov.uk/dataset/f0db0249-f17b-4036-9e65-309148c97ce4/national-lidar-programme)
- 27 -lidar-programme, last access 29 April 2021.

- 1 ESA. (2021). ICEYE commercial satellites join the EU Copernicus programme.
2 [https://www.esa.int/Applications/Observing_the_Earth/Copernicus/
3 ICEYE_commercial_satellites_join_the_EU_Copernicus_programme](https://www.esa.int/Applications/Observing_the_Earth/Copernicus/ICEYE_commercial_satellites_join_the_EU_Copernicus_programme), last ac-
4 cess 28 October 2021.
- 5 EU Science Hub. (2021). The JRC launches a revolutionary tool for monitoring
6 ongoing floods worldwide as part of the copernicus emergency management service.
7 [https://ec.europa.eu/jrc/en/news/jrc-launches-revolutionary-tool-for
8 -monitoring-floods-worldwide-part-copernicus-emergency-management
9 -service](https://ec.europa.eu/jrc/en/news/jrc-launches-revolutionary-tool-for-monitoring-floods-worldwide-part-copernicus-emergency-management-service), last access 28 October 2021.
- 10 Evensen, G. (1994, may). Sequential data assimilation with a nonlinear quasi-
11 geostrophic model using Monte Carlo methods to forecast error statistics.
12 Journal of Geophysical Research: Oceans, 99(C5), 10143–10162. Retrieved from
13 <https://onlinelibrary.wiley.com/doi/full/10.1029/94JC00572>[https://
14 onlinelibrary.wiley.com/doi/abs/10.1029/94JC00572](https://onlinelibrary.wiley.com/doi/abs/10.1029/94JC00572)[https://agupubs
15 .onlinelibrary.wiley.com/doi/10.1029/94JC00572](https://agupubs.onlinelibrary.wiley.com/doi/10.1029/94JC00572) doi: 10.1029/94JC00572
- 16 Farr, T. G., Rosen, P. A., Caro, E., Crippen, R., Duren, R., Hensley, S., ... Alsdorf,
17 D. (2007). The Shuttle Radar Topography Mission. Reviews of Geophysics, 45(2).
18 Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/
19 2005RG000183](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2005RG000183) doi: <https://doi.org/10.1029/2005RG000183>
- 20 Fekete, A., & Sandholz, S. (2021). Here comes the flood, but not failure? Lessons to
21 learn after the heavy rain and pluvial floods in germany 2021. Water (Switzerland),
22 13(21), 1–20. doi: 10.3390/w13213016
- 23 Floodlist. (2017). India – third wave of flooding hits assam, 2 million affected.
24 <http://floodlist.com/asia/india-assam-floods-august-2017>, last access 10th
25 November 2021.
- 26 Floodlist. (2022). Pakistan – almost 1,000 dead, 33 million affected in worst floods in a
27 decade. <https://floodlist.com/asia/pakistan-floods-update-august-2022>,

- 1 last access 23rd May 2023.
- 2 Frasson, R. P. d. M., Schumann, G. J.-P., Kettner, A. J., Brakenridge, G. R., & Krajewski,
3 W. F. (2019). Will the Surface Water and Ocean Topography (SWOT) Satellite
4 Mission Observe Floods? Geophysical Research Letters, 46(17-18), 10435–10445.
5 Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019GL084686)
6 [2019GL084686](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019GL084686) doi: <https://doi.org/10.1029/2019GL084686>
- 7 Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Shukla, S., ... Michaelsen,
8 J. (2015). The climate hazards infrared precipitation with stations - A new
9 environmental record for monitoring extremes. Scientific Data, 2, 1–21. doi:
10 [10.1038/sdata.2015.66](https://doi.org/10.1038/sdata.2015.66)
- 11 Galmiche, N., Hauser, H., Spengler, T., Spensberger, C., Brun, M., & Blaser, N. (2021).
12 Revealing Multimodality in Ensemble Weather Prediction. In D. Archambault,
13 I. Nabney, & J. Peltonen (Eds.), Machine learning methods in visualisation for big
14 data. The Eurographics Association. doi: [10.2312/mlvis.20211073](https://doi.org/10.2312/mlvis.20211073)
- 15 García-Pintado, J., Mason, D. C., Dance, S. L., Cloke, H. L., Neal, J. C., Freer, J., & Bates,
16 P. D. (2015, apr). Satellite-supported flood forecasting in river networks: A real case
17 study. Journal of Hydrology, 523, 706–724. doi: [10.1016/J.JHYDROL.2015.01.084](https://doi.org/10.1016/J.JHYDROL.2015.01.084)
- 18 García-Pintado, J., Neal, J. C., Mason, D. C., Dance, S. L., & Bates, P. D. (2013).
19 Scheduling satellite-based sar acquisition for sequential assimilation of water level
20 observations into flood modelling. Journal of Hydrology, 495. doi: [10.1016/j.jhydrol](https://doi.org/10.1016/j.jhydrol.2013.03.050)
21 [.2013.03.050](https://doi.org/10.1016/j.jhydrol.2013.03.050)
- 22 GFM. (2021). GloFAS Global Flood Monitoring (GFM). [https://www.globalfloods](https://www.globalfloods.eu/technical-information/glofas-gfm/)
23 [.eu/technical-information/glofas-gfm/](https://www.globalfloods.eu/technical-information/glofas-gfm/), last access 28 October 2021.
- 24 Giustarini, L., Hostache, R., Kavetski, D., Chini, M., Corato, G., Schlaffer, S., & Matgen,
25 P. (2016). Probabilistic Flood Mapping Using Synthetic Aperture Radar Data.
26 IEEE Transactions on Geoscience and Remote Sensing, 54(12), 6958–6969. doi:
27 [10.1109/TGRS.2016.2592951](https://doi.org/10.1109/TGRS.2016.2592951)

- 1 GloFAS. (2021). Glofas methods. [https://www.globalfloods.eu/general](https://www.globalfloods.eu/general-information/glofas-methods/)
2 [-information/glofas-methods/](https://www.globalfloods.eu/general-information/glofas-methods/), last access 15th November 2021.
- 3 GloFAS. (2022a). Glofas map viewer. <https://www.globalfloods.eu/>, last access 8th
4 December 2022.
- 5 GloFAS. (2022b). Glofas rapid flood mapping. [https://confluence.ecmwf.int/](https://confluence.ecmwf.int/display/CEMS/GloFAS+Rapid+Flood+Mapping+and+Rapid+Impact+Assessment)
6 [display/CEMS/GloFAS+Rapid+Flood+Mapping+and+Rapid+Impact+Assessment](https://confluence.ecmwf.int/display/CEMS/GloFAS+Rapid+Flood+Mapping+and+Rapid+Impact+Assessment),
7 last access 4th November 2022.
- 8 Google Earth Engine CART. (2021). ee.Classifier.smileCart. [https://developers](https://developers.google.com/earth-engine/apidocs/ee-classifier-smilecart)
9 [.google.com/earth-engine/apidocs/ee-classifier-smilecart](https://developers.google.com/earth-engine/apidocs/ee-classifier-smilecart), last access 29
10 April 2021.
- 11 Google Earth Engine Catalog. (2021). Sentinel Collection. [https://developers.google](https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S1_GRD)
12 [.com/earth-engine/datasets/catalog/COPERNICUS_S1_GRD](https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S1_GRD), last access 4 August
13 2021.
- 14 Google Earth Engine Scale. (2021). Image Pyramids. [https://developers.google.com/](https://developers.google.com/earth-engine/guides/scale)
15 [earth-engine/guides/scale](https://developers.google.com/earth-engine/guides/scale), last access 16 September 2021.
- 16 Gourbesville, P. (1998, 09). Mike 11 gis: interest of gis technology for conception of
17 flood protection systems. In Proceedings of the 3rd hydroinformatics conference -
18 hydroinformatics'98. Copenhagen, Denmark.
- 19 Griffith, H. V., Wade, A. J., Lavers, D. A., & Watts, G. (2020). Atmospheric river
20 orientation determines flood occurrence. Hydrological Processes, 34(23), 4547–4555.
21 doi: 10.1002/hyp.13905
- 22 Grimaldi, S. (2022). Glofas v4.0 hydrological reanalysis. [https://www.globalfloods.eu/](https://www.globalfloods.eu/news/113-glofas-v34-release-and-glofas-v40-hydrological-reanalysis-dataset-announcement/)
23 [news/113-glofas-v34-release-and-glofas-v40-hydrological-reanalysis](https://www.globalfloods.eu/news/113-glofas-v34-release-and-glofas-v40-hydrological-reanalysis-dataset-announcement/)
24 [-dataset-announcement/](https://www.globalfloods.eu/news/113-glofas-v34-release-and-glofas-v40-hydrological-reanalysis-dataset-announcement/), last access 25th November 2022.
- 25 Grimaldi, S., Li, Y., Pauwels, V. R., & Walker, J. P. (2016). Remote Sensing-
26 Derived Water Extent and Level to Constrain Hydraulic Flood Forecasting Mod-
27 els: Opportunities and Challenges. Surveys in Geophysics, 37(5), 977–1034. doi:

- 1 10.1007/s10712-016-9378-y
- 2 Hagen, A. (2003). Fuzzy set approach to assessing similarity of categorical maps.
- 3 International Journal of Geographical Information Science, 17(3), 235–249. doi:
- 4 10.1080/13658810210157822
- 5 Harrigan, S., Zsoter, E., Alfieri, L., Prudhomme, C., Salamon, P., Wetterhall, F., ...
- 6 Pappenberger, F. (2020). GloFAS-ERA5 operational global river discharge re-
- 7 analysis 1979–present. Earth System Science Data, 12(3), 2043–2060. Retrieved
- 8 from <https://essd.copernicus.org/articles/12/2043/2020/> doi: 10.5194/
- 9 essd-12-2043-2020
- 10 Havnø, K., Madsen, M. N., & Dørge, J. (1995). Mike 11 - a generalized river modelling
- 11 package. In V. P. Singh (Ed.), Computer models of watershed hydrology (p. 733-782).
- 12 Colorado, USA: Water Resources Publications.
- 13 Hirpa, F. A., Salamon, P., Beck, H. E., Lorini, V., Alfieri, L., Zsoter, E., & Dadson, S. J.
- 14 (2018). Calibration of the Global Flood Awareness System (GloFAS) using daily
- 15 streamflow data. Journal of Hydrology, 566, 595–606. doi: 10.1016/j.jhydrol.2018
- 16 .09.052
- 17 Hoch, J. M., & Trigg, M. A. (2019). Advancing global flood hazard simulations by
- 18 improving comparability, benchmarking, and integration of global flood models.
- 19 Environmental Research Letters, 14(3). doi: 10.1088/1748-9326/aaf3d3
- 20 Hooker, H., Dance, S. L., Mason, D. C., Bevington, J., & Shelton, K. (2022). Spa-
- 21 tial scale evaluation of forecast flood inundation maps. Journal of Hydrology,
- 22 128170. Retrieved from [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0022169422007399)
- 23 S0022169422007399 doi: <https://doi.org/10.1016/j.jhydrol.2022.128170>
- 24 Hooker, H., Dance, S. L., Mason, D. C., Bevington, J., & Shelton, K. (2023a). Assessing
- 25 the spatial spread–skill of ensemble flood maps with remote-sensing observations.
- 26 Natural Hazards and Earth System Sciences, 23(8), 2769–2785. Retrieved from
- 27 <https://nhess.copernicus.org/articles/23/2769/2023/> doi: 10.5194/nhess

- 1 -23-2769-2023
- 2 Hooker, H., Dance, S. L., Mason, D. C., Bevington, J., & Shelton, K. (2023b, sep).
3 A multi-system comparison of forecast flooding extent using a scale-selective ap-
4 proach. *Hydrology Research*, 54(10), 1115–1133. Retrieved from [https://doi.org/](https://doi.org/10.2166/nh.2023.025)
5 10.2166/nh.2023.025 doi: 10.2166/nh.2023.025
- 6 Hopson, T. M. (2014). Assessing the ensemble spread-error relationship. *Monthly Weather*
7 Review, 142(3), 1125–1142. doi: 10.1175/MWR-D-12-00111.1
- 8 Horritt, M. S., Mason, D. C., & Luckman, A. J. (2001). Flood boundary delin-
9 eation from synthetic aperture radar imagery using a statistical active contour
10 model. *International Journal of Remote Sensing*, 22(13), 2489–2507. doi: 10.1080/
11 01431160116902
- 12 Hossain, S. (2020). *Glofas case study: Bangladesh 2020*. [https://www.globalfloods.eu/](https://www.globalfloods.eu/get-involved/case-study-bangladesh-2020/)
13 get-involved/case-study-bangladesh-2020/, last access 18th October 2022.
- 14 Hossain, S., Cloke, H. L., Ficchi, A., Turner, A. G., & Stephens, E. M. (2021). Hy-
15 drometeorological drivers of flood characteristics in the Brahmaputra river basin
16 in Bangladesh. *Hydrology and Earth System Sciences Discussions*, 2021, 1–28.
17 Retrieved from <https://hess.copernicus.org/preprints/hess-2021-97/> doi:
18 10.5194/hess-2021-97
- 19 Hostache, R., Chini, M., Giustarini, L., Neal, J., Kavetski, D., Wood, M., ... Mat-
20 gen, P. (2018). Near-Real-Time Assimilation of SAR-Derived Flood Maps for
21 Improving Flood Forecasts. *Water Resources Research*, 54(8), 5516–5535. doi:
22 10.1029/2017WR022205
- 23 Hostache, R., Martinis, S., Bauer-Marschallinger, B., Chini, M., Chow, S.,
24 C. Cao, Pelich, R., ... Salamon, P. (2021). A first evaluation of the
25 future CEMS systematic global flood monitoring product. [https://](https://events.ecmwf.int/event/222/contributions/2274/attachments/1280/2347/Hydrological-WS-Hostache.pdf)
26 events.ecmwf.int/event/222/contributions/2274/attachments/1280/2347/
27 Hydrological-WS-Hostache.pdf, last access 4 August 2021.

- 1 Janjić, T., Bormann, N., Bocquet, M., Carton, J. A., Cohn, S. E., Dance, S. L., . . . Weston,
2 P. (2018). On the representation error in data assimilation. Quarterly Journal of
3 the Royal Meteorological Society, 144(713), 1257–1278. Retrieved from [https://](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3130)
4 rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3130 doi: [https://doi](https://doi.org/10.1002/qj.3130)
5 [.org/10.1002/qj.3130](https://doi.org/10.1002/qj.3130)
- 6 JBA. (2021). Storm Ciara, Dennis and Jorge. [https://www.jbarisk.com/](https://www.jbarisk.com/flood-services/event-response/storm-ciara-dennis-and-jorge/)
7 [flood-services/event-response/storm-ciara-dennis-and-jorge/](https://www.jbarisk.com/flood-services/event-response/storm-ciara-dennis-and-jorge/), last access
8 14 September 2021.
- 9 Kendon, M. (2020). Storm Dennis. [https://www.metoffice.gov.uk/binaries/](https://www.metoffice.gov.uk/binaries/content/assets/metofficegovuk/pdf/weather/learn-about/uk-past-events/interesting/2020/2020_03_storm_dennis.pdf)
10 [content/assets/metofficegovuk/pdf/weather/learn-about/uk-past-events/](https://www.metoffice.gov.uk/binaries/content/assets/metofficegovuk/pdf/weather/learn-about/uk-past-events/interesting/2020/2020_03_storm_dennis.pdf)
11 [interesting/2020/2020_03_storm_dennis.pdf](https://www.metoffice.gov.uk/binaries/content/assets/metofficegovuk/pdf/weather/learn-about/uk-past-events/interesting/2020/2020_03_storm_dennis.pdf), last access 29 April 2021.
- 12 Konapala, G., Kumar, S. V., & Khalique Ahmad, S. (2021, oct). Exploring Sentinel-1
13 and Sentinel-2 diversity for flood inundation mapping using deep learning. ISPRS
14 Journal of Photogrammetry and Remote Sensing, 180, 163–173. doi: 10.1016/J
15 [.ISPRSJPRS.2021.08.016](https://doi.org/10.1016/J.ISPRSJPRS.2021.08.016)
- 16 Krullikowski, C., Chow, C., Wieland, M., Martinis, S., Bauer-Marschallinger, B., Roth,
17 F., . . . Salamon, P. (2023). Estimating Ensemble Likelihoods for the Sentinel-1-
18 Based Global Flood Monitoring Product of the Copernicus Emergency Management
19 Service. IEEE Journal of Selected Topics in Applied Earth Observations and Remote
20 Sensing, 16, 6917–6930. doi: 10.1109/JSTARS.2023.3292350
- 21 Lai, X., Liang, Q., Yesou, H., & Daillet, S. (2014). Variational assimilation of remotely
22 sensed flood extents using a 2-D flood model. Hydrology and Earth System Sciences,
23 18(11), 4325–4339. doi: 10.5194/hess-18-4325-2014
- 24 Lauritzen, S. (2023). Fundamentals of Mathematical Statistics. doi: 10.1201/
25 [9781003272359](https://doi.org/10.1201/9781003272359)
- 26 Lavers, D. A., Allan, R. P., Wood, E. F., Villarini, G., Brayshaw, D. J., & Wade, A. J.
27 (2011). Winter floods in Britain are connected to atmospheric rivers. Geophysical

- 1 Research Letters, 38(23), 1–8. doi: 10.1029/2011GL049783
- 2 Lehner, B. (2014a). Hydrobasins global watershed boundaries and sub-basin delineations
3 derived from hydrosheds data at 15 second resolution. [https://www.hydrosheds](https://www.hydrosheds.org/images/inpages/HydroBASINS_TechDoc_v1c.pdf)
4 [.org/images/inpages/HydroBASINS_TechDoc_v1c.pdf](https://www.hydrosheds.org/images/inpages/HydroBASINS_TechDoc_v1c.pdf), last access 5th November
5 2022.
- 6 Lehner, B. (2014b). Hydrobasins global watershed boundaries and sub-basin delineations
7 derived from hydrosheds data at 15 second resolution. [https://www.hydrosheds](https://www.hydrosheds.org/images/inpages/HydroBASINS_TechDoc_v1c.pdf)
8 [.org/images/inpages/HydroBASINS_TechDoc_v1c.pdf](https://www.hydrosheds.org/images/inpages/HydroBASINS_TechDoc_v1c.pdf), last access 5th November
9 2022.
- 10 Lehner, B., & Grill, G. (2013). Global river hydrography and network routing: baseline
11 data and new approaches to study the world’s large river systems. Hydrological
12 Processes, 27(15), 2171-2186. Retrieved from [https://onlinelibrary.wiley.com/](https://onlinelibrary.wiley.com/doi/abs/10.1002/hyp.9740)
13 [doi/abs/10.1002/hyp.9740](https://onlinelibrary.wiley.com/doi/abs/10.1002/hyp.9740) doi: <https://doi.org/10.1002/hyp.9740>
- 14 Leutbecher, M., & Palmer, T. N. (2008, mar). Ensemble forecasting. Journal of
15 Computational Physics, 227(7), 3515–3539. doi: 10.1016/J.JCP.2007.02.014
- 16 LISFLOOD. (2022). Lisflood. [https://ec-jrc.github.io/lisflood-model/1.1](https://ec-jrc.github.io/lisflood-model/1.1_introduction_LISFLOOD/)
17 [_introduction_LISFLOOD/](https://ec-jrc.github.io/lisflood-model/1.1_introduction_LISFLOOD/), last access 4th November 2022.
- 18 Lorenc, A. C., Ballard, S. P., Bell, R. S., Ingleby, N. B., Andrews, P. L. F., Barker, D. M.,
19 ... Saunders, F. W. (2000). The met. office global three-dimensional variational
20 data assimilation scheme. Quarterly Journal of the Royal Meteorological Society,
21 126(570), 2991-3012. Retrieved from [https://rmets.onlinelibrary.wiley.com/](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.49712657002)
22 [doi/abs/10.1002/qj.49712657002](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.49712657002) doi: <https://doi.org/10.1002/qj.49712657002>
- 23 Lorenz, E. N. (1969). The predictability of a flow which possesses many scales of motion.
24 Tellus, 21(3), 289–307. doi: 10.3402/tellusa.v21i3.10086
- 25 Martinis, S., Kersten, J., & Twele, A. (2015). A fully automated TerraSAR-X based flood
26 service. ISPRS Journal of Photogrammetry and Remote Sensing, 104, 203–212.
27 Retrieved from <http://dx.doi.org/10.1016/j.isprsjprs.2014.07.014> doi: 10

- 1 .1016/j.isprsjprs.2014.07.014
- 2 Mason, D. C., Bevington, J., Dance, S. L., Revilla-Romero, B., Smith, R., Vetra-Carvalho,
3 S., & Cloke, H. L. (2021b). Improving urban flood mapping by merging synthetic
4 aperture radar-derived flood footprints with flood hazard maps. Water (Switzerland),
5 13(11). doi: 10.3390/w13111577
- 6 Mason, D. C., Dance, S. L., & Cloke, H. L. (2021a). Floodwater detection in urban areas
7 using Sentinel-1 and WorldDEM data. Journal of Applied Remote Sensing, 15(03),
8 1–22. doi: 10.1117/1.jrs.15.032003
- 9 Mason, D. C., Dance, S. L., Vetra-Carvalho, S., & Cloke, H. L. (2018). Robust al-
10 gorithm for detecting floodwater in urban areas using synthetic aperture radar
11 images. Journal of Applied Remote Sensing, 12(04), 1. Retrieved from [http://](http://centaur.reading.ac.uk/80110/8/045011f_1.pdf)
12 centaur.reading.ac.uk/80110/8/045011f_1.pdf doi: 10.1117/1.jrs.12.045011
- 13 Mason, D. C., Davenport, I. J., Neal, J. C., Schumann, G. J., & Bates, P. D. (2012). Near
14 real-time flood detection in urban and rural areas using high-resolution synthetic
15 aperture radar images. IEEE Transactions on Geoscience and Remote Sensing. doi:
16 10.1109/TGRS.2011.2178030
- 17 Mason, D. C., Horritt, M. S., Dall’Amico, J. T., Scott, T. R., & Bates, P. D. (2007).
18 Improving river flood extent delineation from synthetic aperture radar using airborne
19 laser altimetry. IEEE Transactions on Geoscience and Remote Sensing, 45(12), 3932–
20 3943. doi: 10.1109/TGRS.2007.901032
- 21 Mason, D. C., Schumann, G. J., Neal, J. C., Garcia-Pintado, J., & Bates, P. D. (2012).
22 Automatic near real-time selection of flood water levels from high resolution Syn-
23 thetic Aperture Radar images for assimilation into hydraulic models: A case study.
24 Remote Sensing of Environment. doi: 10.1016/j.rse.2012.06.017
- 25 Matthews, G., Barnard, C., Cloke, H., Dance, S. L., Jurlina, T., Mazzetti, C., & Prud-
26 homme, C. (2022). Evaluating the impact of post-processing medium-range ensem-
27 ble streamflow forecasts from the European Flood Awareness System. Hydrology

- 1 and Earth System Sciences, 26(11), 2939–2968. Retrieved from <https://hess>
2 .[copernicus.org/articles/26/2939/2022/](https://hess.copernicus.org/articles/26/2939/2022/) doi: 10.5194/hess-26-2939-2022
- 3 Met Office. (2020). Record Breaking Rainfall. [https://www.metoffice.gov.uk/](https://www.metoffice.gov.uk/about-us/press-office/news/weather-and-climate/302020/2020-winter)
4 about-us/press-office/news/weather-and-climate/302020/2020-winter
5 -february-stats, last access 29 April 2021.
- 6 Met Office. (2022). Weather Observation Website (WOW). <https://wow.metoffice.gov>
7 .uk/observations/details/20220503d5q856skshe63xibyyb96sm4oy, last access
8 4 May 2022.
- 9 Nanditha, J. S., Kushwaha, A. P., Singh, R., Malik, I., Solanki, H., Chuphal, D. S.,
10 ... Mishra, V. (2023). The Pakistan Flood of August 2022: Causes and Impli-
11 cations. Earth's Future, 11(3), e2022EF003230. Retrieved from <https://agupubs>
12 .[onlinelibrary.wiley.com/doi/abs/10.1029/2022EF003230](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2022EF003230) doi: <https://doi>
13 .org/10.1029/2022EF003230
- 14 National River Flow Archive. (2021). NRFA. <https://nrfa.ceh.ac.uk/data/station/>
15 info/55002, last access 29 April 2021.
- 16 NCAR. (2022). Weather research & forecasting model. [https://www.mmm.ucar.edu/](https://www.mmm.ucar.edu/models/wrf)
17 models/wrf, last access 14th November 2022.
- 18 ndimage.zoom. (2021). Scipy. [https://docs.scipy.org/doc/scipy/reference/](https://docs.scipy.org/doc/scipy/reference/generated/scipy.ndimage.zoom.html)
19 generated/scipy.ndimage.zoom.html, last access 21 September 2021.
- 20 Necker, T., Wolfgruber, L., Serafin, S., Dorninger, M., & Weissmann, M. (2023). How
21 to use the fractions skill score for ensemble forecast verification. EMS Annual
22 Meeting 2023, Bratislava, Slovakia, 4–8 Sep 2023, 20, 5194. Retrieved from
23 <https://meetingorganizer.copernicus.org/EMS2023/EMS2023-35.html>
- 24 NEXTmap. (2016). Nextmap world30 dsm. <https://www.intermap.com/nextmap>, last
25 access 24th January 2023.
- 26 Nguyen, T. H., Ricci, S., Piacentini, A., Fatras, C., Kettig, P., Blanchet, G., ... Bail-
27 larin, S. (2023, January). Assimilation of sar-derived flood extent observations for

- 1 improving fluvial flood forecast – a proof-of-concept. IOP Conference Series: Earth
2 and Environmental Science, 1136(1), 012018. Retrieved from [https://iopscience](https://iopscience.iop.org/article/10.1088/1755-1315/1136/1/012018)
3 [.iop.org/article/10.1088/1755-1315/1136/1/012018](https://iopscience.iop.org/article/10.1088/1755-1315/1136/1/012018) doi: 10.1088/1755-1315/
4 1136/1/012018
- 5 Nguyen, T. H., Ricci, S., Piacentini, A., Simon, E., & Rodriguez-suquet, R. (2023).
6 Gaussian anamorphosis for ensemble kalman filter analysis of sar-derived wet surface
7 ratio observations. , 1–19. Retrieved from www.opentelemac.org
- 8 OCHA. (2020). Bangladesh monsoon flooding 2020: anticipatory action
9 pilot. [https://www.unocha.org/our-work/humanitarian-financing/](https://www.unocha.org/our-work/humanitarian-financing/anticipatory-action/summary-bangladesh-pilot)
10 [anticipatory-action/summary-bangladesh-pilot](https://www.unocha.org/our-work/humanitarian-financing/anticipatory-action/summary-bangladesh-pilot), last access 18th October
11 2022.
- 12 Otsuka, S., Awazu, T., Welzbacher, C. A., Potthast, R., & Miyoshi, T. (2023). Assimilat-
13 ing Precipitation Features Based on the Fractions Skill Score: An Idealized Study
14 with an Intermediate AGCM. In S. K. Park (Ed.), Numerical weather prediction:
15 East asian perspectives (pp. 283–294). Cham: Springer International Publish-
16 ing. Retrieved from https://doi.org/10.1007/978-3-031-40567-9_11 doi:
17 10.1007/978-3-031-40567-9_11
- 18 Palash, W., Akanda, A. S., & Islam, S. (2020). The record 2017 flood in South Asia: State
19 of prediction and performance of a data-driven requisitely simple forecast model.
20 Journal of Hydrology, 589(June). doi: 10.1016/j.jhydrol.2020.125190
- 21 Pappenberger, F., Beven, K. J., Hunter, N. M., Bates, P. D., Gouweleeuw, B. T., Thie-
22 len, J., & de Roo, A. P. (2005). Cascading model uncertainty from medium range
23 weather forecasts (10 days) through a rainfall-runoff model to flood inundation pre-
24 dictions within the European Flood Forecasting System (EFFS). Hydrology and
25 Earth System Sciences, 9(4), 381–393. doi: 10.5194/hess-9-381-2005
- 26 Pappenberger, F., Cloke, H. L., Parker, D. J., Wetterhall, F., Richardson, D. S., & Thielen,
27 J. (2015). The monetary benefit of early flood warnings in Europe. Environmental

- 1 Science and Policy, 51, 278–291. Retrieved from [http://dx.doi.org/10.1016/](http://dx.doi.org/10.1016/j.envsci.2015.04.016)
2 [j.envsci.2015.04.016](http://dx.doi.org/10.1016/j.envsci.2015.04.016) doi: 10.1016/j.envsci.2015.04.016
- 3 Pappenberger, F., Frodsham, K., Beven, K., Romanowicz, R., & Matgen, P. (2007).
4 Fuzzy set approach to calibrating distributed flood inundation models using remote
5 sensing observations. Hydrology and Earth System Sciences, 11(2), 739–752. doi:
6 10.5194/hess-11-739-2007
- 7 Pekel, J. F., Cottam, A., Gorelick, N., & Belward, A. S. (2016). High-resolution mapping
8 of global surface water and its long-term changes. Nature, 540(7633), 418–422. doi:
9 10.1038/nature20584
- 10 Pujol, L., Garambois, P.-A., & Monnier, J. (2022). Multi-dimensional hydrological–
11 hydraulic model with variational data assimilation for river networks and floodplains.
12 Geoscientific Model Development, 15(15), 6085–6113. Retrieved from [https://gmd](https://gmd.copernicus.org/articles/15/6085/2022/)
13 [.copernicus.org/articles/15/6085/2022/](https://gmd.copernicus.org/articles/15/6085/2022/) doi: 10.5194/gmd-15-6085-2022
- 14 Pörtner, H.-O., Roberts, D., Adams, H., Adelekan, I., Adler, C., Adrian, R., . . . Ibrahim,
15 Z. Z. (2022). Climate change 2022: Impacts, adaptation and vulnerability [Book].
16 Cambridge, UK and New York, USA: Cambridge University Press.
- 17 Renner, M., Werner, M. G., Rademacher, S., & Sprokkereef, E. (2009). Verification of
18 ensemble flow forecasts for the River Rhine. Journal of Hydrology, 376(3-4), 463–
19 475. Retrieved from <http://dx.doi.org/10.1016/j.jhydrol.2009.07.059> doi:
20 10.1016/j.jhydrol.2009.07.059
- 21 Rentschler, J., Salhab, M., & Jafino, B. A. (2022). Flood exposure and poverty in 188
22 countries. Nature Communications, 13. doi: 10.1038/s41467-022-30727-4
- 23 Revilla-Romero, B., Shelton, K., Wood, E., Berry, R., Bevington, J., Hankin, B., . . .
24 Huyck, C. (2017, April). Flood Foresight: A near-real time flood monitoring and
25 forecasting tool for rapid and predictive flood impact assessment. In Egu general
26 assembly conference abstracts (p. 1230).
- 27 Revilla-Romero, B., Wanders, N., Burek, P., Salamon, P., & de Roo, A. (2016). Integrating

- 1 remotely sensed surface water extent into continental scale hydrology. Journal of
2 Hydrology, 543, 659–670. Retrieved from <http://dx.doi.org/10.1016/j.jhydrol>
3 [.2016.10.041](http://dx.doi.org/10.1016/j.jhydrol.2016.10.041) doi: 10.1016/j.jhydrol.2016.10.041
- 4 riverlevels.uk. (2020). River levels, River Wye at Hereford Bridge. <https://riverlevels>
5 [.uk/herefordshire-hereford-old-wye-bridge-lvl#.YIFKt31KgUE](https://riverlevels.uk/herefordshire-hereford-old-wye-bridge-lvl#.YIFKt31KgUE), last access
6 29 April 2021.
- 7 Roberts, N. M., & Lean, H. W. (2008). Scale-selective verification of rainfall accumulations
8 from high-resolution forecasts of convective events. Monthly Weather Review. doi:
9 10.1175/2007MWR2123.1
- 10 Savage, J. T. S., Bates, P., Freer, J., Neal, J., & Aronica, G. (2016). When does spatial res-
11 olution become spurious in probabilistic flood inundation predictions? Hydrological
12 Processes, 30(13), 2014–2032. doi: 10.1002/hyp.10749
- 13 Schumann, G. (2019). The need for scientific rigour and accountability in flood mapping
14 to better support disaster response. Hydrological Processes, 33(24), 3138–3142. doi:
15 10.1002/hyp.13547
- 16 Schumann, G., Bates, P. D., Horritt, M. S., Matgen, P., & Pappenberger, F. (2009).
17 Progress in integration of remote sensing-derived flood extent and stage data and
18 hydraulic models. Reviews of Geophysics. doi: 10.1029/2008RG000274
- 19 Schumann, G., Giustarini, L., Tarpanelli, A., & Jarihani, B. (2022). Flood Modeling
20 and Prediction Using Earth Observation Data. Surveys in Geophysics(0123456789).
21 Retrieved from <https://doi.org/10.1007/s10712-022-09751-y> doi: 10.1007/
22 s10712-022-09751-y
- 23 Sefton, C., Muchan, K., Parry, S., Matthews, B., Barker, L. J., Turner, S., & Hannaford,
24 J. (2021). The 2019/2020 floods in the uk: a hydrological appraisal. Weather,
25 76(12), 378–384. Retrieved from [https://rmets.onlinelibrary.wiley.com/doi/](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/wea.3993)
26 [abs/10.1002/wea.3993](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/wea.3993) doi: <https://doi.org/10.1002/wea.3993>
- 27 Skok, G., & Roberts, N. (2016). Analysis of Fractions Skill Score properties for ran-

- 1 dom precipitation fields and ECMWF forecasts. Quarterly Journal of the Royal
2 Meteorological Society, 142(700), 2599–2610. doi: 10.1002/qj.2849
- 3 Skok, G., & Roberts, N. (2018). Estimating the displacement in precipitation forecasts us-
4 ing the Fractions Skill Score. Quarterly Journal of the Royal Meteorological Society,
5 144(711), 414–425. doi: 10.1002/qj.3212
- 6 Sky News. (2022). Pakistan floods: Satellite images and maps show scale
7 of disaster. [https://news.sky.com/story/pakistan-floods-satellite-images](https://news.sky.com/story/pakistan-floods-satellite-images-and-maps-show-scale-of-disaster-12685468)
8 [-and-maps-show-scale-of-disaster-12685468](https://news.sky.com/story/pakistan-floods-satellite-images-and-maps-show-scale-of-disaster-12685468), last access 23 August 2023.
- 9 Speight, L. J., Cranston, M. D., White, C. J., & Kelly, L. (2021). Operational and
10 emerging capabilities for surface water flood forecasting. Wiley Interdisciplinary
11 Reviews: Water, 8(3), 1–24. doi: 10.1002/wat2.1517
- 12 Start Network. (2022). Anticipation and risk financing. [https://startnetwork.org/](https://startnetwork.org/anticipation-and-risk-financing)
13 [anticipation-and-risk-financing](https://startnetwork.org/anticipation-and-risk-financing), last access 25th October 2022.
- 14 Stein, J., & Stoop, F. (2019). Neighborhood-based contingency tables including errors
15 compensation. Monthly Weather Review, 147(1), 329–344. doi: 10.1175/MWR-D
16 -17-0288.1
- 17 Stephens, E., & Cloke, H. (2014a). Improving flood forecasts for better flood preparedness
18 in the UK (and beyond). Geographical Journal, 180(4), 310–316. doi: 10.1111/
19 geoj.12103
- 20 Stephens, E., Schumann, G., & Bates, P. (2014). Problems with binary pattern measures
21 for flood model evaluation. Hydrological Processes. Retrieved from [https://doi](https://doi.org/10.1002/hyp.9979)
22 [.org/10.1002/hyp.9979](https://doi.org/10.1002/hyp.9979) doi: 10.1002/hyp.9979
- 23 Tavus, B., Kocaman, S., Nefeslioglu, H. A., & Gokceoglu, C. (2020). A fusion approach
24 for flood mapping using sentinel-1 and sentinel-2 datasets. International Archives
25 of the Photogrammetry, Remote Sensing and Spatial Information Sciences - ISPRS
26 Archives, 43(B3), 641–648. doi: 10.5194/isprs-archives-XLIII-B3-2020-641-2020
- 27 Tellman, B., Sullivan, J. A., Kuhn, C., Kettner, A. J., Doyle, C. S., Brakenridge, G. R., ...

- 1 Slayback, D. A. (2021). Satellite imaging reveals increased proportion of population
2 exposed to floods. *Nature*, *596*. Retrieved from [http://dx.doi.org/10.1038/
3 s41586-021-03695-w](http://dx.doi.org/10.1038/s41586-021-03695-w) doi: 10.1038/s41586-021-03695-w
- 4 The Guardian. (2022). Pakistan floods cause devastation – in pictures.
5 [https://www.theguardian.com/world/gallery/2022/aug/30/pakistan-floods
6 -cause-devastation-in-pictures](https://www.theguardian.com/world/gallery/2022/aug/30/pakistan-floods), last access 15th June 2023.
- 7 Trepekli, K., Friberg, T., Balstrøm, T., Fog, B., Allotey, A., Kofie, R., & Møller-Jensen,
8 L. (2021). UAV-LiDAR observations increase the precision of urban flood modelling
9 in Accra by detecting critical micro-topographic features. *EGU General Assembly*
10 online 19–30 Apr 2021. doi: 10.5194/egusphere-egu21-10457
- 11 Twele, A., Cao, W., Plank, S., & Martinis, S. (2016). Sentinel-1-based flood mapping: a
12 fully automated processing chain. *International Journal of Remote Sensing*, *37*(13),
13 2990–3004. Retrieved from <http://dx.doi.org/10.1080/01431161.2016.1192304>
14 doi: 10.1080/01431161.2016.1192304
- 15 UK Water Resources Portal. (2022). NRFA UKCEH UK Water Resources Portal.
16 <https://eip.ceh.ac.uk/hydrology/water-resources/>, last access 4 May 2022.
- 17 Ulbrich, U., Leckebusch, G. C., & Pinto, J. G. (2009). Extra-tropical cyclones in the
18 present and future climate: A review. *Theoretical and Applied Climatology*, *96*(1-
19 2), 117–131. doi: 10.1007/s00704-008-0083-8
- 20 UNDRR. (2022a). Early warnings for all the un global early warning initiative for the
21 implementation of climate adaptation. <https://library.wmo.int/idurl/4/58209>,
22 last access November 2023.
- 23 UNDRR. (2022b). Global assessment report on disaster risk reduction. [https://
24 www.undrr.org/gar2022-our-world-risk#container-downloads](https://www.undrr.org/gar2022-our-world-risk#container-downloads), last access 25th
25 October 2022.
- 26 UNDRR. (2023). Gar special report 2023 mapping resilience for the sustainable
27 development goals. <https://www.undrr.org/gar/gar2023-special-report>, last

- 1 access November 2023.
- 2 van Leeuwen, P. J. (2009). Particle Filtering in Geophysical Systems. Monthly Weather
3 Review, 137(12), 4089–4114. Retrieved from [https://journals.ametsoc.org/](https://journals.ametsoc.org/view/journals/mwre/137/12/2009mwr2835.1.xml)
4 [view/journals/mwre/137/12/2009mwr2835.1.xml](https://journals/mwre/137/12/2009mwr2835.1.xml) doi: [https://doi.org/10.1175/](https://doi.org/10.1175/2009MWR2835.1)
5 [2009MWR2835.1](https://doi.org/10.1175/2009MWR2835.1)
- 6 Wagner, W., Freeman, V., Cao, S., Matgen, P., Chini, M., Salamon, P., ... Briese, C.
7 (2020). DATA PROCESSING ARCHITECTURES FOR MONITORING FLOODS
8 USING SENTINEL-1. ISPRS Annals of the Photogrammetry, Remote Sensing
9 and Spatial Information Sciences, V-3-2020, 641–648. Retrieved from [https://](https://isprs-annals.copernicus.org/articles/V-3-2020/641/2020/)
10 isprs-annals.copernicus.org/articles/V-3-2020/641/2020/ doi: [10.5194/](https://doi.org/10.5194/isprs-annals-V-3-2020-641-2020)
11 [isprs-annals-V-3-2020-641-2020](https://doi.org/10.5194/isprs-annals-V-3-2020-641-2020)
- 12 WASDI. (2022). Wasdi. <https://www.wasdi.net/#!/marketplace>, last access 4th
13 November 2022.
- 14 WMO. (2017). South asian climate outlook forum held in bhutan. [https://](https://public.wmo.int/en/media/news/normal-rainfall-likely-much-of-south)
15 public.wmo.int/en/media/news/normal-rainfall-likely-much-of-south
16 [-asia-2017-southwest-monsoon](https://public.wmo.int/en/media/news/normal-rainfall-likely-much-of-south), last access 23rd January 2023.
- 17 World Resources Institute. (2023). How floods in pakistan threaten global
18 security. <https://www.wri.org/insights/pakistan-floods-threaten-global>
19 [-security#:~:text=In%20Pakistan%2C%20torrential%20rains%20and,%20killed%20over%20900%2C000%20livestock.,](https://www.wri.org/insights/pakistan-floods-threaten-global) last access 23 August 2023.
- 20
- 21 World Weather Attribution. (2022). Climate change likely increased extreme
22 monsoon rainfall, flooding highly vulnerable communities in pakistan.
23 <https://www.worldweatherattribution.org/climate-change-likely>
24 [-increased-extreme-monsoon-rainfall-flooding-highly-vulnerable](https://www.worldweatherattribution.org/climate-change-likely)
25 [-communities-in-pakistan/](https://www.worldweatherattribution.org/climate-change-likely), last access 23 August 2023.
- 26 Wu, W., Emerton, R., Duan, Q., Wood, A. W., Wetterhall, F., & Robertson, D. E. (2020).
27 Ensemble flood forecasting: Current status and future opportunities. WIREs Water,

- 1 7(3), 1–32. doi: 10.1002/wat2.1432
- 2 Zappa, M., Jaun, S., Germann, U., Walser, A., & Fundel, F. (2011, may). Superposition
3 of three sources of uncertainties in operational flood forecasting chains. Atmospheric
4 Research, 100(2-3), 246–262. doi: 10.1016/J.ATMOSRES.2010.12.005
- 5 Zsoter, E., Prudhomme, C., Stephens, E., Pappenberger, F., & Cloke, H. (2020). Using en-
6 semble reforecasts to generate flood thresholds for improved global flood forecasting.
7 Journal of Flood Risk Management, 13(4), 1–14. doi: 10.1111/jfr3.12658