



Flood hazard and change impact assessments may profit from rethinking model calibration strategies

Manuela I. Brunner¹, Lieke A. Melsen², Andrew W. Wood^{1,3}, Oldrich Rakovec^{4,5}, Naoki Mizukami¹, Wouter J. M. Knoben⁶, and Martyn P. Clark⁶

¹Research Applications Laboratory, National Center for Atmospheric Research, Boulder CO, USA

²Hydrology and Quantitative Water Management, Wageningen University, Wageningen, Netherlands

³Climate and Global Dynamics Laboratory, National Center for Atmospheric Research, Boulder CO, USA

⁴Department Computational Hydrosystems, Helmholtz Centre for Environmental Research, Leipzig, Germany

⁵Faculty of Environmental Sciences, Czech University of Life Sciences Prague, Praha – Suchbátka, Czech Republic

⁶University of Saskatchewan Coldwater Laboratory, Canmore, Canada

Correspondence: Manuela I. Brunner (manuelab@ucar.edu)

Abstract. Floods cause large damages, especially if they affect large regions. Assessments of current, local and regional flood hazards and their future changes often involve the use of hydrologic models. However, uncertainties in simulated floods can be considerable and yield unreliable hazard and climate change impact assessments. A reliable hydrologic model ideally reproduces both local flood characteristics and spatial aspects of flooding, which is, however, not guaranteed especially when using standard model calibration metrics. In this paper we investigate how flood timing, magnitude and spatial variability are represented by an ensemble of hydrological models when calibrated on streamflow using the Kling–Gupta efficiency metric, an increasingly common metric of hydrologic model performance. We compare how four well-known models (SAC, HBV, VIC, and mHM) represent (1) flood characteristics and their spatial patterns; and (2) how they translate changes in meteorologic variables that trigger floods into changes in flood magnitudes. Our results show that both the modeling of local and spatial flood characteristics is challenging. They further show that changes in precipitation and temperature are not necessarily well translated to changes in flood flow, which makes local and regional flood hazard assessments even more difficult for future conditions. We conclude that models calibrated on integrated metrics such as the Kling–Gupta efficiency alone have limited reliability in flood hazard assessments, in particular in regional and future assessments, and suggest the development of alternative process-based and spatial evaluation metrics.

1 Introduction

To derive flood estimates for current and future conditions, many studies use a hydrological model driven by present or future meteorological forcing data. However, data, model structure, and parameter uncertainties can be considerable (Clark et al., 2016) especially when considering extreme events such as floods (Brunner et al., 2019b; Das and Umamahesh, 2018) and when looking into the future. Reliable assessments of expected flood hazard and future changes in flooding are therefore often challenging.



A reliable model ideally reproduces different aspects of flooding, including local characteristics such as event magnitude and timing. It has been shown, however, that capturing magnitude and timing is challenging when standard calibration metrics are used individually for parameter estimation (Lane et al., 2019; Brunner and Sikorska, 2018; Mizukami et al., 2019). For example, one widely-used metric that is considered integrative compared to others (e.g., bias, correlation) is the Nash–Sutcliffe efficiency (E_{NS} ; Nash and Sutcliffe 1970), but it is formulated so that its optimal value actually underestimates flow variability (Gupta et al., 2009). Using a related metric, the Kling–Gupta efficiency (E_{KG} ; Gupta et al. 2009) can partially overcome this deficiency and improve simulations of peak flows (Mizukami et al., 2019). Yet to achieve further improvement, a broader range of application-specific evaluation metrics is typically required, including objectives that directly characterize hydrologic phenomena (or 'signatures') such as peak flows, flood volumes and timing, recession rates, and seasonal hydrograph shape. Considering multiple objectives in a step-wise calibration sequence, either manual or automated, is common in agencies that implement models for applications such as flood forecasting (Hogue et al., 2000), and strengthens their ability to provide reliable flood predictions. The use of multiple objectives may, however, lead to a decrease in performance with respect to any individual non-flood-related signature (Mizukami et al., 2019). Despite their deficiencies with respect to extremes, individual 'integrative' standard calibration metrics such as E_{NS} or E_{KG} are often used in research modeling studies, even when the focus is on floods and their future changes (for E_{NS} based calibration see e.g. Hundedcha and Merz 2012; Köplin et al. 2014; Vormoor et al. 2015; Wobus et al. 2017 and for E_{KG} based calibration see e.g. Harrigan et al. 2020; Hirpa et al. 2018; Huang et al. 2018; Thober et al. 2018; Brunner and Sikorska 2018)

In addition to timing and magnitude at individual catchments, it is also important to realistically reproduce spatial dependencies, i.e. the relationship of flood occurrence across gauging stations (Keef et al., 2013; De Luca et al., 2017; Berghuijs et al., 2019). An over- or underestimation of spatial dependencies in regional flood hazard and risk assessments has been shown to under- or overestimate regional damage, respectively (Lamb et al., 2010; Metin et al., 2020). Prudhomme et al. (2011) have shown for a set of large-scale hydrological models that simulated high flow episodes are less spatially coherent than observed events. Despite the high relevance for impact, the spatial aspects of flooding has often been overlooked in past simulation studies.

In this paper we explore the suitability of hydrological models for local and regional flood hazard assessments under current and future conditions. We evaluate the extent to which hydrological models calibrated against common individual calibration metrics reproduce (1) local flood characteristics (e.g. flood magnitude and timing at any given gauging station), (2) spatial dependencies in flooding, and (3) relationships between changes in flood triggering variables and changes in flood magnitude. We assess which aspects of hydrological models may need to be improved if we want to bring hazard and change impact assessments to a point where we can make more reliable assessments of regional flood hazard and future changes in local and spatial flood characteristics.

For this assessment, we look at the model output of four widely used hydrological models (Addor and Melsen, 2019), namely, the Sacramento Soil Moisture Accounting model (SAC-SMA) combined with SNOW-17 (Newman et al., 2015), the Hydrologiska Byråns Vattenbalansavdelning model (HBV; Bergström, 1976), the Variable Infiltration Capacity model (VIC; Liang et al., 1994), and the mesoscale hydrologic model (mHM; Kumar et al., 2013; Samaniego et al., 2010). We hope that

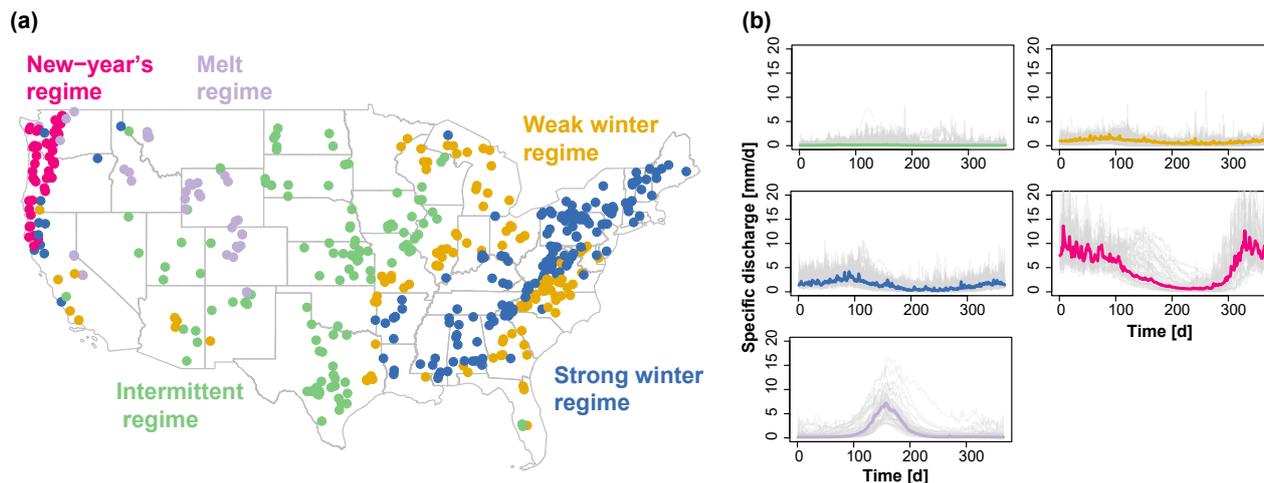


Figure 1. a) Map of the 488 catchments in the conterminous United States belonging to the five regime classes: 1) Intermittent, 2) weak winter, 3) strong winter, 4) New Year's, and 5) melt. b) Median regime per regime class (colored lines) and variability of regimes within a class (on line per catchment, grey) (Brunner et al., 2020b).

identifying and documenting model weaknesses regarding regional and future flooding will highlight avenues for future model development.

2 Data and Methods

To study how local and spatial flood characteristics are reproduced by hydrological models calibrated on streamflow using the individual calibration metric, E_{KG} , we compare observed to simulated flood event characteristics for a set of 488 catchments in the conterminous United States that have minimal human impact and catchment areas ranging from 4 to 2000 km² (Figure 1a) (Newman et al., 2015). The dataset comprises catchments with a wide range of climate and streamflow characteristics ranging from catchments with intermittent regimes and a very weak seasonality to catchments with a very strong seasonal cycle under the influence of snow (New Year's and melt regimes; Figure 1b; Brunner et al. 2020b). Observed streamflow time series are available from the U.S. Geological Survey (USGS, 2019).

2.1 Model simulations

We use daily streamflow simulations for the period 1981-2008 generated with four well-known hydrological models (Addor and Melsen, 2019) offering different model structures and complexity: the lumped SAC model, the lumped HBV model, the lumped version of the VIC model, and the grid-based, distributed mesoscale hydrologic model mHM. The model parameters



70 were calibrated on streamflow observations by minimizing the E_{KG} by Melsen et al. (2018) for SAC, HBV, and VIC and by Mizukami et al. (2019) for mHM. E_{KG} is defined as:

$$E_{KG}(Q) = 1 - \sqrt{[s_{\rho} \cdot (\rho - 1)]^2 + [s_{\alpha} \cdot (\alpha - 1)]^2 + [s_{\beta} \cdot (\beta - 1)]^2}, \quad (1)$$

where ρ is the correlation between observed and simulated runoff, α is the standard deviation of the simulated runoff divided by the standard deviation of observed runoff, and β is the mean of the simulated runoff, divided by the mean of the observed runoff. s_{ρ} , s_{α} , and s_{β} are scaling parameters enabling a weighting of different components. When used individually, E_{KG} has been found to result in a better performance for annual peak flow simulation than the long-standing and related hydrologic model evaluation metric Nash–Sutcliffe efficiency (E_{NS}) (Mizukami et al., 2019).

For SAC, Melsen et al. (2018) calibrated and evaluated 18 out of the 35 parameters available in the coupled Snow-17 and SAC-SMA modeling system, for HBV 15 parameters, for VIC 17 parameters, and for mHM Rakovec et al. (2019) and Mizukami et al. (2019) calibrated and evaluated up to 48 parameters. All the models were driven with spatially lumped meteorological forcing data: SAC, HBV, and VIC were driven with Daymet meteorological forcing (Thornton et al., 2012) and mHM with the forcing by Maurer et al. (2002). SAC, HBV, and VIC were evaluated on the period 1985–2008 while mHM was calibrated on the period 1999–2008 and evaluated on the period 1989–1999.

Model performance in terms of E_{KG} varies spatially and is related to the hydrological regime (Figure 2). It is lowest for catchments with intermittent regimes and a weak seasonality and highest for catchments with a strong seasonality such as a melt and New Year’s regime. The finding that intermittent regimes are challenging to model successfully is well known in hydrology and reproduced in many studies, e.g., Unduche et al. (2018), who show that hydrological modeling on Prairie watersheds is very complex (Hay et al., 2018). Intermittent regimes may suffer in calibration if they rely solely on correlation-type measures because their day to day variation is more difficult to reproduce than a more pronounced and regular seasonality. Overall model performance decreases from mHM (median E_{KG} 0.69), over SAC (median E_{KG} 0.63) and VIC (median E_{KG} 0.60) to HBV (median E_{KG} 0.52). In addition to streamflow, we use areal precipitation and simulated soil moisture to explain potential differences in model performance.

2.2 Model evaluation for floods

We compare local and spatial flood characteristics extracted from the observed time series to those of the series simulated with the four models for the period 1981–2008.

2.2.1 Flood event identification

Flood events are identified for each of the five time series (one observed, four simulated) using a peak-over-threshold (POT) approach similar to the one used in Brunner et al. (2019a, 2020b). This approach consists of two main steps and results in two data sets each, which are used for the local and spatial analysis, respectively: (1) POT events in individual catchments and (2) event occurrences across all catchments. In Step 1, independent POT events are identified in the daily discharge time series of

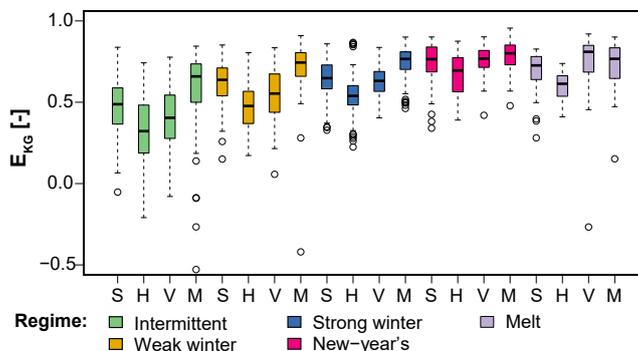


Figure 2. Model performance in terms of E_{KG} for the four models SAC (S), HBV (H), VIC (V), and mHM (M) per hydrological regime: intermittent (114 catchments), weak winter (108), strong winter (176), New Year's (50), and melt (40).

the individual catchments using the 25th percentile of the corresponding time series of annual maxima as a threshold (Schlef et al., 2019) and by prescribing a minimum time lag of 10 days between events (Diederer et al., 2019). In Step 2, a data set consisting of the dates of flood occurrences across all catchments is compiled. This set is converted into a binary matrix which specifies for each catchment (columns) whether or not it is affected by a specific event (rows). We consider a catchment to be affected by a certain event if it experiences an event within a window of ± 2 days of that event to take into account travel times. In addition to a binary matrix over all events, we set up seasonal binary matrices (winter: Dec–Feb, spring: Mar–May, summer: June–Aug, fall: Sept–Nov).

2.2.2 Flood characteristics at individual sites

We use the data sets resulting from Step 1, the POT events at individual catchments, to evaluate how well the models reproduce flood statistics at individual sites. We focus on the total number of events n (actual error: $n_s - n_o$, where s represents simulations and o observations), magnitude in terms of mean peak discharge x (relative error: $(x_s - x_o)/x_o$), and mean timing (absolute error: circular statistics).

2.2.3 Spatial flood dependence

We then use the data sets resulting from Step 2 to evaluate how models reproduce overall and seasonal spatial flood dependence. To do so, we use the connectedness measure introduced by Brunner et al. (2020a), which quantifies the number of catchments with which a specific catchment co-experiences floods. Following the definition used by Brunner et al. (2020a), a catchment is connected to another catchment if they share a certain number of events, i.e. at least 1% of the total or seasonal number of events.



2.2.4 Flood triggers

120 To explain potential differences in model performance, we look at the relationship of simulated peak discharge with the two flood triggers: precipitation and soil moisture on the day of flood occurrence. We focus on the day of occurrence because time of concentration is typically less than one day for small headwater basins.

2.2.5 Floods under change

In addition to assessing model performance under current climate conditions, we would like to understand potential, additional
125 challenges arising when interested in future conditions. To do so, we look at how models translate changes in event temperature and precipitation into changes in POT discharge. To perform this sensitivity analysis, we generate surrogate time series of temperature, precipitation, and streamflow for each catchment by resampling the available hydrological years with replacement (Wood et al., 2004; Brunner et al., 2020b). For each of the surrogate series, we again extract POT flood events using the same
130 procedure as described under Step 1. For each of the extracted events we then determine temperature and precipitation. We use the sets of peak discharge, event temperature and event precipitation to compute mean event discharge, temperature, and precipitation, which enables the derivation of a relationship between mean POT discharge and the two meteorological variables during events. We repeat the resampling $n = 500$ times to derive a relationship between changes in mean event temperature and precipitation and changes in mean POT streamflow. This resampling experiment results in a response surface of POT discharge
135 spanned by mean event temperature and mean event precipitation for each catchment. We summarize the results obtained at individual locations by computing horizontal and vertical sensitivity gradients on these reaction surfaces using a linear regression model. The horizontal gradient describes the strength of POT discharge changes in response to event temperature changes while the vertical gradient describes the strength of change in response to changes in event precipitation. Conducting this experiment for both observed and simulated time series allows for the determination of whether the models react to changes in mean event temperature and precipitation in the same way as the real world system and are therefore suitable for the
140 use in climate change impact assessments on floods. If models produce different climate sensitivities than the ones seen in the observations, the use of models to simulate sets of flood events for future conditions may preclude reliable change assessments.

3 Results

3.1 Flood characteristics at individual sites

Model performance at individual sites with respect to the number of events, event magnitude, and timing varies by model and
145 hydrological regime type (Figure 3). For most catchments, the number of flood events is relatively well simulated by most models (i.e. median deviation close to zero) except by HBV, which overestimates the number of events for catchments with intermittent, weak winter, and melt regimes (Brunner et al., 2020b). Event magnitude in terms of peak discharge is generally underestimated for all regime types independent of the model. Underestimation is in line with previous studies showing that using E_{KG} individually results in an underestimation of peak flow (Mizukami et al., 2019) due to an underestimation of variabil-

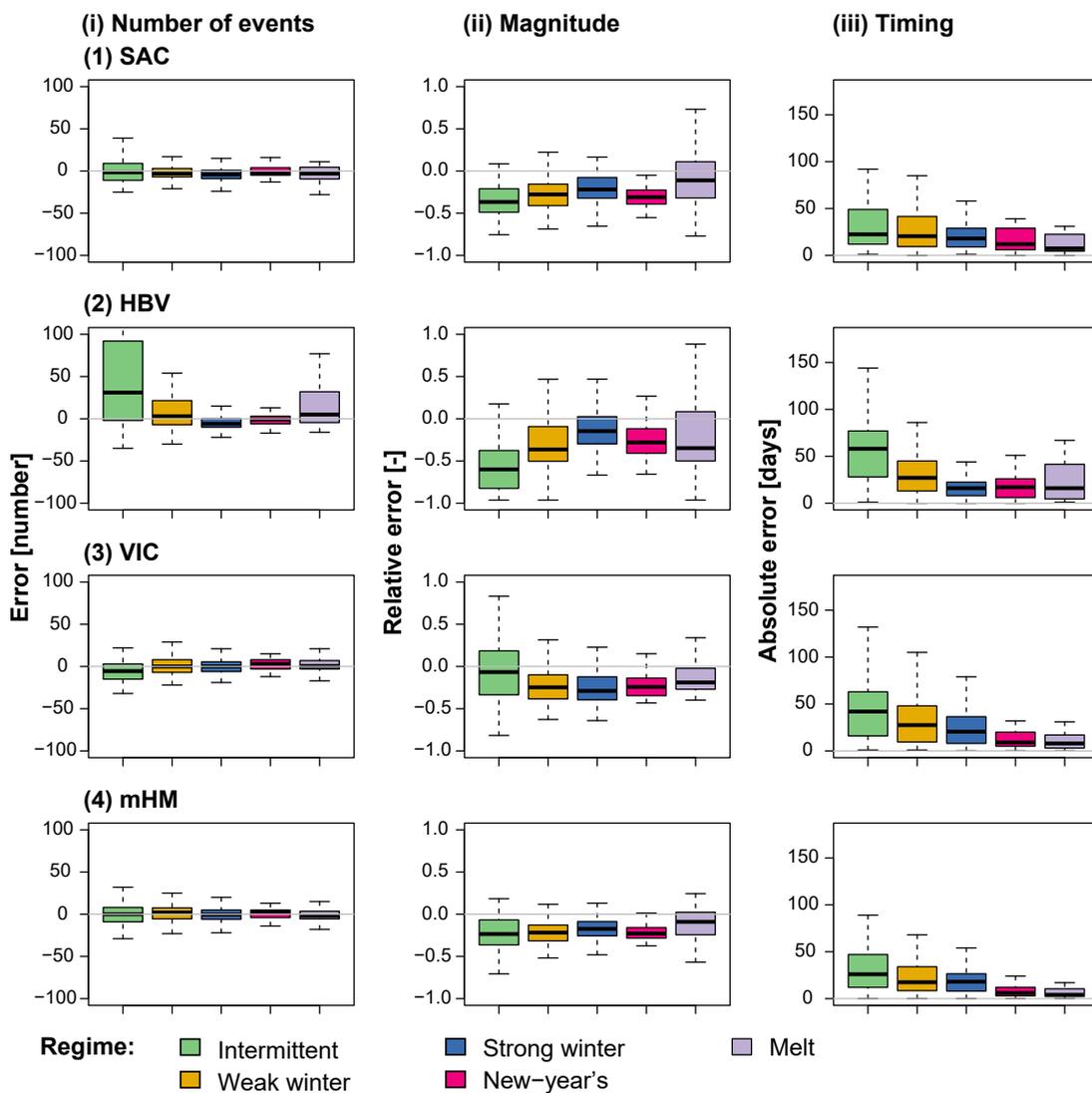


Figure 3. Model errors per regime type: intermittent (114 catchments), weak winter (108), strong winter (176), New Year's (50), and melt (40) (Figure 1). Errors are shown for (i) number of events (error in number of events), (ii) magnitude (relative error), and (iii) timing (absolute error in days) for the four models (1) SAC, (2) HBV, (3) VIC, and (4) mHM.



150 ity, which will result in an under-representation of extremes (Katz and Brown, 1992). Another factor potentially contributing to this underestimation is that the models were forced with spatially lumped instead of distributed data, which may smooth the simulated discharge response. On the other hand, the use of lumped forcings may also artificially synchronize hydrologic response, which would lead to overestimation.

Absolute flood timing errors are present in all models. They are the highest in catchments with intermittent regimes with a high variability in flood timing and low in catchments with a New Year's and melt regime where the flood season is limited to a few months (Brunner et al., 2020a). Over all, there is no clear tendency of one model to perform better than the other ones.

3.2 Spatial flood dependencies

Over all seasons, most models show an acceptable performance (i.e. median error close to zero) in the reproduction of spatial flood dependencies except for the HBV model, which overestimates spatial dependence (Figure 4) particularly in the Western part of the US. Seasonally, however, most models over- or underestimates spatial dependence in certain regions. In winter, connectedness is overestimated by most models except for VIC and the strength of overestimation is strongest for HBV. In spring, most models tend to underestimate spatial dependence except for HBV that results in an overestimation of spatial dependence for catchments with an intermittent regime. The overestimation of spatial dependence in winter is likely related to higher simulated than observed snowmelt as high soil moisture and snow availability have been shown to increase spatial flood connectedness (Brunner et al., 2020a). Related to this, the underestimation of spatial connectedness in spring may be related to the subsequent missing snowmelt contributions. Spatial connectedness in summer has been shown to be generally weak due to the occurrence of localized, convective events (Brunner et al., 2020a), which is reflected by most models except for HBV in the case of intermittent and melt regimes. Spatial flood connectedness has also been shown to be weak in fall (Brunner et al., 2020a) but is overestimated by most models. Connectedness overestimation is most pronounced for catchments with an intermittent regime independent of the model but especially expressed for HBV. The finding that there is room for improvement regarding the representation of spatial flood dependencies is in line with previous studies showing that large-scale hydrological models have a weakness in reproducing regional aspects of floods (Prudhomme et al., 2011).

3.3 Flood triggers

The differences in model performance regarding local and spatial flood characteristics may be partially explained by differences in their structure and how they transform precipitation into runoff. Figure 5 shows how simulated peak discharge is related to event precipitation, event precipitation plus snowmelt, and simulated soil moisture over all catchments for the four hydrologic models. The SAC and VIC models show similar simulated relationships for all three variable pairs. There is a clear positive relationship between peak discharge and precipitation and peak discharge and rainfall plus snowmelt, i.e. the higher the precipitation input or rainfall and snowmelt combined, respectively, the higher the resulting peak discharge. This relationship is slightly more expressed for VIC than for SAC. In both models, soil moisture and event magnitude are also positively related with lower values potentially associated with lower soil moisture states than more severe events. The peak discharge–precipitation relationship of HBV and mHM is less straightforward than the one of SAC and VIC. HBV and mHM

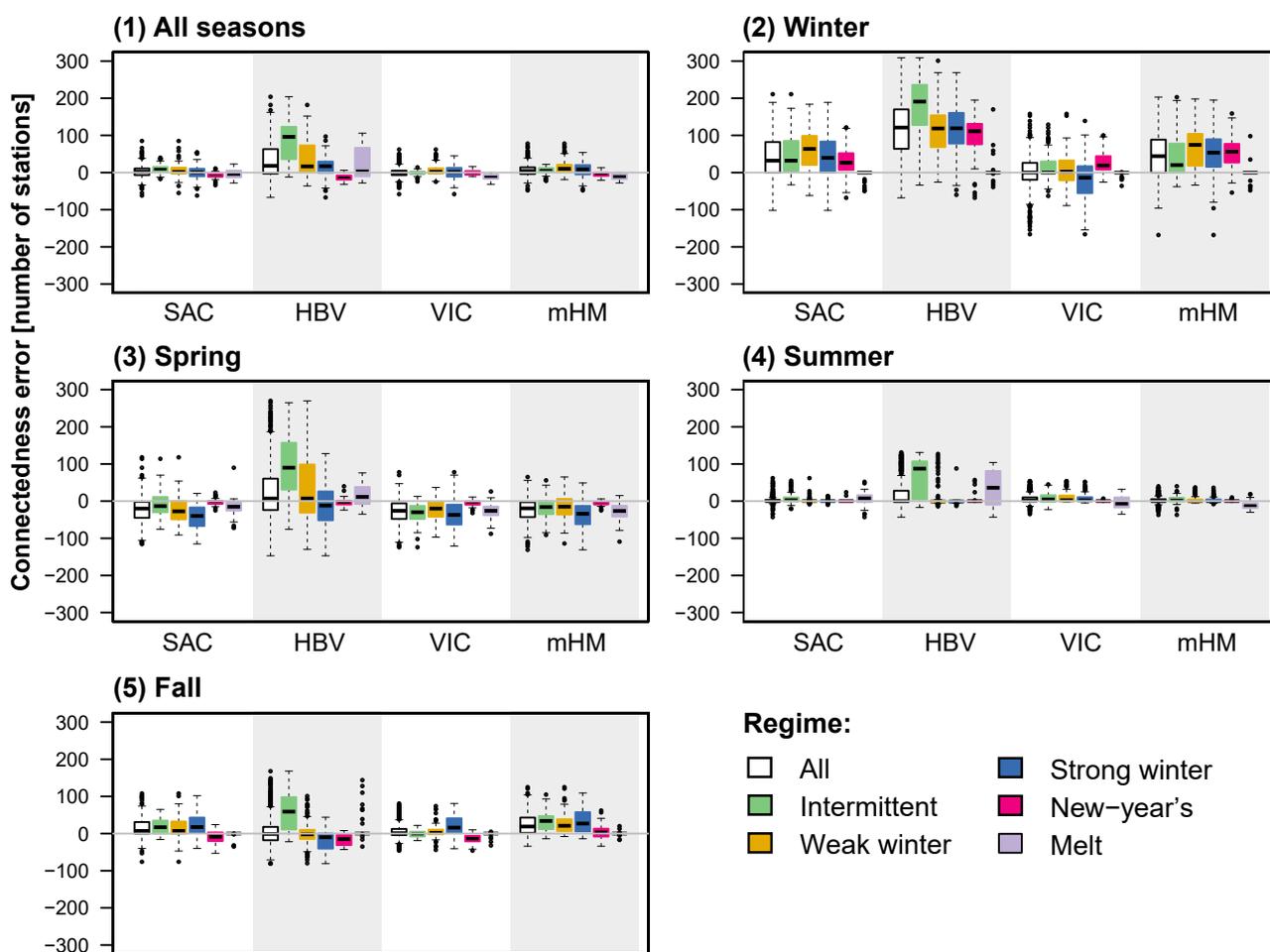


Figure 4. Overall (1) and seasonal (2–5) errors in flood connectedness (simulated minus observed connectedness), i.e. number of catchments a catchment is sharing at least 1% of the total number of flood events with, for the four models SAC, HBV, VIC, and mHM over all regimes and per regime: intermittent (114 catchments), weak winter (108), strong winter (176), New Year’s (50), and melt (40).

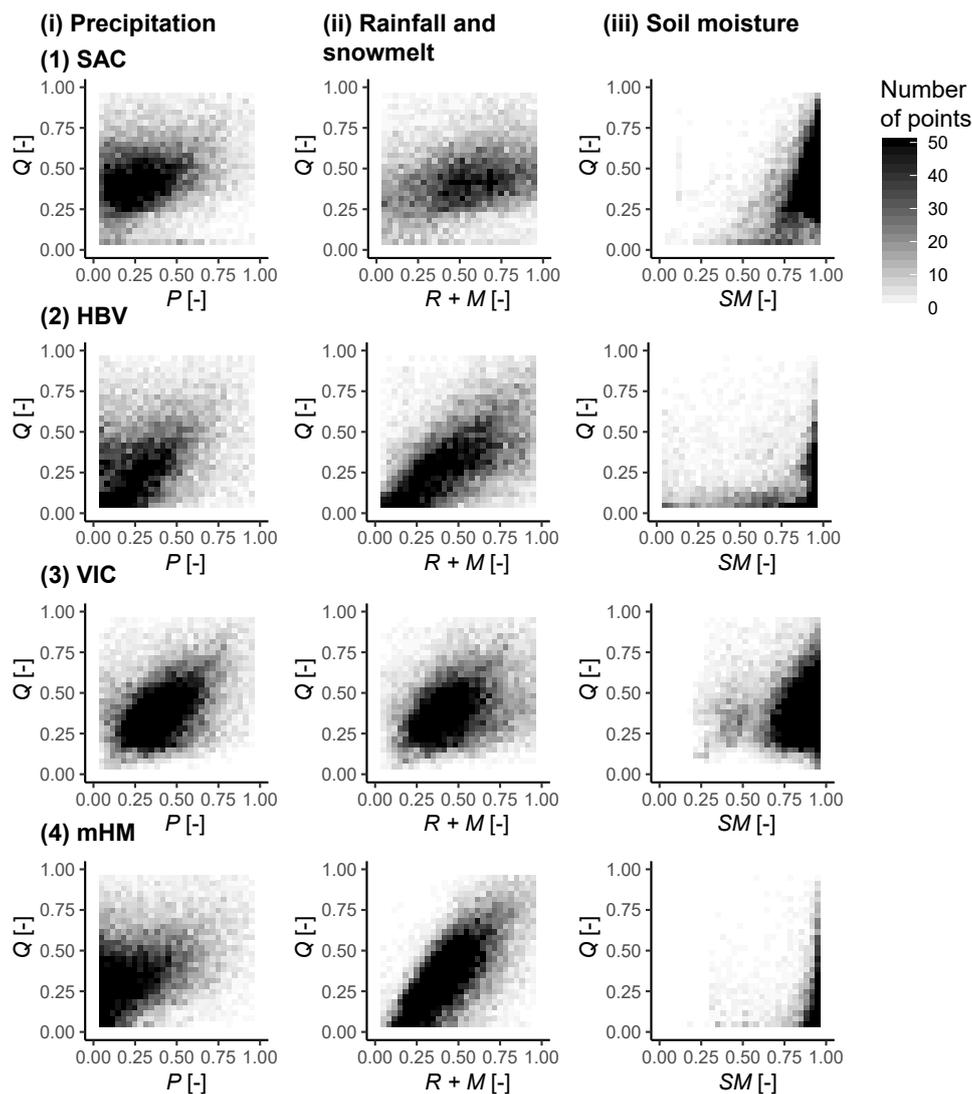


Figure 5. Simulated relationships between normalized flood discharge (Q) and normalized precipitation (i, P), rainfall and snowmelt (ii, $R + M$), and soil moisture (iii, SM , upper two soil layers for mHM) over all catchments represented by a binned scatter plot for the four hydrologic models (1) SAC, (2) HBV, (3) VIC, and (4) mHM. The darker the color, the higher the number of points within a bin (one point per catchment and event).



also show high discharge when precipitation input is high, but may in some cases still produce high discharge values even for low precipitation inputs. However, peak discharge and rainfall plus snowmelt show a strong linear relationship, i.e. the higher the combined rainfall and snowmelt input to the system, the higher is peak discharge. High flows are in most cases related to nearly full storage states but can occasionally also be triggered when soil moisture is low for SAC and VIC. These two types of responses may be related to differences in model behavior. VIC and SAC show more linearity in their event precipitation and peak discharge relationship than HBV and mHM, possibly because VIC and SAC have the capability to generate surface runoff when precipitation intensity exceeds infiltration capacity (Burnash et al., 1973; Liang et al., 1994). In this case, incoming precipitation is directly translated into flood discharge. In contrast, HBV and mHM, which is based on the HBV model structure (Kumar et al., 2013), do not include a surface runoff component and all discharge originates in the model stores (Bergström, 1976). This introduces a non-linearity in model response and may explain why a smaller precipitation input may still generate high peak flows in these models.

We here show that model performance to some degree depends on model choice and that many different combinations of forcing and model states can simulate floods. In addition, model performance may depend on input uncertainty, i.e. the precipitation product used to drive the models (Te Linde et al., 2007), which was here different for mHM than for the other models. These products may underestimate extreme rainfall or the spatial dependence of extreme precipitation at different locations because spatial smoothing or averaging during the gridding process reduces variability (Risser et al., 2019). The importance of input uncertainty is particularly pronounced if we are interested in future changes (Chen et al., 2014).

3.4 Floods under change

In addition to looking at how well local and spatial flood characteristics are represented by models, we look at how changes in temperature and event precipitation are translated into changes in flood flows to assess each models' suitability for climate impact assessments on floods. Our sensitivity analysis shows that the models have difficulty translating changes in event temperature and precipitation into sensitivities of flood flows (Figure 6), which can be problematic if we would like to use such models in climate change assessments. Generally, flood flows show a relatively low sensitivity to changes in mean event precipitation and temperature. This is in contrast to the behavior for mean flow, which is strongly influenced by changes in mean precipitation as demonstrated in a similar experiment by Brunner et al. (2020b). The much stronger relationship between mean precipitation and flow than between event precipitation and flow might arise because mean flow is a climate signal (Knoben et al., 2018), whereas floods are more an event (higher frequency, short-term) signal. However, some catchments, e.g. the Tucca Creek (New Year's regime) show a clear relationship between peak magnitude and both event temperature and precipitation. These relationships are, however, not necessarily captured by the models or the simulated sensitivities may even point in another direction than the observed ones (e.g. Pacific Creek, melt regime). In the case of melt regimes, the misrepresentation of flood sensitivities by models suggests that they may have difficulty simulating snow-influenced flooding.

This relatively poor model performance in capturing observed flood sensitivities can be generalized to the larger set of catchments studied here (Figure 7). Temperature sensitivities are found to be positive or negative, i.e. an increase in temperature could lead to an increase or decrease of peak flow depending on the catchment. In general, these temperature sensitivities are

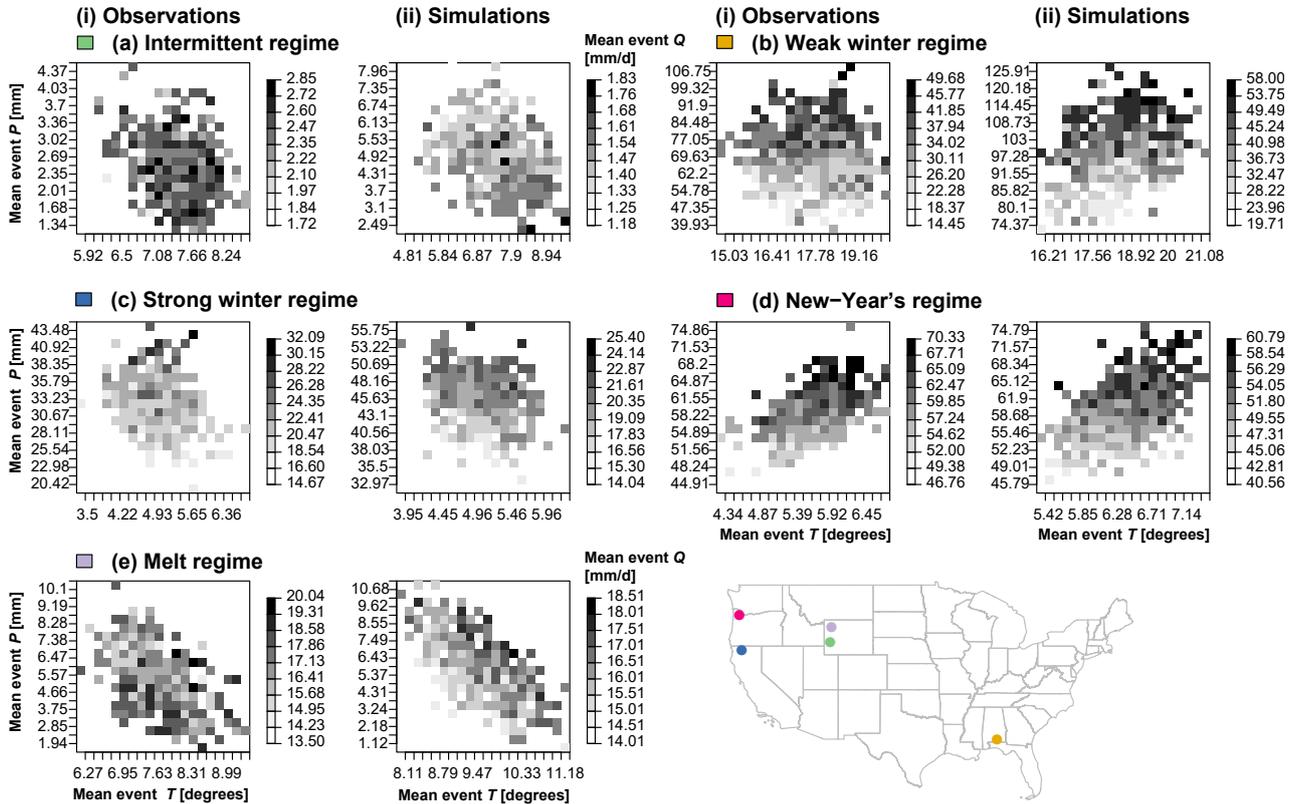


Figure 6. Climate sensitivity analysis for the VIC model: Dependence of mean POT magnitude (Q) on mean flood event precipitation (1-day; P) and mean flood temperature (T) for five example catchments, those with the best E_{KG} per regime type: intermittent regime (green; USGS ID 09210500 Fontanelle Creek near Fontanelle, WY; $E_{KG} = 0.78$), weak winter regime (yellow; USGS ID 02369800 Blackwater River near Bradley, AL; $E_{KG} = 0.83$), strong winter regime (blue; USGS ID 11522500 Salmon River above Somes, CA; $E_{KG} = 0.84$), New Year's regime (pink; USGS ID 14303200 Tucca Creek near Blaine, OR; $E_{KG} = 0.9$), and melt regime (purple; USGS ID 13011500 Pacific Creek at Moran, WY; $E_{KG} = 0.92$).

relatively weak (i.e. gradients are close to zero), which may be the reason why they are difficult to capture. In contrast, precipitation sensitivities are mostly positive, i.e. an increase in event precipitation leads to an increase in peak flow. However, the strength of these sensitivities is underestimated by all models, i.e. a change in precipitation leads to a too small change in peak flow. This underestimation of sensitivity can be understood by the underestimation of flood magnitude in general; namely that precipitation data likely underestimates actual event precipitation.

The results of this study indicate that the limited capability of hydrological models used in this study to reproduce observed hydrologic sensitivities during flooding may be related to insufficient model calibration (Fowler et al., 2016). We illustrate that reliance on an individual calibration metric (E_{KG}) rather than a broader suite of performance metrics can lead to simulation performance deficits for phenomena of interest, including an underestimation of streamflow variability (Mizukami et al., 2019),

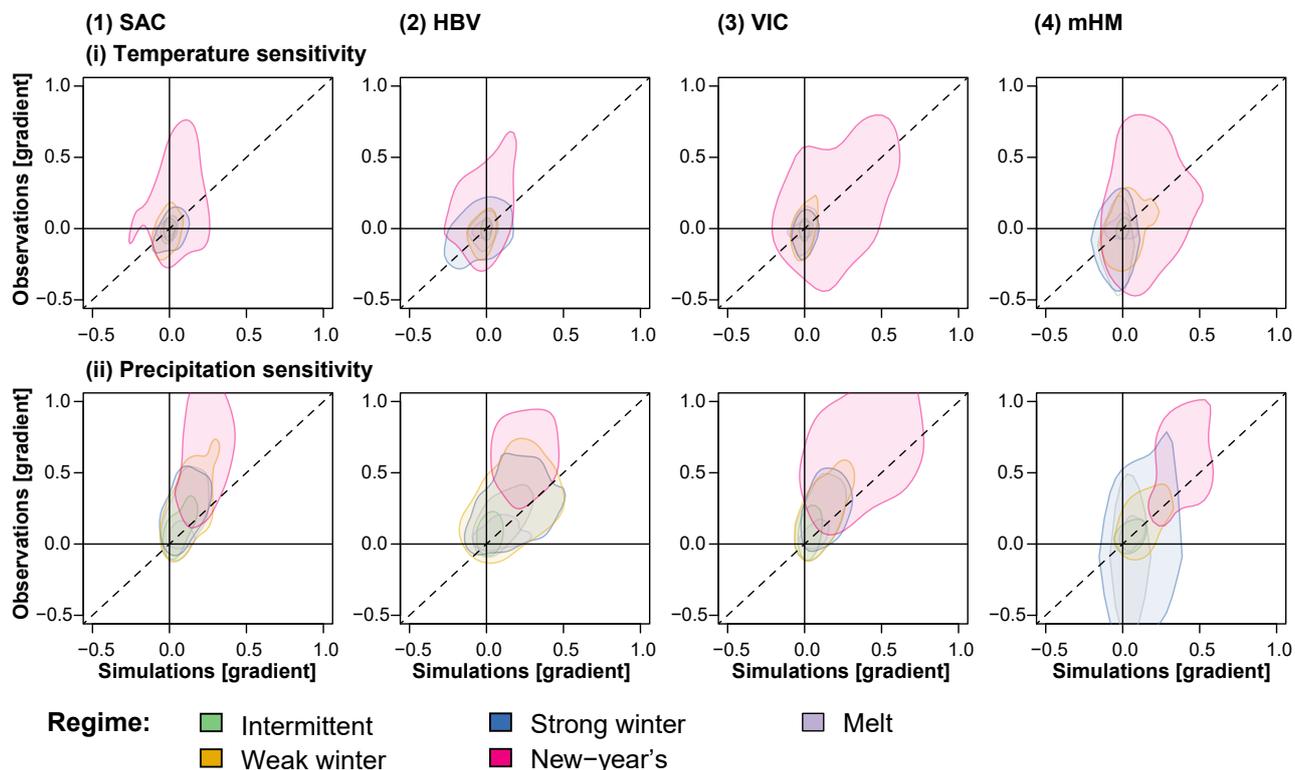


Figure 7. Observed vs. simulated (i) horizontal (temperature) and (ii) vertical (precipitation) climate sensitivities for floods represented by two-dimensional kernel density estimates for the four models (1) SAC, (2) HBV, (3) VIC, and (4) mHM for the five regime types: intermittent (114 catchments), weak winter (108), strong winter (176), New Year's (50), and melt (40) (Figure 1).

peak flood magnitudes and timing. As is evident in some existing practice-oriented applications of hydrological models (Hogue et al., 2000; Unduche et al., 2018; World Meteorological Organization, 2011), the simulation of floods and other hydrologic phenomena is likely to be improved by using more tailored model calibration strategies. The former would include either giving more weight to the variability component of an integrative metric such as the E_{KG} (Pool et al., 2017; Mizukami et al., 2019);
 230 whereas the latter might include optimizing explicitly for key flood characteristics (e.g., peak flow, volume, timing) and/or metrics depicting the fidelity of the model representation of soil moisture and snowmelt, within a multi-objective model calibration process (Moussa and Chahinian, 2009; Sikorska et al., 2018; Sikorska-Senoner et al., 2020). The spatial representation of extremes may also be improved by considering spatially distributed features of model response within a spatial calibration framework (Dembélé et al., 2020; Koch et al., 2018).



235 4 Conclusions

Our model comparison shows that all flood characteristics are not equally well represented by models calibrated with the widely used Kling–Gupta efficiency metric. Flood magnitude and timing are not always well captured by hydrological models in many catchments, in contrast to the number of flood events in a simulation time series, which tend to verify well. Flood magnitudes were underestimated by all models in most catchments, while the ability of the model to accurately reproduce event
240 timing was proportional to the hydroclimatic seasonality. These model deficiencies in reproducing local flood characteristics, especially timing, can lead to a misrepresentation of spatial flood dependencies, particularly in winter, because the temporal and spatial dimension of flooding are closely linked. The limited capability of the models in reproducing local and spatial flood characteristics is partly attributed to a reliance of the calibration on an individual variable (streamflow) and calibration metric (E_{KG}). While E_{KG} is integrative of certain properties (bias, variance, correlation), it does nonetheless not explicitly focus on
245 high flow values, their spatial dependencies, or processes generating high flow values. Such focus could be improved by giving more weight to the variability component of E_{KG} , if a single metric is used, or by using a suite of appropriate and targeted metrics in a multi-objective framework. The spatial concern could be addressed by applying spatial calibration procedures. Such steps are recommended if we would like to improve the reliability of local and regional flood hazard assessments.

Our sensitivity analysis also shows that climate sensitivities of floods, especially to changes in precipitation, are not well
250 represented in models even if the model can be deemed 'well-calibrated' via the individual E_{KG} metric. These sensitivities are generally underestimated by models independent of the geographical areas considered, i.e. an increase in event precipitation may not be translated into a strong enough increase in flood peak. The mis-estimation of these sensitivities may undermine the reliability of future flood hazard assessments relying on such models.

We conclude that calibration using only an individual model performance metric or variable can result in model implemen-
255 tations that have limited value for specific model applications, such as local and in particular spatial flood hazard analyses and change impact assessments. Despite its shortcomings, this practice has become increasingly more common and accepted in the research literature. Yet, our analysis illustrates that the development and adoption of more comprehensive multi-objective and multi-variable calibration strategies are needed to significantly improve model performance regarding floods under both current and future climate conditions.

260 *Data availability.* Observed streamflow measurements were made accessible by the USGS and can be downloaded via the website <https://waterdata.usgs.gov/nwis>. Simulated streamflow, precipitation, and storage time series can be requested from Lieke Melsen (lieke.melsen@wur.nl) for the SAC, HBV, and VIC models and for the mHM model from Oldrich Rakovec (oldrich.rakovec@ufz.de).

Author contributions. MIB and MPC developed the study design. NM, OR, and LAM provided the model simulations and together with MIB, MPC and WK interpreted the model output. AW assisted with the paper's background and messaging and proposed the climate
265 sensitivity strategy. MIB wrote the first draft of the manuscript and all co-authors revised and edited the manuscript.



Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. This work was supported by the Swiss National Science Foundation via a PostDoc.Mobility grant (P400P2_183844, granted to MIB). We acknowledge co-author support by the Bureau of Reclamation (CA R16AC00039) and the US Army Corps of Engineers (CSA 1254557). We also acknowledge support from the Global Water Futures research programme.



270 References

- Addor, N. and Melsen, L. A.: Legacy, rather than adequacy, drives the selection of hydrological models, *Water Resources Research*, 55, 378–390, <https://doi.org/10.1029/2018WR022958>, 2019.
- Berghuijs, W. R., Allen, S. T., Harrigan, S., and Kirchner, J. W.: Growing spatial scales of synchronous river flooding in Europe, *Geophysical Research Letters*, 46, 1423–1428, <https://doi.org/10.1029/2018GL081883>, 2019.
- 275 Bergström, S.: Development and application of a conceptual runoff model for Scandinavian catchments. Swedish Meteorological and Hydrological Institute (SMHI) RHO 7, Tech. Rep. January 1976, Sveriges Meteorologiska och Hydrologiska Institut, Norrköping, 1976.
- Brunner, M. I. and Sikorska, A. E.: Dependence of flood peaks and volumes in modeled runoff time series: effect of data disaggregation and distribution, *Journal of Hydrology*, 572, 620–629, <https://doi.org/10.1016/j.jhydrol.2019.03.024>, 2018.
- Brunner, M. I., Furrer, R., and Favre, A.-C.: Modeling the spatial dependence of floods using the Fisher copula, *Hydrology and Earth System*
280 *Sciences*, 23, 107–124, <https://doi.org/10.5194/hess-23-107-2019>, 2019a.
- Brunner, M. I., Hingray, B., Zappa, M., and Favre, A. C.: Future trends in the interdependence between flood peaks and volumes: Hydroclimatological drivers and uncertainty, *Water Resources Research*, 55, 1–15, <https://doi.org/10.1029/2019WR024701>, 2019b.
- Brunner, M. I., Gilleland, E., Wood, A., Swain, D. L., and Clark, M.: Spatial dependence of floods shaped by spatiotemporal variations in meteorological and land-surface processes, *Geophysical Research Letters*, p. under review, 2020a.
- 285 Brunner, M. I., Newman, A., Melsen, L. A., and Wood, A.: Future streamflow regime changes in the United States: assessment using functional classification, *Hydrology and Earth System Sciences Discussions*, p. under review, <https://doi.org/10.5194/hess-2020-54>, 2020b.
- Burnash, R. J., Ferral, R. L., and McGuire, R. A.: A generalized streamflow simulation system. Conceptual modeling for digital computers, Tech. rep., Joint Federal-State River Forecast Center, Sacramento, 1973.
- Chen, H., Sun, J., and Chen, X.: Projection and uncertainty analysis of global precipitation-related extremes using CMIP5 models, *International Journal of Climatology*, 34, 2730–2748, <https://doi.org/10.1002/joc.3871>, 2014.
- 290 Clark, M. P., Wilby, R. L., Gutmann, E. D., Vano, J. A., Gangopadhyay, S., Wood, A. W., Fowler, H. J., Prudhomme, C., Arnold, J. R., and Brekke, L. D.: Characterizing uncertainty of the hydrologic impacts of climate change, *Current Climate Change Reports*, 2, 55–64, <https://doi.org/10.1007/s40641-016-0034-x>, 2016.
- Das, J. and Umamahesh, N. V.: Assessment of uncertainty in estimating future flood return levels under climate change, *Natural Hazards*,
295 93, 109–124, <https://doi.org/10.1007/s11069-018-3291-2>, 2018.
- De Luca, P., Hillier, J. K., Wilby, R. L., Quinn, N. W., and Harrigan, S.: Extreme multi-basin flooding linked with extra-tropical cyclones, *Environmental Research Letters*, 12, 1–12, <https://doi.org/10.1088/1748-9326/aa868e>, 2017.
- Dembélé, M., Hrachowitz, M., Savenije, H. H. G., and Mariéthoz, G.: Improving the predictive skill of a distributed hydrological model by calibration on spatial patterns with multiple satellite datasets, *Water Resources Research*, p. e2019WR026085,
300 <https://doi.org/10.1029/2019WR026085>, 2020.
- Diederer, D., Liu, Y., Gouldby, B., Diermanse, F., and Vorogushyn, S.: Stochastic generation of spatially coherent river discharge peaks for continental event-based flood risk assessment, *Natural Hazards and Earth System Sciences*, 19, 1041–1053, <https://doi.org/10.5194/nhess-19-1041-2019>, 2019.
- Fowler, K. J. A., Peel, M. C., Western, A. W., Zhang, L., and Peterson, T. J.: Simulating runoff under changing climatic conditions: Revisiting
305 an apparent deficiency of conceptual rainfall-runoff models, *Water Resources Research*, 52, 1820–1846, <https://doi.org/10.1111/j.1752-1688.1969.tb04897.x>, 2016.



- Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling, *Journal of Hydrology*, 377, 80–91, <https://doi.org/10.1016/j.jhydrol.2009.08.003>, 2009.
- Harrigan, S., Zsoter, E., Alfieri, L., Prudhomme, C., Salamon, P., Barnard, C., Cloke, H., and Pappenberger, F.: GloFAS-ERA4 operational global river discharge reanalysis 1979-present, *Earth System Science Data*, pp. 1–23, 2020.
- 310 Hay, L., Norton, P., Viger, R., Markstrom, S., Steven Regan, R., and Vanderhoof, M.: Modelling surface-water depression storage in a Prairie Pothole Region, *Hydrological Processes*, 32, 462–479, <https://doi.org/10.1002/hyp.11416>, 2018.
- Hirpa, F. A., Salamon, P., Beck, H. E., Lorini, V., Alfieri, L., Zsoter, E., and Dadson, S. J.: Calibration of the Global Flood Awareness System (GloFAS) using daily streamflow data, *Journal of Hydrology*, 566, 595–606, <https://doi.org/10.1016/j.jhydrol.2018.09.052>, 2018.
- 315 Hogue, T. S., Sorooshian, S., Gupta, H., Holz, A., and Braatz, D.: A multistep automatic calibration scheme for river forecasting models, *Journal of Hydrometeorology*, 1, 524–542, 2000.
- Huang, S., Kumar, R., Rakovec, O., Aich, V., Wang, X., Samaniego, L., Liersch, S., and Krysanova, V.: Multimodel assessment of flood characteristics in four large river basins at global warming of 1.5, 2.0 and 3.0 K above the pre-industrial level, *Environmental Research Letters*, 13, 124005, <https://doi.org/10.1088/1748-9326/aae94b>, 2018.
- 320 Hundecha, Y. and Merz, B.: Exploring the relationship between changes in climate and floods using a model-based analysis, *Water Resources Research*, 48, <https://doi.org/10.1029/2011WR010527>, 2012.
- Katz, R. W. and Brown, B. G.: Extreme events in a changing climate: variability is more important than averages, *Climatic Change*, 21, 289–302, 1992.
- Keef, C., Tawn, J. A., and Lamb, R.: Estimating the probability of widespread flood events, *Environmetrics*, 24, 13–21, <https://doi.org/10.1002/env.2190>, 2013.
- 325 Knobén, W. J., Woods, R. A., and Freer, J. E.: A quantitative hydrological climate classification evaluated with independent streamflow data, *Water Resources Research*, 54, 5088–5109, <https://doi.org/10.1029/2018WR022913>, 2018.
- Koch, J., Demirel, M. C., and Stisen, S.: The SPAtial EFficiency metric (SPAEF): Multiple-component evaluation of spatial patterns for optimization of hydrological models, *Geoscientific Model Development*, 11, 1873–1886, <https://doi.org/10.5194/gmd-11-1873-2018>, 2018.
- 330 Köplin, N., Schädler, B., Viviroli, D., and Weingartner, R.: Seasonality and magnitude of floods in Switzerland under future climate change, *Hydrological Processes*, 28, 2567–2578, <https://doi.org/10.1002/hyp.9757>, 2014.
- Kumar, R., Samaniego, L., and Attinger, S.: Implications of distributed hydrologic model parameterization on water fluxes at multiple scales and locations, *Water Resources Research*, 49, 360–379, <https://doi.org/10.1029/2012WR012195>, 2013.
- Lamb, R., Keef, C., Tawn, J., Laeger, S., Meadowcroft, I., Surendran, S., Dunning, P., and Batstone, C.: A new method to assess the risk of local and widespread flooding on rivers and coasts, *Journal of Flood Risk Management*, 3, 323–336, <https://doi.org/10.1111/j.1753-318X.2010.01081.x>, 2010.
- 335 Lane, R. A., Coxon, G., Freer, J. E., Wagener, T., Johnes, P. J., Bloomfield, J. P., Greene, S., Macleod, C. J. A., and Reaney, S. M.: Benchmarking the predictive capability of hydrological models for river flow and flood peak predictions across over 1000 catchments in Great Britain, *Hydrology and Earth System Sciences*, 23, 4011–4032, <https://doi.org/10.5194/hess-23-4011-2019>, 2019.
- 340 Liang, X., Lettenmaier, D. P., Wood, E. F., and Burges, S. J.: A simple hydrologically based model of land surface water and energy fluxes for general circulation models, *Journal of Geophysical Research*, 99, 14415, <https://doi.org/10.1029/94JD00483>, 1994.
- Maurer, E. P., Wood, A. W., Adam, J. C., Lettenmaier, D. P., and Nijssen, B.: A long-term hydrologically based dataset of land surface fluxes and states for the conterminous United States, *Journal of Climate*, 15, 3237–3251, 2002.



- Melsen, L., Addor, N., Mizukami, N., Newman, A., Torfs, P., Clark, M., Uijlenhoet, R., and Teuling, R.: Mapping (dis) agreement in hydrologic projections, *Hydrology and Earth System Sciences*, 22, 1775–1791, <https://doi.org/10.5194/hess-22-1775-2018>, 2018.
- Metin, A. D., Dung, N. V., Schröter, K., Vorogushyn, S., Guse, B., Kreibich, H., and Merz, B.: The role of spatial dependence for large-scale flood risk estimation, *Natural Hazards and Earth System Sciences*, 20, 967–979, <https://doi.org/10.5194/nhess-2019-393>, 2020.
- Mizukami, N., Rakovec, O., Newman, A. J., Clark, M. P., Wood, A. W., Gupta, H. V., and Kumar, R.: On the choice of calibration metrics for "high-flow" estimation using hydrologic models, *Hydrology and Earth System Sciences*, 23, 2601–2614, <https://doi.org/10.5194/hess-23-2601-2019>, 2019.
- Moussa, R. and Chahinian, N.: Comparison of different multi-objective calibration criteria using a conceptual rainfall-runoff model of flood events, *Hydrology and Earth System Sciences*, 13, 519–535, <https://doi.org/10.5194/hess-13-519-2009>, 2009.
- Nash, J. E. and Sutcliffe, I. V.: River flow forecasting through conceptual models Part I - A discussion of principles, *Journal of Hydrology*, 10, 282–290, 1970.
- Newman, A. J., Clark, M. P., Sampson, K., Wood, A., Hay, L. E., Bock, A., Viger, R. J., Blodgett, D., Brekke, L., Arnold, J. R., Hopson, T., and Duan, Q.: Development of a large-sample watershed-scale hydrometeorological data set for the contiguous USA: Data set characteristics and assessment of regional variability in hydrologic model performance, *Hydrology and Earth System Sciences*, 19, 209–223, <https://doi.org/10.5194/hess-19-209-2015>, 2015.
- Pool, S., Vis, M. J., Knight, R. R., and Seibert, J.: Streamflow characteristics from modeled runoff time series - Importance of calibration criteria selection, *Hydrology and Earth System Sciences*, 21, 5443–5457, <https://doi.org/10.5194/hess-21-5443-2017>, 2017.
- Prudhomme, C., Parry, S., Hannaford, J., Clark, D. B., Hagemann, S., and Voss, F.: How well do large-scale models reproduce regional hydrological extremes: In Europe?, *Journal of Hydrometeorology*, 12, 1181–1204, <https://doi.org/10.1175/2011JHM1387.1>, 2011.
- Rakovec, O., Mizukami, N., Kumar, R., Newman, A. J., Thober, S., Wood, A. W., Clark, M. P., and Samaniego, L.: Diagnostic evaluation of large-domain hydrologic models calibrated across the Contiguous United States, *Journal of Geophysical Research: Atmospheres*, 124, 13 991–14 007, <https://doi.org/10.1029/2019JD030767>, 2019.
- Risser, M. D., Paciorek, C. J., Wehner, M. F., O'Brien, T. A., and Collins, W. D.: A probabilistic gridded product for daily precipitation extremes over the United States, *Climate Dynamics*, 53, 2517–2538, <https://doi.org/10.1007/s00382-019-04636-0>, 2019.
- Samaniego, L., Kumar, R., and Attinger, S.: Multiscale parameter regionalization of a grid-based hydrologic model at the mesoscale, *Water Resources Research*, 46, 1–25, <https://doi.org/10.1029/2008WR007327>, 2010.
- Schlef, K. E., Moradkhani, H., and Lall, U.: Atmospheric circulation patterns associated with extreme United States floods identified via machine learning, *Scientific Reports*, 9, 1–12, <https://doi.org/10.1038/s41598-019-43496-w>, 2019.
- Sikorska, A. E., Viviroli, D., and Seibert, J.: Effective precipitation duration for runoff peaks based on catchment modelling, *Journal of Hydrology*, 556, 510–522, <https://doi.org/10.1016/j.jhydrol.2017.11.028>, 2018.
- Sikorska-Senoner, A. E., Schaeffli, B., and Seibert, J.: Downsizing parameter ensembles for simulations of extreme floods, *Natural Hazards and Earth System Sciences Discussions*, p. under review, <https://doi.org/10.5194/nhess-2020-79>, 2020.
- Te Linde, A. H., Aerts, J., Dolman, H., and Hurkmans, R.: Comparing model performance of the HBV and VIC models in the Rhine basin, in: *Quantification and Reduction of Predictive Uncertainty for Sustainable Water Resources Management*, 313, pp. 278–285, 2007.
- Thober, S., Kumar, R., Wanders, N., Marx, A., Pan, M., Rakovec, O., Samaniego, L., Sheffield, J., Wood, E. F., and Zink, M.: Multi-model ensemble projections of European river floods and high flows at 1.5, 2, and 3 degrees global warming, *Environmental Research Letters*, 13, <https://doi.org/10.1088/1748-9326/aa9e35>, 2018.



- Thornton, P., Thornton, M., Mayer, B., Wilhelmi, N., Wei, Y., and Cook, R.: Daymet: daily surface weather on a 1 km grid for North America, 1980-2012, 2012.
- Unduche, F., Tolossa, H., Senbeta, D., and Zhu, E.: Evaluation of four hydrological models for operational flood forecasting in a Canadian Prairie watershed, *Hydrological Sciences Journal*, 63, 1133–1149, <https://doi.org/10.1080/02626667.2018.1474219>, 2018.
- 385 USGS: USGS Water Data for the Nation, <https://waterdata.usgs.gov/nwis>, 2019.
- Vormoor, K., Lawrence, D., Heistermann, M., and Bronstert, A.: Climate change impacts on the seasonality and generation processes of floods – projections and uncertainties for catchments with mixed snowmelt/rainfall regimes, *Hydrol. Earth Syst. Sci.*, 19, 913–931, <https://doi.org/10.5194/hess-19-913-2015>, 2015.
- 390 Wobus, C., Gutmann, E., Jones, R., Rissing, M., Mizukami, N., Lorie, M., Mahoney, H., Wood, A. W., Mills, D., and Martinich, J.: Climate change impacts on flood risk and asset damages within mapped 100-year floodplains of the contiguous United States, *Natural Hazards and Earth System Sciences*, 17, 2199–2211, <https://doi.org/10.5194/nhess-17-2199-2017>, 2017.
- Wood, A. W., Leung, L. R., Sridhar, V., and Lettenmaier, D. P.: Hydrologic implications of dynamical and statistical approaches to down-scaling climate model outputs, *Climatic Change*, 62, 189–216, <https://doi.org/10.1023/B:CLIM.0000013685.99609.9e>, 2004.
- World Meteorological Organization: Manual on flood forecasting and warning, Tech. Rep. 1072, WMO, Geneva, 2011.