



# 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

A. J. Newman<sup>1</sup>, M. P. Clark<sup>1</sup>, K. Sampson<sup>1</sup>, A. Wood<sup>1</sup>, L. E. Hay<sup>2</sup>, A. Bock<sup>2</sup>, R. J. Viger<sup>2</sup>, D. Blodgett<sup>3</sup>, L. Brekke<sup>4</sup>, J. R. Arnold<sup>5</sup>, T. Hopson<sup>1</sup>, and Q. Duan<sup>6</sup>

<sup>1</sup>National Center for Atmospheric Research, Boulder CO, USA

<sup>2</sup>United States Geological Survey, Modeling of Watershed Systems, Lakewood CO, USA

<sup>3</sup>United States Geological Survey, Center for Integrated Data Analytics, Middleton WI, USA

<sup>4</sup>US Department of Interior, Bureau of Reclamation, Denver CO, USA

<sup>5</sup>US Army Corps of Engineers, Institute for Water Resources, Seattle WA, USA

<sup>6</sup>Beijing Normal University, Beijing, China

Correspondence to: A. J. Newman (anewman@ucar.edu)

Received: 17 April 2014 – Published in Hydrol. Earth Syst. Sci. Discuss.: 28 May 2014

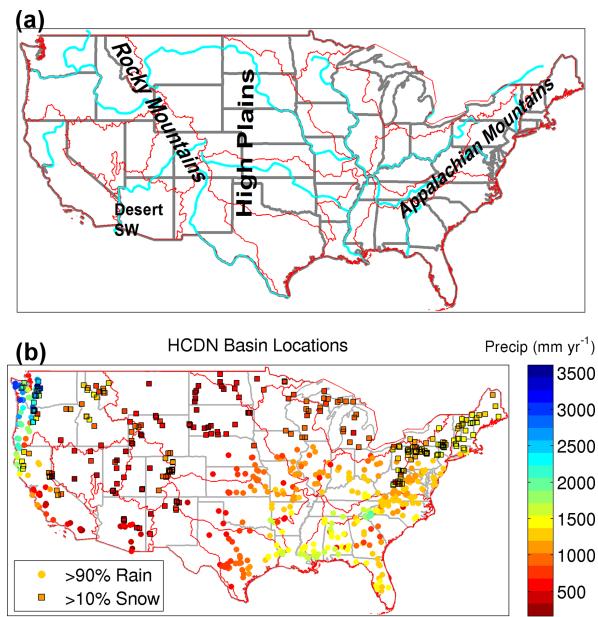
Revised: 23 November 2014 – Accepted: 1 December 2014 – Published: 14 January 2015

**Abstract.** We present a community data set of daily forcing and hydrologic response data for 671 small- to medium-sized basins across the contiguous United States (median basin size of 336 km<sup>2</sup>) that spans a very wide range of hydroclimatic conditions. Area-averaged forcing data for the period 1980–2010 was generated for three basin spatial configurations – basin mean, hydrologic response units (HRUs) and elevation bands – by mapping daily, gridded meteorological data sets to the subbasin (Daymet) and basin polygons (Daymet, Maurer and NLDAS). Daily streamflow data was compiled from the United States Geological Survey National Water Information System. The focus of this paper is to (1) present the data set for community use and (2) provide a model performance benchmark using the coupled Snow-17 snow model and the Sacramento Soil Moisture Accounting Model, calibrated using the shuffled complex evolution global optimization routine. After optimization minimizing daily root mean squared error, 90 % of the basins have Nash–Sutcliffe efficiency scores  $\geq 0.55$  for the calibration period and 34 %  $\geq 0.8$ . This benchmark provides a reference level of hydrologic model performance for a commonly used model and calibration system, and highlights some regional variations in model performance. For example, basins with a more pronounced seasonal cycle generally have a negative low flow bias, while basins with a smaller seasonal cycle have a

positive low flow bias. Finally, we find that data points with extreme error (defined as individual days with a high fraction of total error) are more common in arid basins with limited snow and, for a given aridity, fewer extreme error days are present as the basin snow water equivalent increases.

## 1 Introduction

With the increasing availability of gridded meteorological data sets, streamflow records and computing resources, large-sample hydrology studies have become more common in the last decade or more (i.e., Nathan and McMahon, 1990; Perrin et al., 2001; Maurer et al., 2002; Beldring et al., 2003; Merz and Bloschl, 2004; Andreassian et al., 2004; Lohmann et al., 2004; Duan et al., 2006; Oudin et al., 2006, 2010; Samaniego et al., 2010; Martinez and Gupta, 2010, 2011; Nester et al., 2011, 2012; Livneh and Lettenmaier, 2012, 2013; Kumar et al., 2013; Oubeidillah et al., 2013). Within the United States there have been several studies to produce large-sample hydrometeorological data sets (Maurer et al., 2002; Lohmann et al., 2004; Duan et al., 2006; Thornton et al., 2012; Xia et al., 2012; Livneh et al., 2013). Many of these data sets provide gridded data and may need to be further processed by the end user for their specific hydrologic model configuration.



**Figure 1.** (a) Contiguous United States (CONUS) with states (gray), rivers (blue) and major hydrologic regions (red). Text indicates major geographic regions discussed in text. (b) Location of the 671 HCDN-2009 basins across the contiguous US used in the basin data set with precipitation shaded. Circles denote basins with > 90 % of their precipitation falling as rain, squares with black outlines denote basins with > 10 % of their precipitation falling as snow as determined by using a 0 °C daily mean Daymet temperature threshold. State outlines are in thin gray and hydrologic regions in thin red.

The Model Parameter Estimation Project (MOPEX) data set does provide basin mean hydrometeorological data and observed streamflow records for 438 basins across the contiguous United States (CONUS; Schaake et al., 2006) for over more than 30 years; making it one of the few, high-quality, freely available hydrometeorological data sets with immediate applicability to catchment-type hydrologic models.

Gupta et al. (2014) emphasize that more large-sample hydrologic studies are needed to “balance depth with breadth”; most hydrologic studies have traditionally focused on one or a small number of basins (depth), which hinders the ability to establish general hydrologic concepts applicable across regions (breadth). Gupta et al. (2014) go on to discuss practical considerations for large-sample hydrology studies, noting first and foremost that large data sets of quality basin data need to be available and shared in the community. In support of this philosophy, we present a large-sample hydrometeorological data set and modeling tools to understand regional variability in hydrologic model performance across the contiguous US (Fig. 1). The development of the basin data set presented herein takes advantage of high-quality, freely available data from various US government agencies and re-

search laboratories. It includes (1) daily forcing data for 671 basins for multiple spatial configurations over the 1980–2010 time period; (2) daily streamflow data; (3) basic metadata (e.g., location, elevation, size, and basin delineation shapefiles) and (4) benchmark model performance which contains the final calibrated model parameter sets, model output time series for all basins as well as summary graphics for each basin. This builds on the MOPEX data set by providing basin mean forcing data for 233 more basins along with two other spatial configurations and the benchmark model performance parameter sets and model output.

This data set and benchmark application is intended for the community to use as a test bed to facilitate the evaluation of hydrologic modeling and prediction questions. To this end, the benchmark consists of the calibrated, coupled Snow-17 snow model and the Sacramento Soil Moisture Accounting Model (SAC-SMA) for all 671 basins using the shuffled complex evolution (SCE) global optimization routine. Development of a large-sample hydrologic data set such as this will allow for exploration into many important scientific questions. We provide some basic analysis relating to questions such as (1) what is the model performance across a large sample of basins and how does model performance vary across basin hydroclimatic conditions? (2) How do error characteristics relate to basin calibration performance and hydroclimatic conditions? This basic analysis is intended to highlight some of the important questions that can be answered through large-sample hydrologic studies and provide example results for further exploration.

The next section describes the development of the basin data set from basin selection through forcing data generation. It then briefly describes the modeling system and calibration routine. Next, example results using the basin data set and modeling platform are presented. Finally, concluding thoughts and next steps are discussed.

## 2 Basin data set

The development of a freely available large-sample basin data set requires several choices and subsequent data acquisition. Three major decisions were made and are discussed in this section: (1) the selection process for the basins, (2) the various basin spatial configurations to be developed, and (3) selection of the underlying forcing data set used to develop forcing data time series. Additionally, aggregation of the necessary streamflow data is described.

### 2.1 Basin selection

The United States Geological Survey (USGS) developed an updated version of their Geospatial Attributes of Gages for Evaluating Streamflow (GAGES-II) in 2011 (Falcone et al., 2010; Falcone, 2011). This database contains geospatial information for over 9000 stream gages maintained by the

USGS. As a subset of the GAGES-II database, a portion of the basins with minimal human disturbance (i.e., minimal land use changes or disturbances, minimal human water withdrawals) are noted as “reference” gages. A further subsetting of the reference gages were made as a follow-on to the Hydro-Climatic Data Network (HCDN) 1988 data set (Slack and Landwehr, 1992). These gages, marked HCDN-2009 (Lins, 2012), meet the following criteria: (1) have at least 20 years of complete flow data between 1990 and 2009 and were active as of 2009, (2) are a GAGES-II reference gage, (c) have less than 5 % imperviousness as measured by the National Land Cover Database (NLCD-2011; Jin et al., 2013), and (d) passed a manual survey of human impacts in the basin by local Water Science Center evaluators (Falcone et al., 2010). There are 704 gages in the GAGES-II database that are considered HCDN-2009 across the CONUS. This study uses that portion of the HCDN-2009 basin set as the starting point since they should best represent natural flow conditions. After initial processing and data availability requirements, 671 basins are used for analysis in this study (Fig. 1b). Because these basins have minimal human influence they are almost exclusively smaller, headwater-type basins.

## 2.2 Forcing and streamflow data

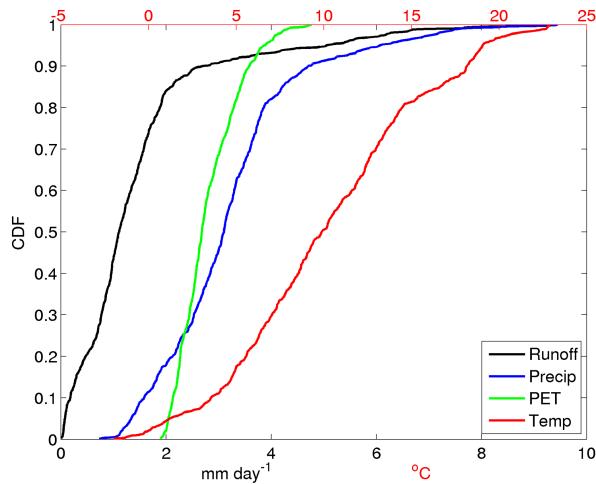
Hydrologic models are run with a variety of spatial configurations, including entire watersheds (lumped), elevation bands, hydrologic response units (HRUs), or grids. For this data set, forcing data were calculated (via areal averaging) for watershed, HRU and elevation band spatial configurations. The basin spatial configurations were created from the base national geospatial fabric for hydrologic modeling developed by the USGS Modeling of Watershed Systems (MoWS) group (Viger, 2014; Viger and Bock, 2014). The geospatial fabric is a watershed-oriented analysis of the National Hydrography Data set that contains points of interest (e.g., USGS streamflow gauges), hydrologic response unit boundaries and simplified stream segments (not used in this study). This geospatial fabric contains points of interest that include USGS streamflow gauges and allowed for the determination of upstream total basin area and basin HRUs (Viger, 2014; Viger and Bock, 2014). A digital elevation model (DEM) was applied to the geospatial fabric data set to create elevation contour polygon shapefiles for each basin. The USGS Geo Data Portal (GDP) developed by the USGS Center for Integrated Data Analytics (CIDA) (Blodgett et al., 2011) was leveraged to produce area-weighted forcing data for the various basin spatial configurations over our time period. The GDP performs all necessary spatial subsetting and weighting calculations and returns the area-weighted time series for the specified inputs.

The Daymet data set was selected as the primary gridded meteorological data set to derive forcing data for our streamflow simulations (Thornton et al., 2012). Daymet was cho-

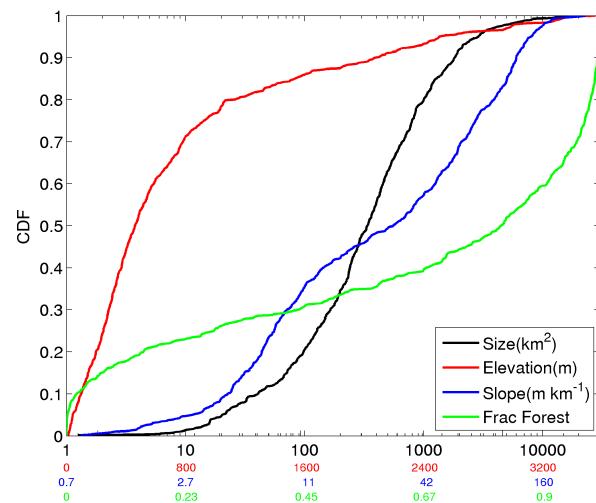
sen because of its high spatial resolution, a necessary requirement to more fully estimate spatial heterogeneity for basins in complex topography. Daymet is a daily, gridded ( $1 \times 1$  km) data set over the CONUS and southern Canada and is available from 1980 to present. It is derived solely from daily observations of temperature and precipitation. The Daymet variables used here are daily maximum and minimum temperature, precipitation, shortwave downward radiation, day length, and humidity; additionally, snow water equivalent is included (not used in this work). These daily values are estimated through the use of an iterative method dependent on local station density and the spatial convolution of a truncated Gaussian filter for station interpolation, and the Mountain Climate Simulator (MT-CLIM) to estimate shortwave radiation and humidity (Thornton et al., 1997; Thornton and Running, 1999; Thornton et al., 2000). Daymet does not include estimates of potential evapotranspiration (PET), a commonly needed input for conceptual hydrologic models or wind speed and direction. Therefore, PET was estimated using the Priestly–Taylor (P–T) method (Priestly and Taylor, 1972) and is discussed further in Sect. 3. Data quality is an ever-present issue in hydrologic modeling, and while the input data to Daymet are subject to rigorous quality control checks (Durre et al., 2008, 2010) potential errors may remain (Menne et al., 2009, 2010; Oubeidillah et al., 2013). Additionally, the Maurer et al. (2002) and National Land Data Assimilation System (NLDAS) (Xia et al., 2012) 12 km gridded data sets were processed to provide daily forcing data for the basin lumped configuration, resulting in three distinct data sets available for future forcing data impact studies.

Daily streamflow data for the HCDN-2009 gages were obtained from the USGS National Water Information System server (<http://waterdata.usgs.gov/usa/nwis/sw>) over the same forcing data time period, 1980–2010. While the period 1980–1990 is not covered by the HCDN-2009 review, it was assumed that these basins would have minimal human disturbances in this time period as well. For the portion of the basins that do not have streamflow records back to 1980, analysis is restricted to the available data records. The USGS provides streamflow data flags to identify periods of estimated flow and are included here. However, other data quality information is unavailable without further investigation and not available in this data set. For reference, 90 % (604) of the basins have 20 % or fewer flow days estimated and 75 % (503 basins) have 10 % or less flow values estimated.

The 671 basins span the entire CONUS and cover a wide range of hydroclimatic conditions. They range from wet, warm basins in the southeastern (SE) US to hot and dry basins in the southwestern (SW) US, to wet, cool basins in the northwestern (NW) and dry, cold basins in the intermountain (Rocky Mountains in Fig. 1a) western US. Figure 1b displays the basin annual precipitation (colored shading) along with symbols to denote rain- and snow-dominated basins. In terms of annual mean CDFs (cumulative density functions),



**Figure 2.** Annual CDFs of runoff ( $\text{mm day}^{-1}$ ) (black, bottom  $x$  axis), precipitation ( $\text{mm day}^{-1}$ ) (blue, bottom  $x$  axis), potential evapotranspiration ( $\text{mm day}^{-1}$ ) (green, bottom  $x$  axis), and temperature ( $^{\circ}\text{C}$ ) (red, top  $x$  axis).



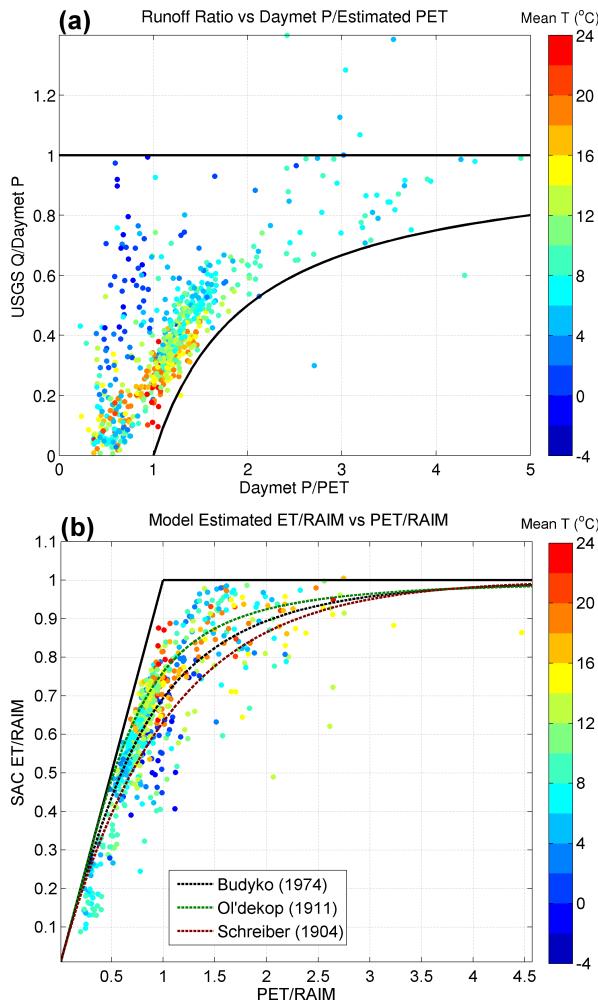
**Figure 3.** Cumulative density functions of basin size ( $\text{km}^2$ ) (black), basin mean elevation (m) (red), mean slope ( $\text{m km}^{-1}$ ) (blue), and fractional forest cover (green) for the basin set.

Daymet-estimated basin mean temperatures range from  $-2$  to  $23\text{ }^{\circ}\text{C}$  with precipitation amounts of  $0.7\text{--}9.4\text{ mm day}^{-1}$  (Fig. 2). Annual observed mean runoff ranges from  $0.01$  to  $9.3\text{ mm day}^{-1}$  with PET estimates ranging from  $1.9$  to  $4.8\text{ mm day}^{-1}$ . Interestingly, this implies that Daymet precipitation itself is not enough to balance the observed runoff in some basins and is consistent with other recent large-sample hydrologic studies (Oubeidillah et al., 2013). Seasonal variations in these four variables are large as well, with some basins reaching mean winter time temperatures lower than  $-10\text{ }^{\circ}\text{C}$  and summer time mean temperatures higher than  $25\text{ }^{\circ}\text{C}$  (not shown). The seasonal water balance varies greatly with some basins experiencing much higher precipitation and runoff rates in one season versus another (e.g., spring runoff peaks in mountain snowmelt-dominated basins). As expected, PET varies seasonally with a minimum in winter and a maximum in summer.

Figure 3 gives CDFs for various physical descriptors of the basin set. The basins range in size from roughly  $1$  to  $25\,800\text{ km}^2$  with the median basin size being about  $335\text{ km}^2$  and have mean elevations spanning from nearly sea level ( $10\text{ m}$ ) to high alpine elevations ( $3570\text{ m}$ ) with a median elevation of  $462\text{ m}$ . Notably,  $75$  basins have mean elevations  $>2000\text{ m}$ . Corresponding to the large range of elevations in the basin set, the mean slopes vary considerably, spanning over  $2$  orders of magnitude from near zero to over  $200\text{ m km}^{-1}$ . The basin set covers a wide range of basin shapes with aspect ratios ranging from  $0.08$  to about  $11$ . Finally, there is a large range of forest covers across the basin set which may have implications for hydrologic similarity (Oudin et al., 2010) with  $20\%$  of the basins having less than

(more than)  $14\%$  ( $98\%$ ) forest cover and the median basin having about  $80\%$  forest cover (NLCD-2011).

This basin set allows us to simulate a variety of energy- and water-limited basins with different snow storage, elevation, slope, and precipitation characteristics. Figure 4a shows runoff ratio (USGS streamflow / Daymet precipitation) versus the aridity index (Daymet Precipitation / PET). Immediately, it can be seen that some basins lie above the water limit line ( $Y = 1$ ) indicating more runoff than precipitation, and many basins are near it ( $Y > 0.9$ ). In these cases the model calibration process would struggle to produce an unbiased calibration, or never in basins above the water limit, because the basic water balance requires nearly zero evapotranspiration (ET) or is not satisfied. This requires a modification to incoming precipitation, which is discussed in the next section. Not coincidentally, the basins near and above the water limit are colder basins (mean annual  $T < 10\text{ }^{\circ}\text{C}$ ) with frozen precipitation during colder months. Additionally, two basins lie to the right of the curved line ( $Y = 1 - 1/\text{aridity}$ ) indicating a surplus of water. These basins may also require modifications to input precipitation, but it is less clear in this case as observations of precipitation are generally underestimates, especially for snowfall (e.g., Yang et al., 1998). Examining the basin set using model output terms in the Budyko framework, there are many energy-limited basins with dryness ratios as small as  $0.2$  and many water-limited basins with model estimated dryness ratios as large as  $4.5$  (Fig. 4b). Note that now no basins lie above the water limit, indicating bulk precipitation corrections were applied as needed during the calibration process. Examination of hydrometeorological forcing data sets across a large spatial extent through the lens of water and energy balance draws attention to gross errors



**Figure 4.** (a) Runoff ratio of observed runoff to Daymet-estimated precipitation versus ratio of Daymet-estimated precipitation to Priestly–Taylor-estimated PET. (b) Model-derived Budyko analysis using model ET, PET and total surface water input (rain plus melt, RAIM) for the 671 basins and three derivations of the Budyko curve (dashed lines). Basin mean temperature is shaded (coloring) in both panels.

in the forcing or streamflow data sets and permits any identified errors to be placed into spatial and temporal context, a benefit of large-sample studies.

As noted above, no additional quality control was performed on the candidate basins before calibration. For completeness and to more fully highlight some of the benefits and tradeoffs made when performing large-sample hydrologic studies, all basins are kept for analysis in this work.

### 3 Hydrologic modeling benchmark

As stated in the introduction, the intended purpose of this data set is a test bed to facilitate assessment of hydrologic modeling and prediction questions across broad hydroclimatic variations, and we focus here on providing a benchmark performance assessment for a widely used calibrated, conceptual hydrologic modeling system. This type of data set can be used for many applications including evaluation of new modeling systems against a well known benchmark system over wide ranging conditions, or as a base for comprehensive predictability experiments exploring the importance of meteorology or initial basin conditions. To this end, we have implemented and tested an initial model and calibration system described below, using the primary models and objective calibration approach that have been used by the US National Weather Service River Forecast Centers (NWSRFCs) in service of operational short-term and seasonal streamflow forecasting.

#### 3.1 Models

The HCDN-2009 basins include those with substantial seasonal snow cover (Fig. 1b), necessitating a snow model in addition to a hydrologic model. Within the NWSRFCs, the coupled Snow-17 and SAC-SMA system is used. Snow-17 is a conceptual air-temperature-index-based snow accumulation and ablation model (Anderson, 1973). It uses near-surface air temperature to determine the energy exchange at the snow-air interface and the only time-varying inputs are typically air temperature and precipitation (Anderson, 1973, 2002). The SAC-SMA model is a conceptual hydrologic model that includes representation of physical processes such as evapotranspiration, percolation, surface flow, and subsurface lateral flow. Required inputs to SAC-SMA are potential evapotranspiration and water input to the soil surface (Burnash et al., 1973; Burnash, 1995). Snow-17 runs first and determines the partition of precipitation into rain and snow and the evolution of the snowpack. Any rain, snowmelt or rain passing unfrozen through the snowpack for a given time step becomes direct input to the SAC-SMA model. Finally, streamflow routing is accomplished through the use of a simple two-parameter, Nash-type instantaneous unit-hydrograph model (Nash, 1957).

#### 3.2 Calibration

We employed a split-sample calibration approach following Klemes (1986): assigning the first 15 years of available streamflow data for calibration and the remainder for validation, then repeating the calibration using the last 15 years and the initial remaining period for validation; thus, approximately 5500 daily streamflow observations were used for each calibration. To initialize the model calibration moisture states on 1 October, we specified an initial wet SAC-SMA

soil moisture state that was allowed to spin down to equilibrium for a given basin by running the first year of the calibration period repeatedly and assumed no initial snowpack. This was done until all SAC-SMA state variables had minimal year over year variations, which is a spin-up approach used by the Project for Intercomparison of Land-Surface Process Schemes (e.g., Schlosser et al., 2000). Determination of optimal calibration sampling and spin-up procedures is an area of active research. Spin-up was performed for every parameter set specified by the optimization algorithm, then the model was integrated for the calibration period and the RMSE (root mean square error) for that parameter set was calculated.

Objective calibration was done by minimizing the RMSE of daily modeled runoff versus observed streamflow using the SCE global search algorithm of Duan et al. (1992, 1993). The SCE algorithm uses a combination of probabilistic and deterministic optimization approaches that systematically spans the allowed parameter search space and also includes competitive evolution of the parameter sets (Duan et al., 1993). Prior applications to the SAC-SMA model have shown good results (Sorooshian et al., 1993; Duan et al., 1994). In the coupled Snow-17 and SAC-SMA modeling system, 35 potential parameters are available for calibration, of which we calibrated 20 parameters having either a priori estimates (Koren et al., 2000) or those found to be most sensitive following Anderson (2002) (Table 1). The SCE algorithm was run using 10 different random seed starts for the initial parameter sets for each basin, in part to evaluate the robustness of the optimum in each case, and the optimized parameter set with the minimum RMSE from the 10 different optimization runs was chosen for evaluation.

For Snow-17, six parameters were chosen for optimization (Table 1): the minimum and maximum melt factors (MFMIN, MFMAX), the wind adjustment for enhanced energy fluxes to the snowpack during rain on snow (UADJ), the rain/snow partition temperature, which may not be 0 °C (PXTEMP), the snow water equivalent for 100 % snow covered area (SI), and the gauge catch correction term for snowfall only (SCF). These six parameters were chosen because MFMIN, MFMAX, UADJ, SCF, and SI are defined as major model parameters by Anderson (2002). PXTEMP was also shown to be important in the Snow-17 model by Mizukami et al. (2013). The SCF is critical in many snow-dominated basins as precipitation is generally underestimated in these types of basins (e.g., Yang et al., 1998) and is certainly underestimated in some basins in Daymet as shown in Figs. 3 and 4.

The areal depletion curve (ADC) is considered a major parameter in Snow-17. However, to avoid expanding the parameter space by the number of ordinates on the curve (typically 10), we manually specified the ADC according to regional variations in latitude, topographic characteristics (e.g., plains, hills or mountains) and typical air mass characteristics (e.g., maritime polar, continental polar) (as suggested in Anderson, 2002). The remaining Snow-17 parame-

ters were set in the same manner. Following the availability of a priori parameter estimates for SAC-SMA from a variety of data sets and various calibration studies with SAC-SMA (Koren et al., 2000; Anderson et al., 2006; Pokhrel and Gupta, 2010; Zhang et al., 2012), 11 parameters from SAC-SMA are included for calibration (Table 1). We use an instantaneous unit hydrograph, represented as a two-parameter gamma distribution for streamflow routing (Sherman, 1932; Clark, 1945; Nash, 1957; Dooge, 1959), the parameters of which were inferred as part of calibration..

Finally, the scaling parameter in the Priestly–Taylor PET estimate is also calibrated. The P–T equation (Priestly and Taylor, 1972) can be written as

$$\text{PET} = \frac{a}{\lambda} \cdot \frac{s \cdot (R_n - G)}{s + \gamma}. \quad (1)$$

Where  $\lambda$  ( $\text{MJ kg}^{-1}$ ) is the latent heat of vaporization,  $R_n$  ( $\text{MJ m}^{-2} \text{ day}^{-1}$ ) is the net radiation estimated using day of year, all Daymet variables and equations to estimate the various radiation terms (Allen et al., 1988; Zotarelli et al., 2009),  $G$  ( $\text{MJ m}^{-2} \text{ day}^{-1}$ ) is the soil heat flux (assumed to be zero in this case),  $s$  ( $\text{kPa}^\circ\text{C}^{-1}$ ) is the slope of the saturation vapor pressure–temperature relationship,  $\gamma$  ( $\text{kPa}^\circ\text{C}^{-1}$ ) is the psychrometric constant and  $a$  (unitless) is the P–T coefficient. The P–T coefficient replaces the aerodynamic term in the Penman–Monteith equation and varies by the typical conditions of the area where the P–T equation is being applied with humid forested basins typically having smaller values and exposed arid basins having larger values (Shuttleworth and Calder, 1979; Morton, 1983; Jensen et al., 1990). Thus, the P–T coefficient was included in the calibration since it should vary from basin to basin.

## 4 Benchmark results

### 4.1 Assessment objectives and metrics

Assessment of the models will focus on overall performance across the basin set, regional variations, and error characteristics. Nash–Sutcliffe efficiency (NSE) (Nash and Sutcliffe, 1970) and two of the decomposition components of NSE, variance bias ( $\alpha$ ) and total volume bias ( $\beta$ ) (Gupta et al., 2009), are the first metrics examined in two variations. Because NSE scores model performance relative to the observed climatological mean, regions in which the model can track a strong seasonal cycle (large flow autocorrelation) perform relatively better when measured by NSE, and this seasonal enhancement may be imparted when using NSE as the objective function for both the calibration and validation phases (e.g., Schaeefli and Gupta, 2007). Additionally, basins with higher streamflow variance and frequent precipitation events have better model performance. Therefore, to give a more standardized picture of model performance across varying hydroclimatologies, the NSE was recomputed using

**Table 1.** Table describing all parameters calibrated and their bounds for calibration.

Parameter	Description	Units	Calibration range
Snow-17			
MFMAX	Maximum melt factor	$\text{mm } ^\circ\text{C}^{-1} 6\text{ h}^{-1}$	0.8–3.0
MFMIN	Minimum melt factor	$\text{mm } ^\circ\text{C}^{-1} 6\text{ h}^{-1}$	0.01–0.79
UADJ	Wind adjustment for enhanced flux during rain on snow	$\text{km } 6\text{ h}^{-1}$	0.01–0.40
SI	SWE for 100 % snow covered area	mm	1.0–3500.0
SCF	Snow gauge undercatch correction factor	—	0.1–5.0
PXTEMP	Temperature of rain/snow transition	$^\circ\text{C}$	–1.0–3.0
SAC-SMA			
UZTWM	Upper zone tension water maximum storage	mm	1.0–800.0
UZFWM	Upper zone free water maximum storage	mm	1.0–800.0
LZTWM	Lower zone tension water maximum storage	mm	1.0–800.0
LZFPM	Lower zone free water primary maximum storage	mm	1.0–1000.0
LZFSM	Lower zone free water secondary maximum storage	mm	1.0–1000.0
UZK	Upper zone free water lateral depletion rate	$\text{day}^{-1}$	0.1–0.7
LZPK	Lower zone primary free water depletion rate	$\text{day}^{-1}$	0.00001–0.025
LZSK	Lower zone secondary free water depletion rate	$\text{day}^{-1}$	0.001–0.25
ZPERC	Maximum percolation rate	—	1.0–250.0
REXP	Exponent of the percolation equation	—	0.0–6.0
PFREE	Fraction percolating from upper to lower zone free water storage	—	0.0–1.0
Others			
USHAPE	Shape of unit hydrograph	—	1.0–5.0
USCALE	Scale of unit hydrograph	—	0.001–150.0
P-T	Priestly–Taylor coefficient	—	1.26–1.74

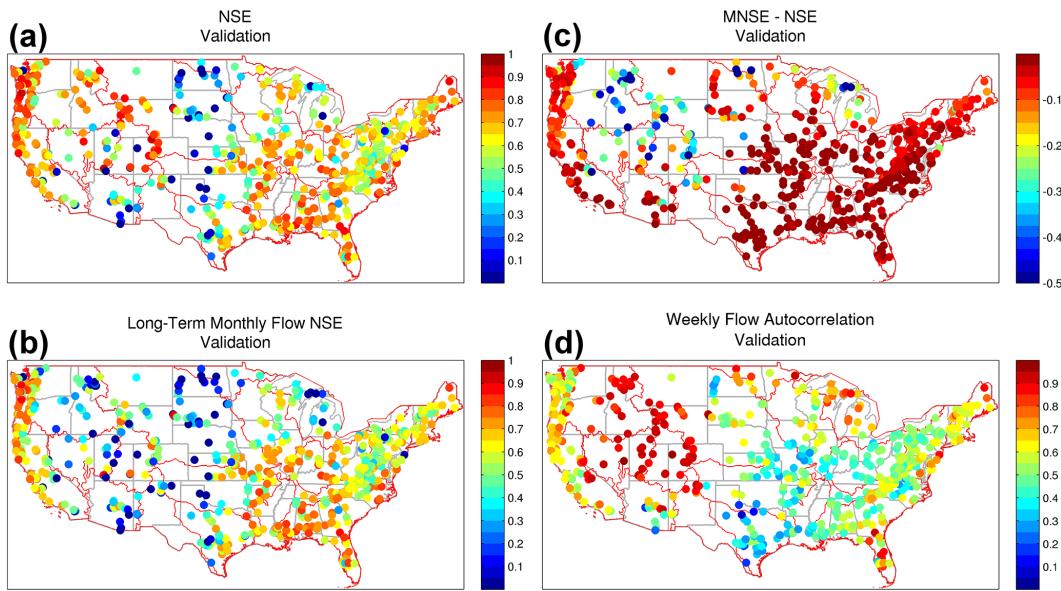
the long-term monthly mean flow instead of mean flow (denoted MNSE hereafter), thus preventing climatological seasonality from inflating the NSE and more accurately ranking basins by the degree to which the model added value over climatology in response to weather events (Garrick et al., 1978; Martinec and Rango, 1989; Schaeefli et al., 2005). MNSE in this context is defined for each day of year (DOY) via a 31-day window centered on a given DOY. The long-term flow for that 31-day “month” is computed giving rise to a “monthly” mean flow. Using this type of climatology as the base for an NSE-type analysis provides improved standardization in basins with large flow autocorrelations. This definition is similar to the one proposed by Garrick et al. (1978) but with the addition of the 31-day smoother, which is done to provide a smoother reference climatology.

Also, several other advanced, more physically based, metrics of model performance are provided. First, three diagnostic signatures based on the flow duration curve (FDC) from Yilmaz et al. (2008) are computed: (1) the top 2 % flow bias, (2) the bottom 30 % flow bias, and (3) the bias of the slope of the middle portion (20–70 percentile) of the FDC. Second, examination of the time series of squared error contribution to the RMSE statistic was performed to highlight events in which the model performs poorly following Clark et al. (2008). This analysis was performed to gauge

the representativeness of performance metrics over the model record by using the sorted (highest to lowest) time series of squared error to identify the  $N$  number of the largest error days and determine their fractional error contribution to the total. Finally, we extend this analysis to introduce a simple, normalized general error index for application and comparison across varying modeling and calibration studies. We coin the index, E50, the fraction of calibration points contributing 50 % of the error. This captures the number of points determining the majority of the error and thus the optimal parameter set.

#### 4.2 Spatial variability

It is informative to examine spatial patterns of the aforementioned metrics to elucidate factors leading to weak (and strong) model performance. This also allows for identification of outlier basins and characterization of contributing factors (i.e., forcing or streamflow data issues or poor calibration). Poorly performing basins are most common along the high plains and desert southwest (Fig. 5a, Sect. 3c). When examining MNSE (Fig. 5b), basins with high nonseasonal streamflow variance and frequent precipitation events (SE and NW US) have the highest model MNSE, while most of the snowmelt-dominated basins see MNSE scores reduced



**Figure 5.** (a) Spatial distribution of NSE, (b) NSE using MNSEs rather than the long-term mean flow, (c) MNSE – NSE for the validation period, and (d) weekly flow autocorrelation.

relative to NSE, particularly in the validation phase (Fig. 5c). This indicates that RMSE as an objective function may not be well suited for model calibration in basins with high flow autocorrelation (Kavetski and Fenicia, 2011; Evin et al., 2014). This is confirmed by comparing Fig. 5d to Fig. 5c, basins with large flow autocorrelations (1 week mean flow for example) generally have lower MNSE scores.

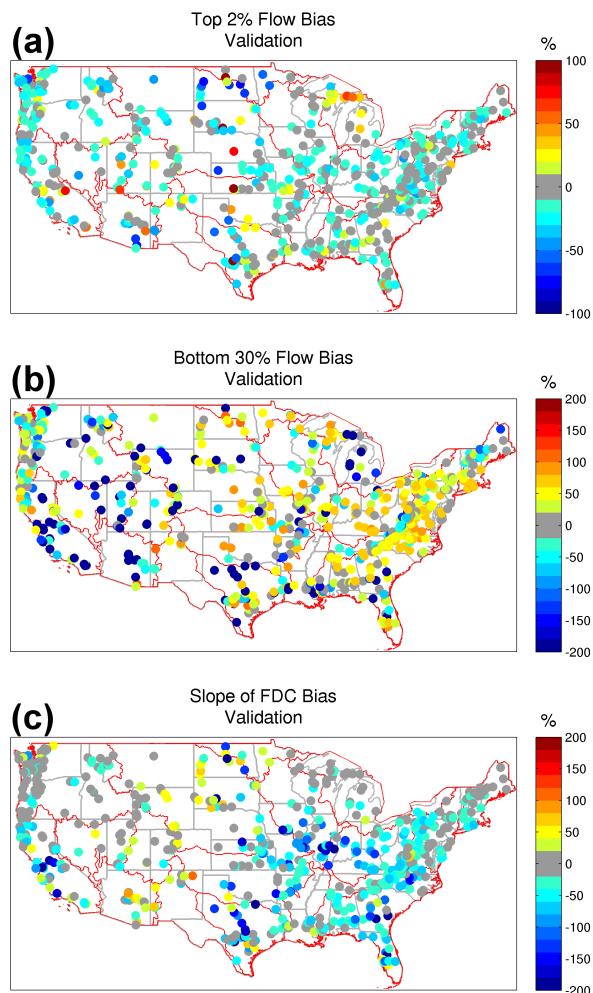
Areas with low-validation NSE and MNSE scores have generally large biases when looking at FDC metrics as well (Fig. 6). Focusing on the high plains, high flow biases of  $\pm 50\%$  are common. Extreme negative low flow biases are also present along the high plains and desert SW along with a general model trend to have large negative FDC slope biases, consistent with a poorly calibrated model. For the 72 % of basins with validation NSE  $> 0.55$  (basins with yellow-green to dark red colors in Fig. 6a), there is no noticeable spatial pattern across the CONUS in regard to high flow periods. However, basins with a more pronounced seasonal cycle (e.g., snowpack-dominated watersheds, central west coast) generally have a negative low flow bias, while basins with a smaller seasonal cycle have a positive low flow bias (Fig. 6b). Correspondingly, basins with a pronounced seasonal cycle generally have a near zero or positive slope of the FDC bias, while basins with a smaller seasonal cycle have a negative slope bias (Fig. 6c).

Past applications with similar conceptual snow and hydrologic modeling systems across the CONUS have shown comparable spatial performance patterns. Clark et al. (2008) applied many conceptual models to a subset of the MOPEX basin set and found poor performance in arid regions. Mar-

tinez and Gupta (2010), using a monthly water balance model, found the best performance generally along the east coast, most of SE CONUS, and along the west coast with scattered good performance in the Rocky Mountains. They found that many basins along the high plains and north side of the Appalachian Mountains perform poorly. They also note that arid regions have high variability error (variability bias term in KGE – Kling–Gupta efficiency).

#### 4.3 Cumulative performance

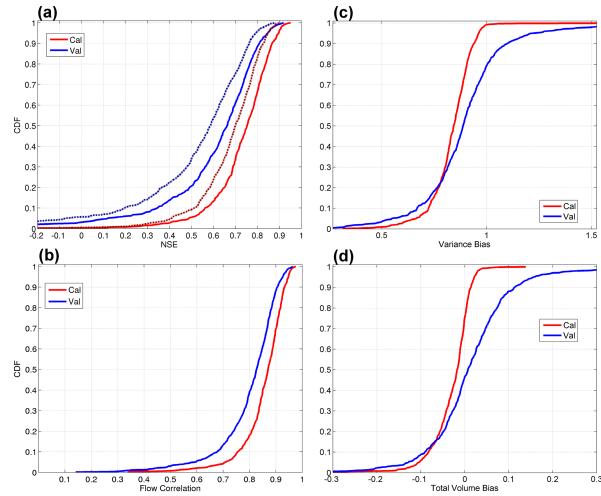
Two basic cumulative thresholds for model performance are highlighted here, NSE values of 0.55 and 0.8. An NSE of 0.55 indicates some model skill, and an NSE of 0.8 suggests reasonably good model performance. For the calibration period, 90 % (604) of the basins have a NSE greater than 0.55, while 72 % (484) of the basins had a validation period NSE  $> 0.55$  (Fig. 7a). At the NSE  $> 0.8$  level, 34 % (225) of the basin models perform better during calibration and 12 % (78) of the basin models meet that criteria during the validation phase. When using MNSE, 85 and 57 % (568 and 385) of the basins lie above 0.55, and 1 and 4 % (114 and 29) of the basins lie above 0.8 during the calibration and validation phases. The decomposition of the NSE (Gupta et al., 2009) shows that 90 % of the basins have a calibration (validation) model–observation flow correlation  $> 0.75$  (0.68) and 30 % (12 %) of the basins have a model–observation flow correlation  $> 0.9$  (Fig. 7b). However, nearly all basins have too little modeled variance (values less than one) for both the calibration and validation phases (Fig. 7c). The total volume biases are generally small with 94 % (79 %) of the basins having a



**Figure 6.** (a) Spatial distribution of the high flow bias, (b) low flow bias, and (c) flow duration curve bias for the validation period.

calibration (validation) period total flow bias within 10 % of the observed flow (Fig. 7d). These are expected results when using RMSE for the objective function (Gupta et al., 2009) and reaffirm that our implementation of the SCE calibrates the model properly.

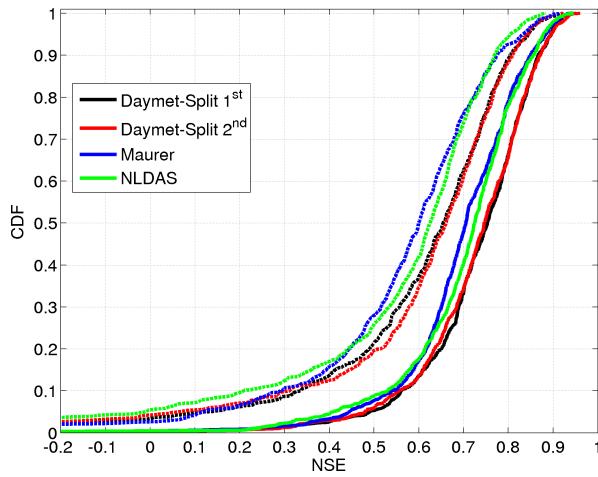
Figure 8 highlights the full split sample approach for calibration following Klemes (1986). It is seen that the calibration and validation statistics give quite similar results regardless of which time period is used for calibration and validation using the Daymet data. This could indicate that both halves of the data are equally challenging to model with this modeling system. We have also included basin calibrations using only the first 15 years for the Maurer et al. (2002) and NLDAS-II (Xia et al., 2012) data sets. It can be seen that the Daymet forcing provides better model performance overall than both Maurer et al. and NLDAS forcing data. This likely relates to the coarser resolution of the Maurer et



**Figure 7.** (a) CDFs of the model NSE (solid) for the calibration (red) and validation periods (blue) and NSE using the MNSEs (dark shaded and dashed). CDFs for (b) simulated–observed flow correlation in the decomposition of the NSE, (c) for the variance bias in the decomposition of the NSE, and (d) total volume bias in the decomposition of the NSE.

al. (2002) and NLDAS data (12 km) and the somewhat small basin sizes in this basin set. More importantly, the inclusion of the Klemes (1986) split-sample approach provides users of this data set two parameter estimates for each basin using different calibration periods, while the inclusion of three total forcing data sets begins to allow for ensemble-type forcing data impact studies across a large basin sample size. In the remaining discussion, only model performance results using the first half of the split sample for calibration are presented.

With respect to advanced diagnostics, the model underpredicts high flow events in nearly all basins during calibration and slightly less so for the validation period (Fig. 9a). This is an expected result when using RMSE as the objective function because the optimal calibration underestimates flow variability (Gupta et al., 2009). Low flow periods are more evenly over- and underpredicted (Fig. 9b) for both the calibration and validation time frames with 58 and 61 % of basins having more modeled low flow. Finally, the bias in the slope of the FDC is generally underpredicted with about 75 % of the basins having a negative model bias (FDC slope is negative, thus a negative bias indicates the model slope is more positive and that the modeled flow variability is too compressed). The slope of the FDC indicates the variance of daily flows, which primarily relate to the seasonal cycle or the “flashiness” of a basin. Again this indicates model variability is less than that observed, at both short and longer timescales. In aggregate, these results agree with Fig. 5 and are expected based on the analysis of Gupta et al. (2009). Optimization using RMSE or NSE as the objective function generally results in underprediction of flow variance and



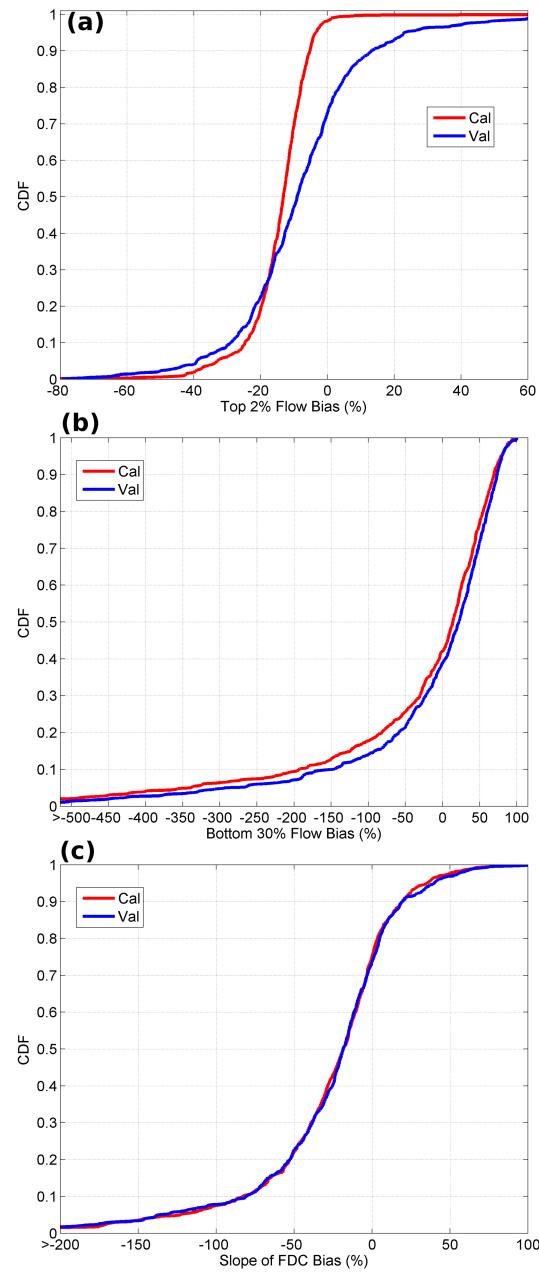
**Figure 8.** Cumulative density functions for model Nash–Sutcliffe efficiency for the calibration (solid) and validation (dashed) periods using three different forcing data sets (Daymet, Maurer, NLDAS). The Daymet data set was calibrated using the first 15 years (Split 1<sup>st</sup>) and validated against the remaining data and also calibrated using the last 15 years (Split 2<sup>nd</sup>) and validated against the initial streamflow data. Maurer and NLDAS calibrations performed using the first 15 years of observed streamflow only.

near-zero total flow bias (Fig. 7). This manifests itself in the simulated hydrograph as underpredicted high flows, generally overpredicted low flows and a more positive slope to the middle portion of the FDC (Fig. 9). It is worth repeating that the goal of this initial application is to provide to community with a benchmark of model performance using well known models, calibration systems and widely used, simple objective functions, thus the use of RMSE.

#### 4.4 Error characteristics

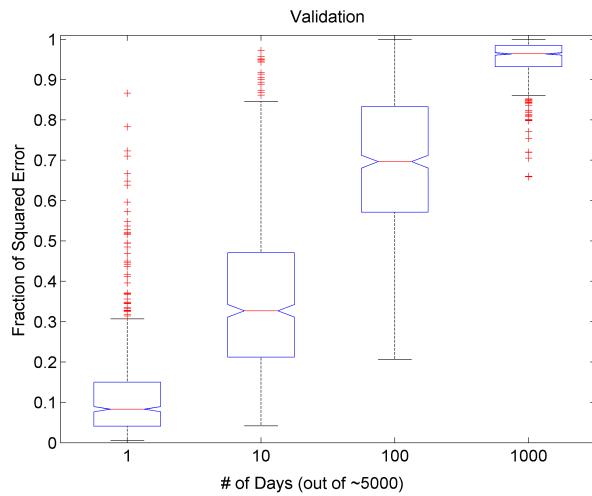
When examining fractional error statistics for the basin set, 15 basins have single days that contribute at least half the total squared error (potential outlier basins), whereas at the median, the largest error day contributes 8.3 % of the total squared error for the median basin (Fig. 10). The fractional error contribution for the 10, 100 and 1000 largest error days for the median basin are 33, 70 and 96 % of the total squared error respectively. This indicates that for nearly all basins, there are 100 or fewer points that drive the RMSE and therefore optimal model parameters. This type of analysis can be undertaken for any objective function to identify the most influential points and allow for more in-depth examination of forcing data, streamflow records, and calibration strategies (i.e., Kavetski et al., 2006; Vrugt et al., 2008; Beven and Westerberg, 2011; Beven et al., 2011; Kauffeldt et al., 2013), or if different model physics are warranted.

The spatial distribution of fractional error contributions show that the issue of model performance being explained



**Figure 9.** (a) CDFs for model high flow bias for the calibration (red) and validation periods (blue), (b) model low flow bias, and (c) model flow duration curve slope bias.

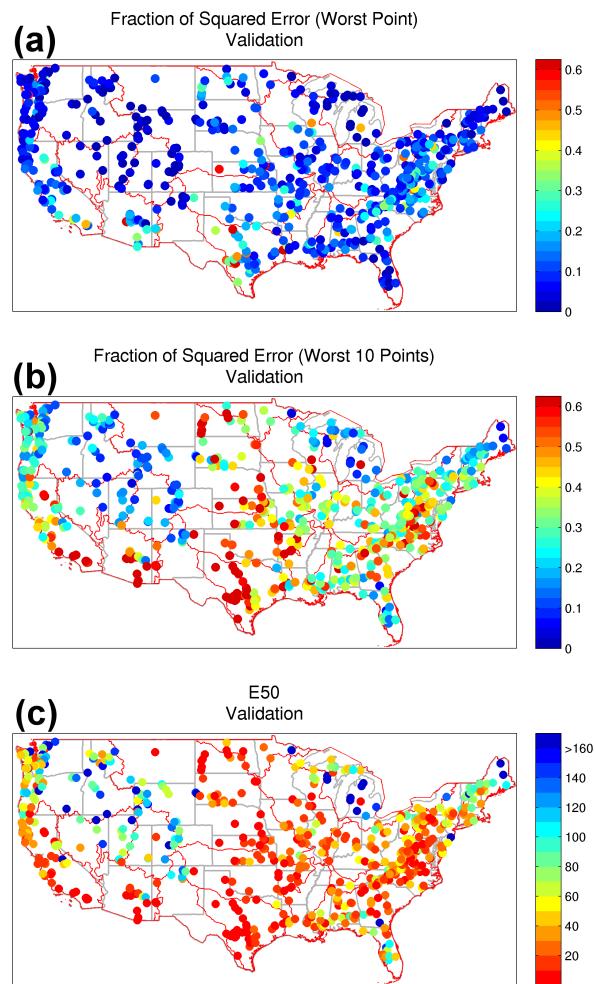
by a relatively small set of days is more prevalent in arid regions of the CONUS (desert SW US and high plains) as well as basins slightly inland from the east coast of the CONUS (Fig. 11a, b). The arid basins are generally dry with sporadic high precipitation (and flow) events, while the Appalachian basins are wetter (Fig. 1b) with extreme precipitation events



**Figure 10.** Fractional contribution of the total squared error for the 1, 10, 100, 1000 largest error days. The box plots represent the 671 basins with the blue area defining the interquartile range, the whiskers representing reasonable values and the red crosses denoting outliers. The median is given by the red horizontal line with the notch in the box denoting the 95 % confidence interval of the median value.

interspersed throughout the record. Basins with significant snowpack tend to have lower error contributions from the largest error days (Fig. 11a, b). The E50 metric highlights mean peak snow water equivalent (SWE) and frequent precipitation basins as well. These regions contain and order of magnitude more days than the high plains and desert SW, giving insight into how representative of the entire streamflow time series the optimal model parameter set really is.

Additionally, ranking the basins using their fractional error characteristics provides a similar insight. As the aridity index increases, the fractional error contribution increases for basins with little to no mean peak SWE. For basins with significant SWE, the fractional error contribution decreases with increasing aridity (Fig. 12). Alternatively, for a given aridity index the fractional error contribution for  $N$  days will decrease with increasing SWE. This dynamic arises because more arid basins with SWE produce a relatively greater proportion of their runoff from snowmelt, without intervening rainfall. This implies that the optimized model produces a more uniform error distribution with less heteroscedasticity in basins with more SWE. Moreover, as the fractional error contribution for the 10 largest error days increases, model NSE generally decreases in the validation phase (Fig. 13). This indicates fractional error metrics are related to overall model performance and that calibration methods to reduce extreme error days should improve model performance. This is not unexpected due to the fact that the residuals from an RMSE-type calibration are heteroscedastic. Arid basins typically have few high flow events, which are generally subject

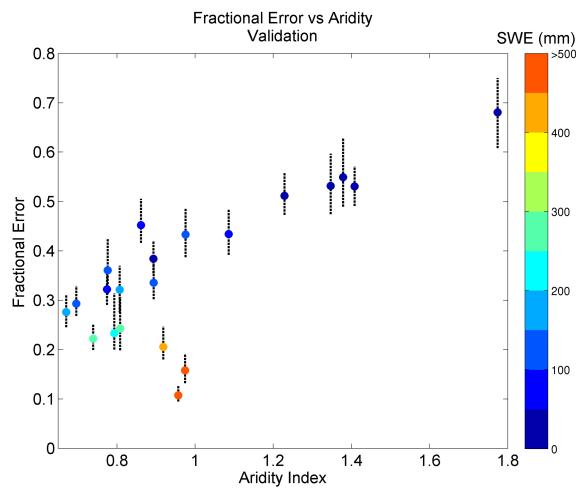


**Figure 11.** (a) Spatial distribution of the fractional contribution of total squared error for the largest day during the validation period, the (b) 10 largest error days, and (c) the number of days contributing 50 % of the total objective function error, E50.

to larger errors when minimizing the RMSE. Using advanced calibration methodologies that account for heteroscedasticity (Kavetski and Fenicia, 2011; Evin et al., 2014) may produce improved calibrations for arid basins in this basin set and provide different insights into model behavior using this type of analysis.

#### 4.5 Limitations and uncertainties

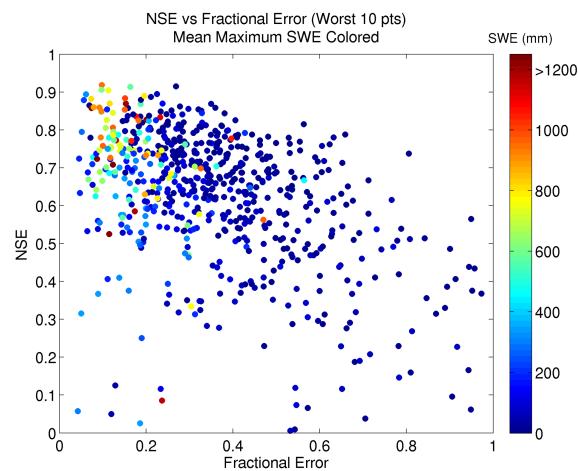
One interesting example of the usefulness (and a potential limitation) of large-sample hydrology stemming from this work lies in the identification of issues with forcing data sets. Figures 3 and 4 show Daymet has too little precipitation in certain regions, which is also seen in Oubeidillah et al. (2013). When examining calibrated model performance in



**Figure 12.** Ranked fractional squared error contribution for the 100 largest error days for the 671 basins versus the aridity index with mean maximum SWE shaded. Each dot represents a 32-basin bin defined by the rank of the fractional error contribution for the 100 largest error days for all basins. The dashed vertical black lines denote the 95 % confidence interval for the mean of the fractional error contribution for a given bin.

the Pacific northwest, it is seen that several basins along the west coast have low outlier NSE scores. Tracing this unexpected result, we find the Daymet forcing data available for those basins has a negative temperature bias, preventing mid-winter rain and melt episodes in the modeling system, identifying scope for improvement in the Daymet forcing. Moreover, winter periods of observed precipitation and streamflow rises coincide with subzero  $T_{\max}$  in the Daymet data set, also suggesting areas to improve the Daymet forcing. The large-sample of basins in this region (91) allowed for identification of the outlier basins and the underlying causes.

This may also limit interpretation of these results and other large-sample hydrologic studies. As noted by Gupta et al. (2014), large-sample hydrology requires a tradeoff between breadth and depth. The lack of depth may inhibit discovery and identification of all data quality issues and the underlying causes of outliers in any analysis (e.g., Fig. 13). Explanation of these outliers is sometimes difficult and not complete in the initial development and analysis due to the lack of familiarity with specific basins and any forcing or validation data peculiarities. However, providing forcing data, model parameters and model output permits additional focused studies and helps reduce these limitations. Additional prescreening using the methods of Martinez and Gupta (2011) can also help identify outliers due to data quality issues and help identify basins and regions where model physics errors are present.



**Figure 13.** Nash–Sutcliffe efficiency versus the fractional error of the 10 largest error days for the validation period for all basins with basin mean peak snow water equivalent (mm) colored.

## 5 Summary and discussion

Most hydrologic studies focus in detail on a small number of watersheds, providing comprehensive but highly local insights, and may be limited in their ability to inform general hydrologic concepts applicable across regions (Gupta et al., 2014). To facilitate large-sample hydrologic studies, large-sample basin data sets and corresponding benchmarks of model performance using standard methodology across all basins need to be freely available to the community. To that end, we have compiled a community data set of daily forcing and streamflow data for 671 basins and provide a benchmark of performance using a widely used conceptual hydrologic modeling and calibration scheme over a wide range of conditions.

Overall, application of the basin set to assessing an objectively calibrated conceptual hydrologic model representation of the 671 watersheds yielded calibration NSE scores of  $> 0.55$  (0.8) for 90 % (34 %) of the basins. Performance of the models varied regionally, and the main factors influencing this variation were found to be aridity and precipitation intermittency, contribution of snowmelt, and runoff seasonality. Analysis of the cumulative fractional error contributions from the largest error days showed that the presence of significant SWE offset the negative impact of increasing aridity on simulation performance. This study has identified potential outlier basins for this modeling system and has provided insights into potential forcing data limitations. Although this modeling application utilized a conceptual hydrologic model with a single-objective calibration strategy, the findings provide a baseline for assessing more complex strategies in each area, including multiobjective calibration of more highly distributed hydrologic models (e.g., in Shi et al., 2008). The unusually broad variation of hydroclimatologies represented by

the data set, which contains forcing and streamflow data obtained by consistent methodology and retains outlier basins, makes it a notable resource for these and other future large-sample watershed-scale hydrologic analysis efforts.

This data set and the applications presented are made available to the community (see <http://ral.ucar.edu/projects/hap/flowpredict/subpages/modelvar.php> or <http://dx.doi.org/10.5065/D6MW2F4D>).

**Acknowledgements.** This work is funded by the US Army Corps of Engineers Climate Preparedness and Resilience Programs and the US Department of the Interior Bureau of Reclamation. The authors would like to thank the USGS Modeling of Watershed Systems (MoWS) group, specifically for providing technical support and the national geospatial fabric data to generate all the basin spatial configurations. We would also like to thank Jordan Read and Tom Kunicki of the USGS Center for Integrated Data Analytics for their help with the USGS Geodata Portal.

Edited by: S. Archfield

## References

- Allen, R. G., Pereira, L. S., Raes, D., and Smith, M.: Crop evapotranspiration: guidelines for computing crop water requirements. Food and Agriculture Organization of the United Nations, Rome, 15 pp., 1988.
- Anderson, E. A.: National Weather Service River Forecast System – Snow accumulation and ablation model. NOAA Technical Memorandum, NWS, HYDRO-17, US Department of Commerce, Silver Spring, MD, 217 pp., 1973.
- Anderson, E. A.: Calibration of conceptual hydrologic models for use in river forecasting. NOAA Technical Report, NWS 45, Hydrology Laboratory, Silver Spring, MD, 2002.
- Anderson, R. M., Koren, V. I., and Reed, S. M.: Using SSURGO data to improve Sacramento Model a priori parameter estimates. *J. Hydrol.*, 320, 103–116, 2006.
- Andreassian, V., Oddos, A., Michel, C., Anctil, F., Perrin, C., and Loumange, C.: Impact of spatial aggregation of inputs and parameters on the efficiency of rainfall-runoff models: A theoretical study using chimera watersheds, *Water Resour. Res.*, 40, W05209, doi:10.1029/2003WR002854, 2004.
- Beldring, S., Engeland, K., Roald, L. A., Sælthun, N. R., and Voksø, A.: Estimation of parameters in a distributed precipitation-runoff model for Norway, *Hydrol. Earth Syst. Sci.*, 7, 304–316, doi:10.5194/hess-7-304-2003, 2003.
- Beven, K. and Westerberg, I.: On red herrings and real herrings: disinformation and information in hydrological inference, *Hydrol. Process.*, 25, 1676–1680, 2011.
- Beven, K., Smith, P. J., and Wood, A.: On the colour and spin of epistemic error (and what we might do about it), *Hydrol. Earth Syst. Sci.*, 15, 3123–3133, doi:10.5194/hess-15-3123-2011, 2011.
- Blodgett, D. L., Booth, N. L., Kunicki, T. C., Walker, J. L., and Viger, R. J.: Description and testing of the geo data portal: A data integration framework and web processing services for environmental science collaboration. US Geological Survey, Open-File Report 2011-1157, 9 pp., Middleton WI, USA, 2011.
- Burnash, R. J. C.: The NWS River Forecast System – Catchment model, in: Computer Models of Watershed Hydrology, edited by: Singh, V. P., 311–366, Water Resources Publications, Highlands Ranch, Colo, 1995.
- Burnash, R. J. C., Ferral, R. L., McGuire, R. A.: A generalized streamflow simulation system conceptual modeling for digital computers, US Department of Commerce National Weather Service and State of California Department of Water Resources, 1973.
- Clark, C. O.: Storage and the unit hydrograph. *Proc. Am. Soc. Civ. Eng.*, 9, 1333–1360, 1945.
- Clark, M. P., Slater, A. G., Rupp, D. E., Woods, R. A., Vrugt, J. A., Gupta, H. V., Wagener, T., and Hay, L. E.: Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences between hydrologic models, *Water Resour. Res.*, 44, W00B02, doi:10.1029/2007WR006735, 2008.
- Dooge, J. C. I.: A general theory of the unit hydrograph, *J. Geophys. Res.*, 64, 241–256, 1959.
- Duan, Q., Sorooshian, S., and Gupta, V. K.: Effective and efficient global optimization for conceptual rainfall-runoff models, *Water Resour. Res.*, 28, 1015–1031, 1992.
- Duan, Q., Gupta, V. K., and Sorooshian, S.: A shuffled complex evolution approach for effective and efficient optimization, *J. Optimiz. Theor. Appl.*, 76, 501–521, 1993.
- Duan, Q., Sorooshian, S., and Gupta, V. K.: Optimal use of the SCE-UA global optimization method for calibrating watershed models, *J. Hydrol.*, 158, 265–284, 1994.
- Duan, Q., Schaake, J., Andreassian, V., Franks, S., Goteti, G., Gupta, H. V., Gusev, Y. M., Habets, F., Hall, A., Hay, L., Houge, T., Huang, M., Leavesley, G., Liang, X., Nasonova, O. N., Noilhan, J., Oudin, L., Sorooshian, S., Wagener, T., and Wood, E. F.: Model Parameter Estimation Experiment (MOPEX): An overview of science strategy and major results from the second and third workshops, *J. Hydrol.*, 320, 3–17, 2006.
- Durre, I., Menne, M. J., and Vose, R. S.: Strategies for evaluating quality assurance procedures, *J. Appl. Meteor. Climatol.*, 47, 1785–1791, doi:10.1175/2007JAMC1706.1, 2008.
- Durre, I., Menne, M. J., Gleason, B. E., Houston, T. G., and Vose, R. S.: Comprehensive Automated Quality Assurance of Daily Surface Observations, *J. Appl. Meteor. Climatol.*, 49, 1615–1633, doi:10.1175/2010JAMC2375.1, 2010.
- Evin, G., Thyer, M., Kavetski, D., McInerney, D., and Kuczera, G.: Comparison of joint versus postprocessor approaches for hydrological uncertainty estimation accounting for error autocorrelation and heteroscedasticity, *Water Resour. Res.*, 50, 2350–2375, doi:10.1002/2013WR014185, 2014.
- Falcone, J. A.: GAGES-II: Geospatial Attributes of Gages for Evaluating Streamflow. Digital spatial data set 2011, available at: [http://water.usgs.gov/GIS/metadata/usgswrd/XML/gagesII\\_Sept2011.xml](http://water.usgs.gov/GIS/metadata/usgswrd/XML/gagesII_Sept2011.xml) (last access: 10 October 2013), 2011.
- Falcone, J. A., Carlisle, D. M., Wolock, D. M., and Meador, M. R.: GAGES: A stream gage database for evaluating natural and altered flow conditions in the conterminous United States. *Ecology*, 91, p. 621, A data paper in Ecological Archives E091-045-D1, available at: <http://esapubs.org/Archive/ecol/E091/045/metadata.htm> (last access: 5 April 2014), 2010.

- Garrick, M., Cunnane, C., and Nash, J. E.: A criterion of efficiency for rainfall-runoff models, *J. Hydrology*, 36, 375–381, 1978.
- Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez-Barquero, G. F.: Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modeling, *J. Hydrol.*, 377, 80–91, doi:10.1016/j.jhydrol.2009.08.003, 2009.
- Gupta, H. V., Perrin, C., Blöschl, G., Montanari, A., Kumar, R., Clark, M., and Andréassian, V.: Large-sample hydrology: a need to balance depth with breadth, *Hydrol. Earth Syst. Sci.*, 18, 463–477, doi:10.5194/hess-18-463-2014, 2014.
- Jensen, M. E., Burman, R. D., and Allen, R. G.: Evapotranspiration and irrigation water requirements. American Society of Civil Engineers, ASCE Manual and Reports on Engineering Practice, 332 p., New York, NY, 1990.
- Jin, S., Yang, L., Danielson, P., Homer, C., Fry, J., and Xian, G.: A comprehensive change detection method for updating the National Land Cover Database to circa 2011, *Remote Sens. Environ.*, 132, 159–175, 2013.
- Kauffeldt, A., Halldin, S., Rodhe, A., Xu, C.-Y., and Westerberg, I. K.: Disinformative data in large-scale hydrological modelling, *Hydrol. Earth Syst. Sci.*, 17, 2845–2857, doi:10.5194/hess-17-2845-2013, 2013.
- Kavetski, D. and Fenicia, F.: Elements of a flexible approach for conceptual hydrological modeling: 2. Application and experimental insights, *Water Resour. Res.*, 47, W11511, doi:10.1029/2011WR010748, 2011.
- Kavetski, D., Kuczera, G., and Franks, S. W.: Bayesian analysis of input uncertainty in hydrological modeling: 2. Application, *Water Resour. Res.*, 42, W03407, doi:10.1029/2005WR004376, 2006.
- Klemes, V.: Operational testing of hydrological simulation models, *Hydrol. Sci. J.*, 31, 13–24, 1986.
- Koren, V. I., Smith, M., Wang, D., and Zhang, Z.: Use of soil property data in the derivation of conceptual rainfall-runoff model parameters. American Meteorological Society 15th Conference on Hydrology, Long Beach, CA, 103–106, 2000.
- Kumar, R., Samaniego, L., and Attinger, S.: Implications of distributed hydrologic model parameterization on water fluxes at multiple scales and locations, *Water Resour. Res.*, 49, 360–379, doi:10.1029/2012WR012195, 2013.
- Lins, H. F.: USGS Hydro-Climatic Data Network 2009 (HCDN-2009), US Geological Survey, Fact Sheet 2012-3047, Reston VA, USA, 2012.
- Livneh, B. and Lettenmaier, D. P.: Multi-criteria parameter estimation for the Unified Land Model, *Hydrol. Earth Syst. Sci.*, 16, 3029–3048, doi:10.5194/hess-16-3029-2012, 2012.
- Livneh, B. and Lettenmaier, D. P.: Regional parameter estimation for the Unified Land Model, *Water Resour. Res.*, 49, 100–114, doi:10.1029/2012WR012220, 2013.
- Livneh, B., Rosenberg, E. A., Lin, C., Nijssen, B., Mishra, V., Andreadis, K. M., Maurer, E. P., and Lettenmaier, D. P.: A Long-Term Hydrologically Based Dataset of Land Surface Fluxes and States for the Conterminous United States: Update and Extensions, *J. Climate*, 26, 9384–9392, doi:10.1175/JCLI-D-12-00508.1, 2013.
- Lohmann, D., Mitchell, K. E., Houser, P. R., Wood, E. F., Schaake, J. C., Robock, A., Cosgrove, B. A., Sheffield, J., Duan, Q., Luo, L., Higgins, R. W., Pinker, R. T., and Tarpley, J. D.: Streamflow and water balance intercomparisons of four land surface models in the North American Land Data Assimilation System project, *J. Geophys. Res.*, 109, D07S91, doi:10.1029/2003JD003517, 2004.
- Martinec, J. and Rango, A.: Merits of statistical criteria for the performance of hydrological models, *Water Resour. B.*, 25, 421–432, 1989.
- Martinez, G. and Gupta, H. V.: Toward improved identification of hydrologic models: A diagnostic evaluation of the “abcd” monthly water balance model for the conterminous United States, *Water Resour. Res.*, 46, W08507, doi:10.1029/2009WR008294, 2010.
- Martinez, G. and Gupta, H. V.: Hydrologic consistency as a basis for assessing complexity of monthly water balance models for the continental United States, *Water Resour. Res.*, 47, W12540, doi:10.1029/2011WR011229, 2011.
- Maurer, E. P., Wood, A. W., Adam, J. C., Lettenmaier, D. P., and Nijssen, B.: A long-term hydrologically-based data set of land surface fluxes and states for the conterminous United States, *J. Climate*, 15, 3237–3251, 2002.
- Menne, M. J., Williams Jr., C. N., and Vose, R. S.: The U.S. Historical Climatology Network monthly temperature data, version 2, *Bull. Am. Meteor. Soc.*, 90, 993–1007, doi:10.1175/2008BAMS2613.1, 2009.
- Menne, M. J., Williams, C. N., and Palecki, M. A.: On the reliability of the U.S. surface temperature record, *J. Geophys. Res.*, 115, D11108, doi:10.1029/2009JD013094, 2010.
- Merz, R. and Bloschl, G.: Regionalization of catchment model parameters, *J. Hydrol.*, 287, 95–123, 2004.
- Mizukami, N., Koren, V., Smith, M., Kingsmill, D., Zhang, Z., Cosgrove, B., and Cui, Z.: The impact of precipitation type discrimination on hydrologic simulation: Rain-snow partitioning derived from HMT-West radar-detected brightband height versus surface temperature data, *J. Hydrometeorol.*, 14, 1139–1158, doi:10.1175/JHM-D-12-035.1, 2013.
- Morton, F. I.: Operational estimates of actual evapotranspiration and their significance to the science and practice of hydrology, *J. Hydrol.*, 66, 1–76, 1983.
- Nash, J. E.: The form of the instantaneous unit hydrograph, International Association of Scientific Hydrology Publication, 45, 114–121, Toronto ON, CA, 1957.
- Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models. Part I: A discussion of principles, *J. Hydrol.*, 10, 282–290, doi:10.1016/0022-1694(70)90255-6, 1970.
- Nathan, R. J. and McMahon, T. A.: The SFB model, Part I – Validation of fixed model parameters, *Civil Eng. Trans., CE32*, 157–161, 1990.
- Nester, T., Kirnbauer, R., Gutknecht, D., and Bloschl, G.: Climate and catchment controls on the performance of regional flood simulations, *J. Hydrol.*, 402, 340–356, 2011.
- Nester, T., Kirnbauer, R., Parajka, J., and Bloschl, G.: Evaluating the snow component of a flood forecasting model, *Hydrol. Res.*, 43, 762–779, 2012.
- Oubeidillah, A. A., Kao, S.-C., Ashfaq, M., Naz, B. S., and Tootle, G.: A large-scale, high-resolution hydrological model parameter data set for climate change impact assessment for the conterminous US, *Hydrol. Earth Syst. Sci.*, 18, 67–84, doi:10.5194/hess-18-67-2014, 2014.
- Oudin, L., Andréassian, V., Mathevet, T., Perrin, C., and Michel, C.: Dynamic averaging of rainfall-runoff model simulations from

- complementary model parameterizations, *Water Resour. Res.*, 42, W07410, doi:10.1029/2005WR004636, 2006.
- Oudin, L., Kay, A. L., Andreassian, V., and Perrin, C.: Are seemingly physically similar catchments truly hydrologically similar?, *Water Resour. Res.*, 46, W11558, doi:10.1029/2009WR008887, 2010.
- Perrin, C., Michel, C., and Andreassian, V.: Does a large number of parameters enhance model performance? Comparative assessment of common catchment model structures on 429 catchments, *J. Hydrol.*, 242, 275–301, doi:210.1016/S0022-1694(1000)00393-00390, 2001.
- Pokhrel, P. and Gupta, H. V.: On the use of spatial regularization strategies to improve calibration of distributed watershed models, *Water Resour. Res.*, 46, W01505, doi:10.1029/2009WR008066, 2010.
- Priestly, C. H. B. and Taylor, R. J.: On the assessment of surface heat flux and evaporation using large-scale parameters, *Mon. Weather Rev.*, 100, 81–82, 1972.
- Samaniego, L., Bardossy, A., and Lumar, R.: Streamflow prediction in ungauged catchments using copula-based dissimilarity measures, *Water Resour. Res.*, 46, W02506, doi:10.1029/2008WR007695, 2010.
- Schaake, J., Cong, S., Duan, Q.: U.S. MOPEX data set. Report UCRL-JRNL-221228, Lawrence Livermore National Laboratory, Livermore CA, USA, available at: <https://e-reports-ext.llnl.gov/pdf/333681.pdf> (last access: 10 September 2014), 2006.
- Schaeefli, B., Hingray, B., Niggli, M., and Musy, A.: A conceptual glacio-hydrological model for high mountainous catchments, *Hydrol. Earth Syst. Sci.*, 9, 95–109, doi:10.5194/hess-9-95-2005, 2005.
- Schaeefli, B. and Gupta, H. V.: Do Nash values have value?, *Hydrol. Process.*, 21, 2075–2080, doi:10.1002/hyp.6825, 2007.
- Schlosser, C. A., Slater, A. G., Robock, A., Pitman, A. J., Vinnikov, K. Y., Henderson-Sellers, A., Speranskaya, N. A., Mitchell, K., and the PILPS 2(d) contributors: Simulations of a boreal grassland hydrology at Valdai, Russia: PILPS phase 2(d), *Mon. Weather Rev.*, 128, 301–321, 2000.
- Sherman, L. K.: Streamflow from rainfall by the unit graph method, *Eng. News Rec.*, 108, 501–505, 1932.
- Shi, X., Wood, A. W., and Lettenmaier, D. P.: How essential is hydrologic model calibration to seasonal streamflow forecasting?, *J. Hydrometeorol.*, 9, 1350–1363, 2008.
- Shuttleworth, W. J. and Calder, I. R.: Has the Priestly-Taylor equation any relevance to forest evaporation?, *J. Appl. Meteorol.*, 18, 639–646, 1979.
- Slack, J. R. and Landwehr, J. M.: Hydro-Climatic Data Network (HCDN): A US Geological Survey streamflow data set for the United States for the study of climate variations, 1874–1988, US Geological Survey, Open-File Report 92-129, Reston VA, USA, 1992.
- Sorooshian, S., Duan, Q., and Gupta, V. K.: Calibration of conceptual rainfall-runoff models using global optimization: application to the Sacramento soil moisture accounting model, *Water Resour. Res.*, 29, 1185–1194, 1993.
- Thornton, P. E. and Running, S. W.: An improved algorithm for estimating incident daily solar radiation from measurements of temperature, humidity and precipitation, *Agr. Forest Meteorol.*, 93, 211–228, 1999.
- Thornton, P. E., Running, S. W., and White, M. A.: Generating surfaces of daily meteorological variables over large regions of complex terrain, *J. Hydrol.*, 190, 214–251, doi:10.1016/S0022-1694(96)03128-9, 1997.
- Thornton, P. E., Hasenauer, H., and White, M. A.: Simultaneous estimation of daily solar radiation and humidity from observed temperature and precipitation: An application over complex terrain in Austria, *Agr. Forest Meteorol.*, 104, 255–271, 2000.
- Thornton, P. E., Thornton, M. M., Mayer, B. W., Wilhelm, N., Wei, Y., and Cook, R. B.: Daymet: Daily surface weather on a 1 km grid for North America, 1980–2012, available at: <http://daymet.ornl.gov/> (last access: 15 July 2013) from Oak Ridge National Laboratory Distributed Active Archive Center, Oak Ridge, Tennessee, USA, 2012.
- Viger, R. J.: Preliminary spatial parameters for PRMS based on the Geospatial Fabric, NLCD2001 and SSURGO, US Geological Survey, doi:10.5066/F7WM1BF7, 2014.
- Viger, R. J. and Bock, A.: GIS Features of the Geospatial Fabric for National Hydrologic Modeling, US Geological Survey, doi:10.5066/F7542KMD, 2014.
- Vrugt, J. A., ter Braak, C. J. F., Clark, M. P., Hyman, J. M., and Robinson, B. A.: Treatment of input uncertainty in hydrologic modeling: Doing hydrology backward with Markov chain Monte Carlo simulation, *Water Resour. Res.*, 44, W00B09, doi:10.1029/2007WR006720, 2008.
- Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Luo, L., Alonge, C., Wei, H., Meng, J., Livneh, B., Lettenmaier, D., Koren, V., Duan, Q., Mo, K., Fan, Y., and Mocko, D.: Continental-scale water and energy flux analysis and validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2): 1. Intercomparison and application of model products, *J. Geophys. Res.*, 117, D03109, doi:10.1029/2011JD016048, 2012.
- Yang, D., Goodison, B. E., Metcalfe, J. R., Golubev, V. S., Bates, R., Pangburn, T., and Hanson, C. L.: Accuracy of NWS 8" standard nonrecording precipitation gauge: Results and application of WMO intercomparison, *J. Atmos. Ocean. Technol.*, 15, 54–68, 1998.
- Yilmaz, K. K., Gupta, H. V., and Wagener, T.: A process-based diagnostic approach to model evaluation: Application to the NWS distributed hydrologic model, *Water Resour. Res.*, 44, W09417, doi:10.1029/2007WR006716, 2008.
- Zhang, Z., Koren, V., Reed, S., Smith, M., Zhang, Y., Moreda, F., and Cosgrove, B.: SAC-SMA a priori parameter differences and their impact on distributed hydrologic model simulations, *J. Hydrol.*, 420–421, 216–227, 2012.
- Zotarelli, L., Dukes, M. D., Romero, C. C., Migliaccio, K. W., and Morgan, K. T.: Step by step calculation of the Penman-Monteith Evapotranspiration (FAO-56 Method). University of Florida Extension, AE459, available at: <http://edis.ifas.ufl.edu> (last access: 1 April 2014), 10 pp., 2009.