

Development of a large-sample watershed-scale hydrometeorological dataset for the contiguous USA: Dataset characteristics and assessment of regional variability in hydrologic model performance

A. J. Newman¹, M. P. Clark¹, K. Sampson¹, A. Wood¹, L. E. Hay², A. Bock², R. J. Viger², D. Blodgett³, L. Brekke⁴, J. R. Arnold⁵, T. Hopson¹ and Q. Duan⁶

¹ National Center for Atmospheric Research, Boulder CO, USA

² United States Geological Survey, Modeling of Watershed Systems, Lakewood CO, USA

³ United States Geological Survey, Center for Integrated Data Analytics, Middleton WI, USA

⁴ U.S. Department of Interior, Bureau of Reclamation, Denver CO, USA

⁵ US Army Corps of Engineers, Institute for Water Resources, Seattle WA, USA

⁶ Beijing Normal University, Beijing, China

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

1 **Abstract**

2 We present a community dataset of daily forcing and hydrologic response data for 671
3 small- to medium-sized basins across the contiguous United States (median basin size of 336
4 km²) that spans a very wide range of hydroclimatic conditions. Areal averaged forcing data for
5 the period 1980-2010 was generated for three basin spatial configurations -- basin mean,
6 Hydrologic Response Units (HRUs) and elevation bands -- by mapping daily, gridded
7 meteorological datasets to the sub-basin (Daymet) and basin polygons (Daymet, Maurer and
8 NLDAS). Daily streamflow data was compiled from the United States Geological Survey
9 National Water Information System. The focus of this paper is to (1) present the dataset for
10 community use; and (2) provide a model performance benchmark using the coupled Snow-17
11 snow model and the Sacramento Soil Moisture Accounting conceptual hydrologic model,
12 calibrated using the Shuffled Complex Evolution global optimization routine. After optimization
13 minimizing daily root mean squared error, 90% of the basins have Nash-Sutcliffe Efficiency
14 scores ≥ 0.55 for the calibration period and 34% ≥ 0.8 . This benchmark provides a reference
15 level of hydrologic model performance for a commonly used model and calibration system, and
16 highlights some regional variations in model performance. For example, basins with a more
17 pronounced seasonal cycle generally have a negative low flow bias, while basins with a smaller
18 seasonal cycle have a positive low flow bias. Finally, we find that data points with extreme error
19 (defined as individual days with a high fraction of total error) are more common in arid basins
20 with limited snow, and, for a given aridity, fewer extreme error days are present as basin snow
21 water equivalent increases.

22

23 1. Introduction

24 With the increasing availability of gridded meteorological datasets, streamflow records
25 and computing resources, large sample hydrology studies have become more common in the last
26 decade or more (i.e. Nathan and McMahon 1990; Perrin et al. 2001; Maurer et al, 2002; Beldring
27 et al. 2003; Merz and Blöschl 2004; Andreassian et al. 2004; Lohmann et al. 2004; Duan et al.
28 2006; Oudin et al. 2006; Oudin et al. 2010; Samaniego et al. 2010; Martinez and Gupta 2010;
29 Nester et al. 2011; Martinez and Gupta 2011; Nester et al. 2012; Livneh and Lettenmaier 2012,
30 2013; Kumar et al. 2013; Oubeidillah et al. 2013). Within the United States there have been
31 several studies to produce large sample hydrometeorological datasets (Maurer et al. 2002;
32 Lohmann et al. 2004; Duan et al. 2006; Thornton et al. 2012; Xia et al. 2012; Livneh et al. 2013).
33 Many of these datasets provide gridded data and may need to be further processed by the end
34 user for their specific hydrologic model configuration. The Model Parameter Estimation Project
35 (MOPEX) dataset does provide basin mean hydrometeorological data and observed streamflow
36 records for 438 basins across the contiguous United States (Schaake et al. 2006) over 30+ years;
37 making it one of the few, high quality, freely available hydrometeorological datasets with
38 immediate applicability to catchment type hydrologic models.

39 Gupta et al. (2014) emphasize that more large-sample hydrologic studies are needed to
40 “balance depth with breadth”; most hydrologic studies have traditionally focused on one or a
41 small number of basins (depth), which hinders the ability to establish general hydrologic
42 concepts applicable across regions (breadth). Gupta et al. (2014) go on to discuss practical
43 considerations for large sample hydrology studies, noting first and foremost that large datasets of
44 quality basin data need to be available and shared in the community. In support of this
45 philosophy, we present a large-sample hydrometeorological dataset and modeling tools to
46 understand regional variability in hydrologic model performance across the contiguous USA
47 (Fig. 1). The development of the basin dataset presented herein takes advantage of high quality
48 freely-available data from various US government agencies and research laboratories. It
49 includes (1) daily forcing data for 671 basins for multiple spatial configurations over the 1980-
50 2010 time period; (2) daily streamflow data; (3) basic metadata (e.g. location, elevation, size, and
51 basin delineation shapefiles) and (4) benchmark model performance which contains the final
52 calibrated model parameter sets, model output timeseries for all basins as well as summary
53 graphics for each basin. This builds on the MOPEX dataset by providing basin mean forcing
54 data for 233 more basins along with two other spatial configurations and the benchmark model
55 performance parameter sets and model output.

56 This dataset and benchmark application is intended for the community to use as a test-bed
57 to facilitate the evaluation of hydrologic modeling and prediction questions. To this end, the
58 benchmark consists of the calibrated, coupled Snow-17 snow model and the Sacramento Soil
59 Moisture Accounting conceptual hydrologic model for all 671 basins using the Shuffled
60 Complex Evolution global optimization routine. Development of a large sample hydrologic
61 dataset such as this will allow for exploration into many important scientific questions. We

62 provide some basic analysis relating to questions such as: 1) What is the model performance
63 across a large sample of basins and how does model performance vary across basin hydro-
64 climatic conditions? 2) How do error characteristics relate to basin calibration performance and
65 hydro-climatic conditions? This basic analysis is intended to highlight some of the important
66 questions that can be answered through large-sample hydrologic studies and provide example
67 results for further exploration.

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

72 **2. Basin Dataset**

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

78 **2.1 Basin Selection**

79 The United States Geological Survey (USGS) developed an updated version of their
80 Geospatial Attributes of Gages for Evaluating Streamflow (GAGES-II) in 2011 (Falcone et al.
81 2010; Falcone 2011). This database contains geospatial information for over 9,000 stream gages
82 maintained by the USGS. As a subset of the GAGES-II database, a portion of the basins with
83 minimal human disturbance (i.e. minimal land use changes or disturbances, minimal human
84 water withdrawals) are noted as “reference” gages. A further sub-setting of the reference gages
85 were made as a follow-on to the Hydro-Climatic Data Network (HCDN) 1988 dataset (Slack and
86 Landwehr 1992). These gages, marked HCDN-2009 (Lins 2012), meet the following criteria: 1)
87 have at least 20 years of complete flow data between 1990-2009 and were active as of 2009, 2)
88 are a GAGES-II reference gage, c) have less than 5 percent imperviousness as measured by the
89 National Land Cover Database (NLCD-2011, Jin et al. 2013), and d) passed a manual survey of
90 human impacts in the basin by local Water Science Center evaluators (Falcone et al. 2010).
91 There are 704 gages in the GAGES-II database that are considered HCDN-2009 across the
92 contiguous United States (CONUS). This study uses that portion of the HCDN-2009 basin set as
93 the starting point since they should best represent natural flow conditions. After initial
94 processing and data availability requirements, 671 basins are used for analysis in this study (Fig.
95 1b). Because these basins have minimal human influence they are almost exclusively smaller,
96 headwater-type basins.

97 **2.2 Forcing and Streamflow Data**

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

116 The Daymet dataset was selected as the primary gridded meteorological dataset to derive
117 forcing data for our streamflow simulations (Thornton et al. 2012). Daymet was chosen because
118 of its high spatial resolution, a necessary requirement to more fully estimate spatial heterogeneity
119 for basins in complex topography. Daymet is a daily, gridded (1x1 km) dataset over the CONUS
120 and southern Canada and is available from 1980 to present. It is derived solely from daily
121 observations of temperature and precipitation. The Daymet variables used here are daily
122 maximum and minimum temperature, precipitation, shortwave downward radiation, day length,
123 and humidity; additionally snow water equivalent is included (not used in this work). These
124 daily values are estimated through the use of an iterative method dependent on local station
125 density and the spatial convolution of a truncated Gaussian filter for station interpolation, and the
126 Mountain Climate simulator (MT-CLIM) to estimate shortwave radiation and humidity
127 (Thornton et al. 1997; Thornton and Running 1999; Thornton et al. 2000). Daymet does not
128 include estimates of potential evapotranspiration (PET), a commonly needed input for conceptual
129 hydrologic models or wind speed and direction. Therefore, PET was estimated using the
130 Priestly-Taylor method (Priestly and Taylor 1972) and is discussed further in section 3. Data
131 quality is an ever-present issue in hydrologic modeling, and while the input data to Daymet are
132 subject to rigorous quality control checks (Durre et al. 2008; 2010) potential errors may remain
133 (Menne et al. 2009; 2010; Oubeidillah et al. 2013). Additionally, the Maurer et al. (2002) and
134 National Land Data Assimilation System (NLDAS) (Xia et al. 2012) 12 km gridded datasets
135 were processed to provide daily forcing data for the basin lumped configuration, resulting in
136 three distinct datasets available for future forcing data impact studies.

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

147 The 671 basins span the entire CONUS and cover a wide range of hydro-climatic
148 conditions. They range from wet, warm basins in the Southeast (SE) US to hot and dry basins
149 in the Southwest (SW) US, to wet cool basins in the Northwest (NW) and dry cold basins in the
150 intermountain (Rocky Mountains in Fig. 1a) western US. Figure 1b displays the basin annual
151 precipitation (colored shading) along with symbols to denote rain and snow dominated basins.
152 In terms of annual mean CDFs, Daymet estimated basin mean temperatures range from -2 °C to
153 23 °C with precipitation amounts of 0.7 to 9.4 mm day⁻¹ (Fig. 2). Annual observed mean runoff
154 ranges from 0.01 to 9.3 mm day⁻¹ with PET estimates ranging from 1.9 to 4.8 mm day⁻¹.
155 Interestingly, this implies that Daymet precipitation itself is not enough to balance the observed
156 runoff in some basins and is consistent with other recent large sample hydrologic studies
157 (Oubeidillah et al. 2013). Seasonal variations in these four variables are large as well, with some
158 basins reaching mean winter time temperatures lower than -10 °C and summer time mean
159 temperatures higher than 25 °C (not shown). The seasonal water balance varies greatly with
160 some basins experiencing much higher precipitation and runoff rates in one season versus
161 another (e.g. spring runoff peaks in mountain snowmelt dominated basins). As expected, PET
162 varies seasonally with a minimum in winter and a maximum in summer.

163 Figure 3 gives cumulative density functions (CDFs) for various physical descriptors of
164 the basin set. The basins range in size from roughly 1 to 25,800 km² with the median basin size
165 being about 335 km² and have mean elevations spanning from nearly sea level (10 m) to high
166 alpine elevations (3570 m) with a median elevation of 462 m. Notably, 75 basins have mean
167 elevations > 2000 m. Corresponding to the large range of elevations in the basin set, the mean
168 slopes vary considerably, spanning over 2 orders of magnitude from near zero to over 200 m km⁻¹.
169 The basin set covers a wide range of basin shapes with aspect ratios ranging from 0.08 to
170 about 11. Finally, there is a large range of forest covers across the basin set which may have
171 implications for hydrologic similarity (Oudin et al. 2010) with 20% of the basins having less than
172 (more than) 14% (98%) forest cover and the median basin having about 80% forest cover
173 (NLCD-2011).

174

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

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

201 **3. Hydrologic modeling benchmark**

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

213 3.1 Models

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

229 3.2 Calibration

230 We employed a split-sample calibration approach following Klemes (1986), assigning the
231 first 15 years of available streamflow data for calibration and the remainder for validation then
232 repeating the calibration using the last 15 years and the initial remaining period for validation;
233 thus, approximately 5500 daily streamflow observations were used for each calibration. To
234 initialize the model calibration moisture states on 1 October, we specified an initial wet SAC-
235 SMA soil moisture state that was allowed to spin down to equilibrium for a given basin by
236 running the first year of the calibration period repeatedly and assumed no initial snow pack. This
237 was done until all SAC-SMA state variables had minimal year over year variations, which is a
238 spin-up approach used by the Project for Intercomparison of Land-Surface Process Schemes (e.g.
239 Schlosser et al. 2000). Determination of optimal calibration sampling and spin-up procedures is
240 an area of active research. Spin-up was performed for every parameter set specified by the
241 optimization algorithm, then the model was integrated for the calibration period and the RMSE
242 for that parameter set was calculated.

243 Objective calibration was done by minimizing the root mean squared error (RMSE) of
244 daily modeled runoff versus observed streamflow using the Shuffled Complex Evolution (SCE)
245 global search algorithm of Duan et al. (1992, 1993). The SCE algorithm uses a combination of
246 probabilistic and deterministic optimization approaches that systematically spans the allowed
247 parameter search space and also includes competitive evolution of the parameter sets (Duan et al.
248 1993). Prior applications to the SAC-SMA model have shown good results (Sorooshian et al.
249 1993; Duan et al. 1994). In the coupled Snow-17 and SAC-SMA modeling system, 35 potential
250 parameters are available for calibration, of which we calibrated 20 parameters having either a

251 priori estimates (Koren et al. 2000) or those found to be most sensitive following Anderson
 252 (2002) (Table 1). The SCE algorithm was run using 10 different random seed starts for the
 253 initial parameter sets for each basin, in part to evaluate the robustness of the optimum in each
 254 case, and the optimized parameter set with the minimum RMSE from the ten different
 255 optimization runs was chosen for evaluation.

256 For Snow-17, six parameters were chosen for optimization (Table 1): The minimum and
 257 maximum melt factors (*MFMIN*, *MFMAX*), the wind adjustment for enhanced energy fluxes to
 258 the snow pack during rain on snow (*UADJ*), the rain/snow partition temperature, which may not
 259 be 0°C (*PXTEMP*), the snow water equivalent for 100% snow covered area (*SI*), and the gauge
 260 catch correction term for snowfall only (*SCF*). These six parameters were chosen because
 261 *MFMIN*, *MFMAX*, *UADJ*, *SCF*, and *SI* are defined as major model parameters by Anderson
 262 (2002). *PXTEMP* was also shown to be important in the Snow-17 model by Mizukami et al.
 263 (2013). The *SCF* is critical in many snow dominated basins as precipitation is generally
 264 underestimated in these types of basins (e.g. Yang et al. 1998) and is certainly underestimated in
 265 some basins in Daymet as shown in Figures 3 and 4.

266 The areal depletion curve (ADC) is considered a major parameter in Snow-17. However, to
 267 avoid expanding the parameter space by the number of ordinates on the curve (typically 10), we
 268 manually specified the ADC according to regional variations in latitude, topographic
 269 characteristics (e.g. plains, hills or mountains) and typical air mass characteristics (e.g. maritime
 270 polar, continental polar) (as suggested in Anderson, 2002). The remaining Snow-17 parameters
 271 were set in the same manner. Following the availability of a priori parameter estimates for SAC-
 272 SMA from a variety of datasets and various calibration studies with SAC-SMA (Koren et al.
 273 2000; Anderson et al. 2006; Pokhrel and Gupta 2010; Zhang et al. 2012) 11 parameters from
 274 SAC-SMA are included for calibration (Table 1). We use an instantaneous unit hydrograph,
 275 represented as a two-parameter Gamma distribution for streamflow routing (Sherman 1932;
 276 Clark 1945; Nash 1957; Dooge 1959), the parameters of which were inferred as part of
 277 calibration. .

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

$$280 \quad PET = \frac{a}{\lambda} \cdot \frac{s \cdot (R_n - G)}{s + \gamma} \quad (1)$$

281 Where λ (MJ kg⁻¹) is the latent heat of vaporization, R_n (MJ m⁻² day⁻¹) is the net radiation
 282 estimated using day of year, all Daymet variables and equations to estimate the various radiation
 283 terms (Allen et al. 1998; Zotarelli et al. 2009), G (MJ m⁻² day⁻¹) is the soil heat flux (assumed to
 284 be zero in this case), s (kPa °C⁻¹) is the slope of the saturation vapor pressure-temperature
 285 relationship, γ (kPa °C⁻¹) is the psychrometric constant and a (unitless) is the P-T coefficient. The
 286 P-T coefficient replaces the aerodynamic term in the Penman-Monteith equation and varies by

287 the typical conditions of the area where the P-T equation is being applied with humid forested
288 basins typically having smaller values and exposed arid basins having larger values
289 (Shuttleworth and Calder 1979; Morton 1983; ASCE 1990). Thus the P-T coefficient was
290 included in the calibration since it should vary from basin to basin.

291 **4. Benchmark results**

292 **4.1 Assessment Objectives and Metrics**

293 Assessment of the models will focus on overall performance across the basin set, regional
294 variations, and error characteristics. Nash-Sutcliffe efficiency (NSE) (Nash and Sutcliffe 1970)
295 and two of the decomposition components of NSE, variance bias (α) and total volume bias (β)
296 (Gupta et al. 2009) are the first metrics examined in two variations. Because NSE scores model
297 performance relative to the observed climatological mean, regions in which the model can track
298 a strong seasonal cycle (large flow autocorrelation) perform relatively better when measured by
299 NSE, and this seasonal enhancement may be imparted when using NSE as the objective function
300 for both the calibration and validation phases (e.g. Schaepli et al. 2007). Additionally, basins
301 with higher streamflow variance and frequent precipitation events have better model
302 performance. Therefore, to give a more standardized picture of model performance across
303 varying hydroclimatologies, the NSE was recomputed using the long-term monthly mean flow
304 instead of mean flow (denoted MNSE hereafter), thus preventing climatological seasonality from
305 inflating the NSE and more accurately ranking basins by the degree to which the model added
306 value over climatology in response to weather events (Garrick et al. 1978; Martinec and Rango
307 1989; Schaepli et al. 2005). MNSE in this context is defined for each day of year (DOY) via a
308 31-day window centered on a given DOY. The long-term flow for that 31-day “month” is
309 computed giving rise to a “monthly” mean flow. Using this type of climatology as the base for
310 an NSE type analysis provides improved standardization in basins with large flow
311 autocorrelations. This definition is similar to the one proposed by Garrick et al. (1978) but with
312 the addition of the 31-day smoother, which is done to provide a smoother reference climatology.

313 Also, several other advanced, more physically based, metrics of model performance are
314 provided. First, three diagnostic signatures based on the flow duration curve (FDC) from Yilmaz
315 et al. (2008) are computed: 1) the top 2% flow bias, 2) the bottom 30% flow bias and 3) the bias
316 of the slope of the middle portion (20-70 percentile) of the FDC. Second, examination of the
317 time series of squared error contribution to the RMSE statistic was performed to highlight events
318 in which the model performs poorly following Clark et al. (2008). This analysis was performed
319 to gauge the representativeness of performance metrics over the model record by using the sorted
320 (highest to lowest) time series of squared error to identify the N number of the largest error days
321 and determine their fractional error contribution to the total. Finally, we extend this analysis to
322 introduce, a simple, normalized general error index for application and comparison across
323 varying modeling and calibration studies. We coin the index, E50, the fraction of calibration

324 points contributing 50% of the error. This captures the number of points determining the
325 majority of the error and thus the optimal parameter set.

326 **4.2 Spatial variability**

327 It is informative to examine spatial patterns of the aforementioned metrics to elucidate
328 factors leading to weak (and strong) model performance. This also allows for identification of
329 outlier basins and characterization of contributing factors (i.e. forcing or streamflow data issues
330 or poor calibration). Poor performing basins are most common along the high plains and desert
331 southwest (Fig 5a, section 3c). When examining MNSE (Fig 5b), basins with high non-seasonal
332 streamflow variance and frequent precipitation events (SE and NW US) have the highest model
333 MNSE, while most of the snowmelt dominated basins see MNSE scores reduced relative to NSE,
334 particularly in the validation phase (Fig. 5c). This indicates that RMSE as an objective function
335 may not be well suited for model calibration in basins with high flow autocorrelation (Kavetski
336 and Fenicia 2011; Evin et al. 2014). This is confirmed by comparing Fig. 5d to Fig. 5c, basins
337 with large flow autocorrelations (one week mean flow for example) generally have lower MNSE
338 scores.

339 Areas with low validation NSE and MNSE scores have generally large biases when
340 looking at FDC metrics as well (Fig. 6). Focusing on the high plains, high flow biases of $\pm 50\%$
341 are common. Extreme negative low flow biases are also present along the high plains and desert
342 SW along with a general model trend to have large negative FDC slope biases, consistent with a
343 poorly calibrated model. For the 72% of basins with validation NSE > 0.55 (basins with yellow-
344 green to dark red colors in Fig. 6a), there is no noticeable spatial pattern across CONUS in
345 regard to high flow periods. However, basins with a more pronounced seasonal cycle (e.g.
346 snowpack dominated watersheds, central West coast) generally have a negative low flow bias,
347 while basins with a smaller seasonal cycle have a positive low flow bias (Fig. 6b).
348 Correspondingly, basins with a pronounced seasonal cycle generally have a near zero or positive
349 slope of the FDC bias, while basins with a smaller seasonal cycle have a negative slope bias (Fig.
350 6c).

351 Past applications similar conceptual snow and hydrologic modeling systems across the
352 CONUS have shown comparable spatial performance patterns. Clark et al. (2008) applied many
353 conceptual models to a subset of the MOPEX basin set and found poor performance in arid
354 regions. Martinez and Gupta (2010), using a monthly water balance model found the best
355 performance generally along the east coast, most of SE CONUS, and along the west coast with
356 scattered good performance in the Rocky Mountains. They found that many basins along the
357 High Plains and north side of the Appalachian Mountains perform poorly. They also note that
358 arid regions have high variability error (variability bias term in KGE).

359 **4.3 Cumulative Performance**

360 Two basic cumulative thresholds for model performance are highlighted here, NSE
361 values of 0.55 and 0.8. An NSE of 0.55 indicates some model skill, and an NSE of 0.8 suggests
362 reasonably good model performance. For the calibration period, 90% (604) of the basins have a
363 NSE greater than 0.55, while 72% (484) of the basins had a validation period NSE > 0.55 (Fig
364 7a). At the NSE > 0.8 level, 34% (225) basin models perform better during calibration and 12%
365 (78) basin models meet that criteria during the validation phase. When using MNSE, 85% and
366 57% (568 and 385) of the basins lie above 0.55 and 17% and 4% (114 and 29) of the basins lie
367 above 0.8 during the calibration and validation phases. The decomposition of the NSE (Gupta et
368 al. 2009) shows that and 90% of basins have a calibration (validation) model-observation flow
369 correlation > 0.75 (0.68) and 30% (12%) of basins have a model-observation flow correlation >
370 0.9 (Fig 7b). However, nearly all basins have too little modeled variance (values less than one)
371 for both the calibration and validation phases (Fig. 7c). The total volume biases are generally
372 small with 94% (79%) of the basins having a calibration (validation) period total flow bias
373 within 10% of observed (Fig. 7d). These are expected results when using RMSE for the
374 objective function (Gupta et al. 2009) and reaffirm that our implementation of SCE is calibrating
375 the model properly.

376 Figure 8 highlights the full split sample approach for calibration following Klemes
377 (1986). It is seen that the calibration and validation statistics give quite similar results regardless
378 of which time period is used for calibration and validation using the Daymet data. This could
379 indicate that both halves of the data are equally challenging to model with this modeling system.
380 We have also included basin calibrations using the first 15 years only for the Maurer et al. (2002)
381 and NLDAS-II (Xia et al. 2012) datasets. It can be seen that the Daymet forcing provides better
382 model performance overall than both Maurer et al. and NLDAS forcing data. This likely relates
383 to the coarser resolution of the Maurer et al. and NLDAS data (12 km) and the somewhat small
384 basin sizes in this basin set. More importantly the inclusion of the Klemes (1986) split-sample
385 approach provides users of this dataset two parameter estimates for each basin using different
386 calibration periods, while the inclusion of three total forcing datasets begins to allow for
387 ensemble type forcing data impact studies across a large basin sample size. In the remaining
388 discussion, only model performance results using the first half of the split sample for calibration
389 are presented.

390 With respect to advanced diagnostics, the model under predicts high flow events in
391 nearly all basins during calibration and slightly less so for the validation period (Fig. 9a). This is
392 an expected result when using RMSE as the objective function because the optimal calibration
393 underestimates flow variability (Gupta et al. 2009). Low flow periods are more evenly over and
394 under predicted (Fig. 9b) for both the calibration and validation time frames with 58% and 61%
395 of basins having more modeled low flow. Finally, the bias in the slope of the FDC is generally
396 under predicted with about 75% of basins having a negative model bias (FDC slope is negative,
397 thus a negative bias indicates the model slope is more positive and that the modeled flow
398 variability is too compressed). The slope of the FDC indicates the variance of daily flows, which

399 primarily relate to the seasonal cycle or the “flashiness” of a basin. Again this indicates model
400 variability is less than observed, at both short and longer time scales. In aggregate, these results
401 agree with Figure 5 and are expected based on the analysis of Gupta et al. (2009). Optimization
402 using RMSE or NSE as the objective function generally results in under prediction of flow
403 variance and near zero total flow bias (Fig 7). This manifests itself in the simulated hydrograph
404 as under predicted high flows, generally over predicted low flows and a more positive slope to
405 the middle portion of the FDC (Fig. 9). It is worth repeating that the goal of this initial
406 application is to provide to community with a benchmark of model performance using well
407 known models, calibration systems and widely used, simple objective functions, thus the use of
408 RMSE.

409 **4.4 Error Characteristics**

410 When examining fractional error statistics for the basin set, 15 basins have single days
411 that contribute at least half the total squared error (potential outlier basins), whereas at the
412 median, the largest error day contributes 8.3% of the total squared error for the median basin (Fig
413 10). The fractional error contribution for the 10, 100 and 1000 largest error days for the median
414 basin are 33%, 70% and 96% of the total squared error respectively. This indicates that for
415 nearly all basins, there are 100 or fewer points that drive the RMSE and therefore optimal model
416 parameters. This type of analysis can be undertaken for any objective function to identify the
417 most influential points and allow for more in-depth examination of forcing data, streamflow
418 records, calibration strategies (i.e. Kavetski et al. 2006; Vrugt et al. 2008; Beven and Westerberg
419 2011; Beven et al. 2011; Kauffeldt et al. 2013), or if different model physics are warranted.

420 The spatial distribution of fractional error contributions show that the issue of model
421 performance being explained by a relatively small set of days is more prevalent in arid regions of
422 CONUS (desert SW US and high plains) as well as basins slightly inland from the east coast of
423 CONUS (Fig 11a-b). The arid basins are generally dry with sporadic high precipitation (and
424 flow) events, while the Appalachian basins are wetter (Fig. 1b) with extreme precipitation events
425 interspersed throughout the record. Basins with significant snowpack tend to have lower error
426 contributions from the largest error days (Fig. 11a-b). The E50 metric highlights mean peak
427 snow water equivalent (SWE) and frequent precipitation basins as well. These regions contain
428 and order of magnitude more days than the high plains and desert SW, giving insight into how
429 representative of the entire streamflow timeseries the optimal model parameter set really is.

430 Additionally, ranking the basins using their fractional error characteristics provides a
431 similar insight. As the aridity index increases, the fractional error contribution increases for
432 basins with little to no mean peak SWE. For basins with significant SWE, the fractional error
433 contribution decreases with increasing aridity (Fig. 12). Alternatively, for a given aridity index
434 the fractional error contribution for N days will decrease with increasing SWE. This dynamic
435 arises because more arid basins with SWE produce a relatively greater proportion of their runoff
436 from snowmelt, without intervening rainfall. This implies that the optimized model produces a

437 more uniform error distribution with less heteroscedasticity in basins with more SWE. Moreover,
438 as the fractional error contribution for the 10 largest error days increases, model NSE generally
439 decreases in the validation phase (Fig. 13). This indicates fractional error metrics are related to
440 overall model performance and that calibration methods to reduce extreme error days should
441 improve model performance. This is not unexpected due to the fact that the residuals from an
442 RMSE type calibration are heteroscedastic. Arid basins typically have few high flow events,
443 which are generally subject to larger errors when minimizing RMSE. Using advanced
444 calibration methodologies that account for heteroscedasticity (Kavetski and Fenicia 2011; Evin et
445 al. 2014) may produce improved calibrations for arid basins in this basin set and provide
446 different insights into model behavior using this type of analysis.

447 **4.5 Limitations and Uncertainties**

448 One interesting example of the usefulness (and a potential limitation) of large sample
449 hydrology stemming from this work lies in the identification of issues with forcing datasets.
450 Figures 3 and 4 show Daymet has too little precipitation in certain regions which is also seen in
451 Oubeidillah et al. (2013). When examining calibrated model performance in the Pacific
452 Northwest, it is seen that several basins along the west coast have low outlier NSE scores.
453 Tracing this unexpected result, we find the Daymet forcing data available for those basins has a
454 negative temperature bias, preventing mid-winter rain and melt episodes in the modeling system,
455 identifying scope to improve the Daymet forcing. Moreover, winter periods of observed
456 precipitation and streamflow rises coincide with subzero T_{\max} in the Daymet dataset, also
457 suggesting areas to improve the Daymet forcing. The large sample of basins in this region (91)
458 allowed for identification of the outlier basins and the underlying causes.

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

469 **5. Summary and Discussion**

470 Most hydrologic studies focus in detail on a small number of watersheds, providing
471 comprehensive but highly local insights, and may be limited in their ability to inform general
472 hydrologic concepts applicable across regions (Gupta et al. 2014). To facilitate large-sample
473 hydrologic studies, large-sample basin datasets and corresponding benchmarks of model

474 performance using standard methodology across all basins need to be freely available to the
475 community. To that end, we have compiled a community dataset of daily forcing and
476 streamflow data for 671 basins and provide a benchmark of performance using a widely used
477 conceptual a hydrologic modeling and calibration scheme over a wide range of conditions.

478 Overall, application of the basin set to assessing an objectively calibrated conceptual
479 hydrologic model representation of the 671 watersheds yielded calibration Nash-Sutcliffe
480 Efficiency (NSE) scores of > 0.55 (0.8) for 90 (34) percent of the basins. Performance of the
481 models varied regionally, and the main factors influencing this variation were found to be aridity
482 and precipitation intermittency, contribution of snowmelt, and runoff seasonality. Analysis of
483 the cumulative fractional error contributions from the largest error days showed that the presence
484 of significant snow water equivalent (SWE) offset the negative impact of increasing aridity on
485 simulation performance. This study has identified potential outlier basins for this modeling
486 system and has provided insights into potential forcing data limitations. Although this modeling
487 application utilized a conceptual hydrologic model with a single-objective calibration strategy,
488 the findings provide a baseline for assessing more complex strategies in each area, including
489 multi-objective calibration of more highly distributed hydrologic models (e.g., in Shi et al 2008).
490 The unusually broad variation of hydroclimatologies represented by the dataset, which contains
491 forcing and streamflow data obtained by consistent methodology and retains outlier basins,
492 makes it a notable resource for these and other future large-sample watershed-scale hydrologic
493 analysis efforts.

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

497 **6. Acknowledgements**

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

504 **7. References**

505 Allen, R. G., L. S. Pereira, D. Raes, and M. Smith (1988). Crop evapotranspiration: guidelines
506 for computing crop water requirements. Food and Agriculture Organization of the United
507 Nations, Rome, 15 pp.

508 Anderson, E. A. (1973). National Weather Service River Forecast System – Snow accumulation
509 and ablation model. NOAA Technical Memorandum, NWS, HYDRO-17, US Department of
510 Commerce, Silver Spring, MD, 217 pp.

511 Anderson, E. A. (2002). Calibration of conceptual hydrologic models for use in river
512 forecasting. NOAA Technical Report, NWS 45, Hydrology Laboratory, Silver Spring, MD.

513 Anderson, R. M., V. I. Koren, S. M. Reed (2006). Using SSURGO data to improve Sacramento
514 Model a priori parameter estimates. *Journal of Hydrology*, 320, 103-116.

515 Andreassian, V., A. Oddos, C. Michel, F. Anctil, C. Perrin and C. Loumange (2004). Impact of
516 spatial aggregation of inputs and parameters on the efficiency of rainfall-runoff models: A
517 theoretical study using chimera watersheds. *Water Resources Research*, 40(5): W05209, doi:
518 10.1029/2003WR002854.

519 Beldring, S., Engeland, K., Roald, L. A., Saelthun, N. R., and Vokso, A (2003). Estimation of
520 Parameters in a distributed precipitation-runoff model for Norway, *Hydrology and
521 Earth System Sciences*, 7, 304G316.

522

523 Beven, K. and I. Westerberg (2011). On red herrings and real herrings: disinformation and
524 information in hydrological inference. *Hydrologic Processes*, 25, 1676-1680.

525 Beven, K, P. J. Smith, and A. Wood (2011). On the colour and spin of epistemic error (and what
526 we might do about it). *Hydrology and Earth System Sciences*, 15, 3123-3133,
527 doi:10.5194/hess-15-3123.

528 Blodgett, D. L., N. L. Booth, T. C. Kunicki, J. L. Walker, and R. J. Viger (2011). Description
529 and testing of the geo data portal: A data integration framework and web processing services
530 for environmental science collaboration. US Geological Survey, Open-File Report 2011-
531 1157, 9 pp., Middleton WI, USA.

532 Burnash, R.J.C., R. L. Ferral, R. A. McGuire (1973). A generalized streamflow simulation
533 system conceptual modeling for digital computers, U. S. Department of Commerce National
534 Weather Service and State of California Department of Water Resources.

535 Burnash, R. J. C. (1995). The NWS River Forecast System – Catchment model. In *Computer
536 Models of Watershed Hydrology*, edited by V. P. Singh, pp. 311-366, Water Resources
537 Publications, Highlands Ranch, Colo.

538 Clark, C. O. (1945). Storage and the unit hydrograph. *Proc. Am. Soc. Civ. Eng.*, vol. 9, pp 1333-
539 1360.

540 Clark, M. P., A. G. Slater, D. E. Rupp, R. A. Woods, J. A. Vrugt, H. V. Gupta, T. Wagener, and
541 L. E. Hay (2008). Framework for Understanding Structural Errors (FUSE): A modular

542 framework to diagnose differences between hydrologic models. *Water Resources Research*,
543 44, W00B02, doi:10.1029/2007WR006735.

544 Clark, M. P., D. Kavetski, and F. Fenicia (2011). Pursuing the method of multiple working
545 hypotheses for hydrological modeling. *Water Resources Research*, 47, W09301,
546 doi:10.1029/2010WR009827.

547 Dooge, J. C. I. (1959). A general theory of the unit hydrograph. *Journal of Geophysical*
548 *Research*, 64(2): 241-256.

549 Duan, Q., S. Sorooshian, and V. K. Gupta (1992). Effective and efficient global optimization for
550 conceptual rainfall-runoff models. *Water Resources Research*, 28(4): 1015-1031.

551 Duan, Q., V. K. Gupta, and S. Sorooshian (1993). A shuffled complex evolution approach for
552 effective and efficient optimization. *Journal of Optimization Theory and Applications*, 76(3):
553 501-521.

554 Duan, Q., S. Sorooshian, V. K. Gupta (1994). Optimal use of the SCE-UA global optimization
555 method for calibrating watershed models. *Journal of Hydrology*, 158, 265-284.

556 Duan, Q., J. Schaake, V. Andreassian, S. Franks, G. Goteti, H. V. Gupta, Y. M. Gusev, F.
557 Habets, A. Hall, L. Hay, T. Houge, M. Huang, G. Leavesley, X. Liang, O. N. Nasonova, J.
558 Noilhan, L. Oudin, S. Sorooshian, T. Wagener, and E. F. Wood (2006). Model Parameter
559 Estimation Experiment (MOPEX): An overview of science strategy and major results from
560 the second and third workshops. *J. Hydrology*, **320**, 3-17.

561 Evin, G., M. Thyer, D. Kavetski, D. McInerney, and G. Kuczera (2014). Comparison of joint
562 versus postprocessor approaches for hydrological uncertainty estimation accounting for error
563 autocorrelation and heteroscedasticity. *Water Resources Research*, 50, 2350-2375, doi:
564 10.1002/2013WR014185.

565 Falcone, J. A., D. M. Carlisle, D. M. Wolock, and M. R. Meador (2010). GAGES: A stream gage
566 database for evaluating natural and altered flow conditions in the conterminous United
567 States. *Ecology*, 91(2), p. 621. A data paper in Ecological Archives E091-045-D1, available
568 at <http://esapubs.org/Archive/ecol/E091/045/metadata.htm>, (last access: 05 April 2014),
569 2010.

570 Falcone, J. A. (2011). GAGES-II: Geospatial Attributes of Gages for Evaluating Streamflow.
571 Digital spatial data set 2011. Available at:
572 http://water.usgs.gov/GIS/metadata/usgswrd/XML/gagesII_Sept2011.xml (last access: 10
573 Oct 2013), 2011.

574 Garrick, M., C. Cunnane, and J. E. Nash (1978). A criterion of efficiency for rainfall-runoff
575 models. *J. Hydrology*, **36**(3-4), 375-381.

576 Gupta, H. V., H. Kling, K. K. Yilmaz, and G. F. Martinez-Barquero (2009). Decomposition of
577 the mean squared error and NSE performance criteria: Implications for improving
578 hydrological modeling. *Journal of Hydrology*, 377, 80-91,
579 doi:10.1016/j.jhydrol.2009.08.003.

580 Gupta, H. V., C. Perrin, R. Kumar, G. Bloschl, M. Clark, A. Montanari, and V. Andreassian
581 (2014). Large-sample hydrology: A need to balance depth with breadth. *Hydrology and
582 Earth System Sciences-Earth System Discussions*.

583 Jensen, M. E., R. D. Burman, and R. G. Allen (1990). Evapotranspiration and irrigation water
584 requirements. American Society of Civil Engineers, ASCE Manual and Reports on
585 Engineering Practice, 332 p., New York, NY.

586 Jin, S., Yang, L., Danielson, P., Homer, C., Fry, J., and Xian, G. 2013. [A comprehensive change
587 detection method for updating the National Land Cover Database to circa 2011](#). *Remote
588 Sensing of Environment*, 132: 159 – 175.

589 Kauffeldt, A., S. Halldin, A. Rodhe, C.-Y. Xu, and I. K. Westerberg (2013). Disinformative data
590 in large-scale hydrological modeling. *Hydrology and Earth System Sciences*, 17, 2845-2857,
591 doi:10.5194/hess-17-2845-2013.

592 Kavetski, D., and F. Fenicia (2011). Elements of a flexible approach for conceptual hydrological
593 modeling: 2. Application and experimental insights. *Water Resources Research*, 47,
594 W11511,doi:10.1029/2011WR010748.

595 Kavetski, D., G. Kuczera, and S. W. Franks (2006). Bayesian analysis of input uncertainty in
596 hydrological modeling: 2. Application. *Water Resources Research*, 42, W03407,
597 doi:10.1029/2005WR004376.

598 Klemes, V. (1986). Operational testing of hydrological simulation models. *Hydrol. Sci. J.*, **31**(1),
599 13-24.

600 Koren, V. I., M. Smith, D. Wang, and Z. Zhang (2000). Use of soil property data in the
601 derivation of conceptual rainfall-runoff model parameters. American Meteorological Society
602 15th Conference on Hydrology, Long Beach, CA, pp. 103-106.

603 Kumar, R., L. Samaniego, and S. Attinger (2013). Implications of distributed hydrologic model
604 parameterization on water fluxes at multiple scales and locations. *Water Resources Research*,
605 49, 360-379, doi:10.1029/2012WR012195.

606 Lins, H. F. (2012). USGS Hydro-Climatic Data Network 2009 (HCDN-2009), U. S. Geological
607 Survey, Fact Sheet 2012-3047, Reston VA, USA.

608 Livneh, B. and D. P. Lettenmaier (2012). Multi-criteria parameter estimation for the Unified
609 Land Model, *Hydrology and Earth System Sciences*, 16, 3029-3048, doi:10.5194/hess-16-
610 3029-2012.

611 Livneh, B. and D. P. Lettenmaier (2013). Regional parameter estimation for the Unified Land
612 Model. *Water Resources Research*, 49, 100-114, 10.1029/2012WR012220.

613 Livneh, B., E. A. Rosenberg, C. Lin, B. Nijssen, V. Mishra, K. M. Andreadis, E. P. Maurer, D. P.
614 Lettenmaier (2013): A Long-Term Hydrologically Based Dataset of Land Surface Fluxes and
615 States for the Conterminous United States: Update and Extensions. *J. Climate*, **26**, 9384–
616 9392. doi: <http://dx.doi.org/10.1175/JCLI-D-12-00508.1>
617

618 Lohmann, D., Mitchell, K.E., Houser, P.R., Wood, E.F., Schaake, J.C., Robock, A., Cosgrove,
619 B.A., Sheffield, J., Duan, Q., Luo, L., Higgins, R. W., Pinker, R.T., and Tarpley, J.D. (2004).
620 Streamflow and water balance intercomparisons of four land surface models in the North
621 American Land Data Assimilation System project. *Journal Geophysical Research*, 109,
622 D07S91, doi:10.1029/2003ID003517.
623

624 Martinec, J. and A. Rango (1989). Merits of statistical criteria for the performance of
625 hydrological models. *Water Resources Bulletin*, **25**(2), 421-432.
626

627 Martinez, G. and H. V. Gupta (2010). Toward improved identification of hydrologic models: A
628 diagnostic evaluation of the “abcd” monthly water balance model for the conterminous
629 United States. *Water Resour. Res.*, **46**, W08507, doi:10.1029/2009WR008294.
630

631 Martinez, G. and H. V. Gupta (2011). Hydrologic consistency as a basis for assessing
632 complexity of monthly water balance models for the continental United States. *Water*
633 *Resour. Res.*, **47**, W12540, doi:10.1029/2011WR011229.
634

635 Maurer, E. P., A. W. Wood, J. C. Adam, D. P. Lettenmaier, and B. Nijssen (2002). A long-term
636 hydrologically-based data set of land surface fluxes and states for the conterminous United
637 States. *Journal of Climate*, 15(22), 3237-3251.

638 Merz, R. and G. Bloschl (2004). Regionalization of catchment model parameters. *Journal of*
639 *Hydrology*, 287(1-4): 95-123.

640 Mizukami, N., V. Koren, M. Smith, D. Kingsmill, Z. Zhang, B. Cosgrove, and Z. Cui (2013).
641 The impact of precipitation type discrimination on hydrologic simulation: Rain-snow
642 partitioning derived from HMT-West radar-detected brightband height versus surface
643 temperature data. *Journal of Hydrometeorology*, 14, 1139-1158, doi:10.1175/JHM-D-12-
644 035.1.

645 Morton, F. I. (1983). Operational estimates of actual evapotranspiration and their significance to
646 the science and practice of hydrology. *Journal of Hydrology*, 66: 1-76.

- 647 Nash, J. E. (1957). The form of the instantaneous unit hydrograph. International Association of
648 Scientific Hydrology Publication, 45(3), 114-121, Toronto ON, CA.
- 649 Nash, J. E., and J. V. Sutcliffe (1970). River flow forecasting through conceptual models. Part I:
650 A discussion of principles. *Journal of Hydrology*, 10(3), 282-290, doi:10.1016/0022-
651 1694(70)90255-6.
- 652 Nathan, R. J., and T. A. McMahon (1990). The SFB model, Part I – Validation of fixed model
653 parameters. *Civil Engineering Transactions*, CE32, 157-161.
- 654 Nester, T., R. Kirnbauer, D. Gutknecht and G. Bloschl (2011). Climate and catchment controls
655 on the performance of regional flood simulations. *Journal of Hydrology*, 402, 340-356.
- 656 Nester, T., R. Kirnbauer, J. Parajka and G. Bloschl (2012). Evaluating the snow component of a
657 flood forecasting model. *Hydrology Research*, 43(6), 762-779.
- 658 Oubeidillah, A. A., S.-C. Kao, M. Ashfaq, B. Naz, and G. Tootle (2013). A large-scale, high-
659 resolution hydrological model parameter dataset for climate change impact assessment for
660 the conterminous United States. *Hydrology and Earth System Sciences*, 10, 9575-9613,
661 doi:10.5194/hessd-10-9575-2013.
- 662 Oudin, L., V. Andreassian, T. Mathevet, C. Perrin, and C. Michel (2006). Dynamic averaging of
663 rainfall-runoff model simulations from complementary model parameterizations. *Water
664 Resources Research*, 42(7): W07410, doi:10.1029/2005WR004636.
- 665 Oudin, L., A. L. Kay, V. Andreassian, and C. Perrin (2010). Are seemingly physically similar
666 catchments truly hydrologically similar? *Water Resources Research*, 46, W11558,
667 doi:10.1029/2009WR008887.
- 668 Perrin, C., C. Michel, and V. Andreassian (2001). Does a large number of parameters enhance
669 model performance? Comparative assessment of common catchment model structures on 429
670 catchments. *J. Hydrology*, 242(3-4), 275-301, doi:10.1016/S0022-1694(1000)00393-00390.
- 671 Pokhrel, P. and H. V. Gupta (2010). On the use of spatial regularization strategies to improve
672 calibration of distributed watershed models. *Water Resources Research*, 46, W01505,
673 doi:10.1029/2009WR008066.
- 674 Priestly, C. H. B. and R. J. Taylor (1972). On the assessment of surface heat flux and evaporation
675 using large-scale parameters. *Monthly Weather Review*, 100:81-82.
- 676 Samaniego, L., A. Bardossy, and R. Lumar (2010). Streamflow prediction in ungauged
677 catchments using copula-based dissimilarity measures. *Water Resources Research*, 46,
678 W02506, doi:10.1029/2008WR007695.

679 Schaake, J., S. Cong, Q. Duan (2006). U.S. MOPEX data set. Report UCRL-JRNL-221228,
680 Lawrence Livermore National Laboratory, Livermore CA, USA. Available online at:
681 <https://e-reports-ext.llnl.gov/pdf/333681.pdf> (last access: 10 September 2014).

682 Schaefli, B., B. Hingray, M. Niggli, and A. Musy (2005). A conceptual glacio-hydrological
683 model for high mountainous catchments. *Hydrology and Earth System Sciences*, **9**, 157-171.

684 Schaefli, B., and H. V. Gupta (2007). Do Nash values have value? *Hydrological Processes*, **21**,
685 2075-2080, doi:10.1002/hyp.6825.

686 Schlosser, C.A., Slater, A.G., Robock, A., Pitman, A.J., Vinnikov, K.Y., Henderson-Sellers, A.,
687 Speranskaya, N.A., Mitchell, K., and the PILPS 2(d) contributors (2000). Simulations of a
688 boreal grassland hydrology at Valdai, Russia: PILPS phase 2(d). *Monthly Weather Review*,
689 **128**, 301-321.

690 Sherman, L. K. (1932). Streamflow from rainfall by the unit graph method. *Eng. News Rec.*,
691 **108**, 501-505.

692 Shi, X, A. W. Wood, and D. P. Letenmaier (2008). How essential is hydrologic model
693 calibration to seasonal streamflow forecasting? *Journal of Hydrometeorology*, **9**, 1350-1363.

694 Shuttleworth, W. J., and I. R. Calder (1979). Has the Priestly-Taylor equation any relevance to
695 forest evaporation? *Journal of Applied Meteorology*, **18**, 639-646.

696 Slack, J. R., and J. M. Landwehr (1992). Hydro-Climatic Data Network (HCDN): A U. S.
697 Geological Survey streamflow data set for the United States for the study of climate
698 variations, 1874-1988. U. S. Geological Survey, Open-File Report 92-129, Reston VA, USA.

699 Sorooshian, S., Q. Duan, and V. K. Gupta (1993). Calibration of conceptual rainfall-runoff
700 models using global optimization: application to the Sacramento soil moisture accounting
701 model. *Water Resources Research*, **29**(4): 1185-1194.

702 Thornton, P. E., S. W. Running, and M. A. White (1997). Generating surfaces of daily
703 meteorological variables over large regions of complex terrain. *Journal of Hydrology*, **190**:
704 214-251. [http://dx.doi.org/10.1016/S0022-1694\(96\)03128-9](http://dx.doi.org/10.1016/S0022-1694(96)03128-9).

705 Thornton, P. E., and S. W. Running (1999). An improved algorithm for estimating incident daily
706 solar radiation from measurements of temperature, humidity and precipitation. *Agriculture
707 and Forest Meteorology*, **93**:211-228.

708 Thornton, P. E., H. Hasenauer, and M. A. White (2000). Simultaneous estimation of daily solar
709 radiation and humidity from observed temperature and precipitation: An application over
710 complex terrain in Austria. *Agricultural and Forest Meteorology*, **104**:255-271.

711 Thornton, P. E., M. M. Thornton, B. W. Mayer, N. Wilhelmi, Y. Wei, and R. B. Cook (2012).
712 Daymet: Daily surface weather on a 1 km grid for North America, 1980-2012. Acquired
713 online (<http://daymet.ornl.gov/>) on 15/07/2013 from Oak Ridge National Laboratory
714 Distributed Active Archive Center, Oak Ridge, Tennessee, USA.
715 http://doi:10.3334/ORNLDAAAC/Daymet_V2

716 Viger, R. J, and A. Bock (2014). GIS Features of the Geospatial Fabric for National Hydrologic
717 Modeling, US Geological Survey, <http://dx.doi.org/doi:10.5066/F7542KMD>

718 Viger, R. J. (2014). Preliminary spatial parameters for PRMS based on the Geospatial Fabric,
719 NLCD2001 and SSURGO, US Geological Survey,
720 <http://dx.doi.org/doi:10.5066/F7WM1BF7>

721 Vrugt, J. A., C. J. F. ter Braak, M. P. Clark, J. M. Hyman, and B. A. Robinson (2008). Treatment
722 of input uncertainty in hydrologic modeling: Doing hydrology backward with Markov chain
723 Monte Carlo simulation. *Water Resources Research*, 44, W00B09,
724 doi:10.1029/2007WR006720.

725 Xia, Y., K. Mitchell, M. Ek, J. Sheffield, B. Cosgrove, E. Wood, L. Luo, C. Alonge, H. Wei, J.
726 Meng, B. Livneh, D. Lettenmaier, V. Koren, Q. Duan, K. Mo, Y. Fan and D. Mocko (2012).
727 Continental-scale water and energy flux analysis and validation for the North American Land
728 Data Assimilation System project phase 2 (NLDAS-2): 1. Intercomparison and application of
729 model products. *J. Geophys. Res.*, **117**, D03109, doi:10.1029/2001JD016048.

730 Ylimaz, K. K., H. V. Gupta, and T. Wagener (2008). A process-based diagnostic approach to
731 model evaluation: Application to the NWS distributed hydrologic model. *Water Resources*
732 *Research*, 44, W09417, doi:10.1029/2007WR006716.

733 Zhang, Z., V. Koren, S. Reed, M. Smith, Y. Zhang, F. Moreda, and B. Cosgrove (2012). SAC-
734 SMA a priori parameter differences and their impact on distributed hydrologic model
735 simulations. *Journal of Hydrology*, 420-421 (2012), 216-227.

736 Zotarelli, L., M. D. Dukes, C. C. Romero, K. W. Migliaccio, and K. T. Morgan (2009). Step by
737 step calculation of the Penman-Monteith Evapotranspiration (FAO-56 Method). University
738 of Florida Extension, AE459, <http://edis.ifas.ufl.edu> (last access: 01 April 2014), 10 pp.

739

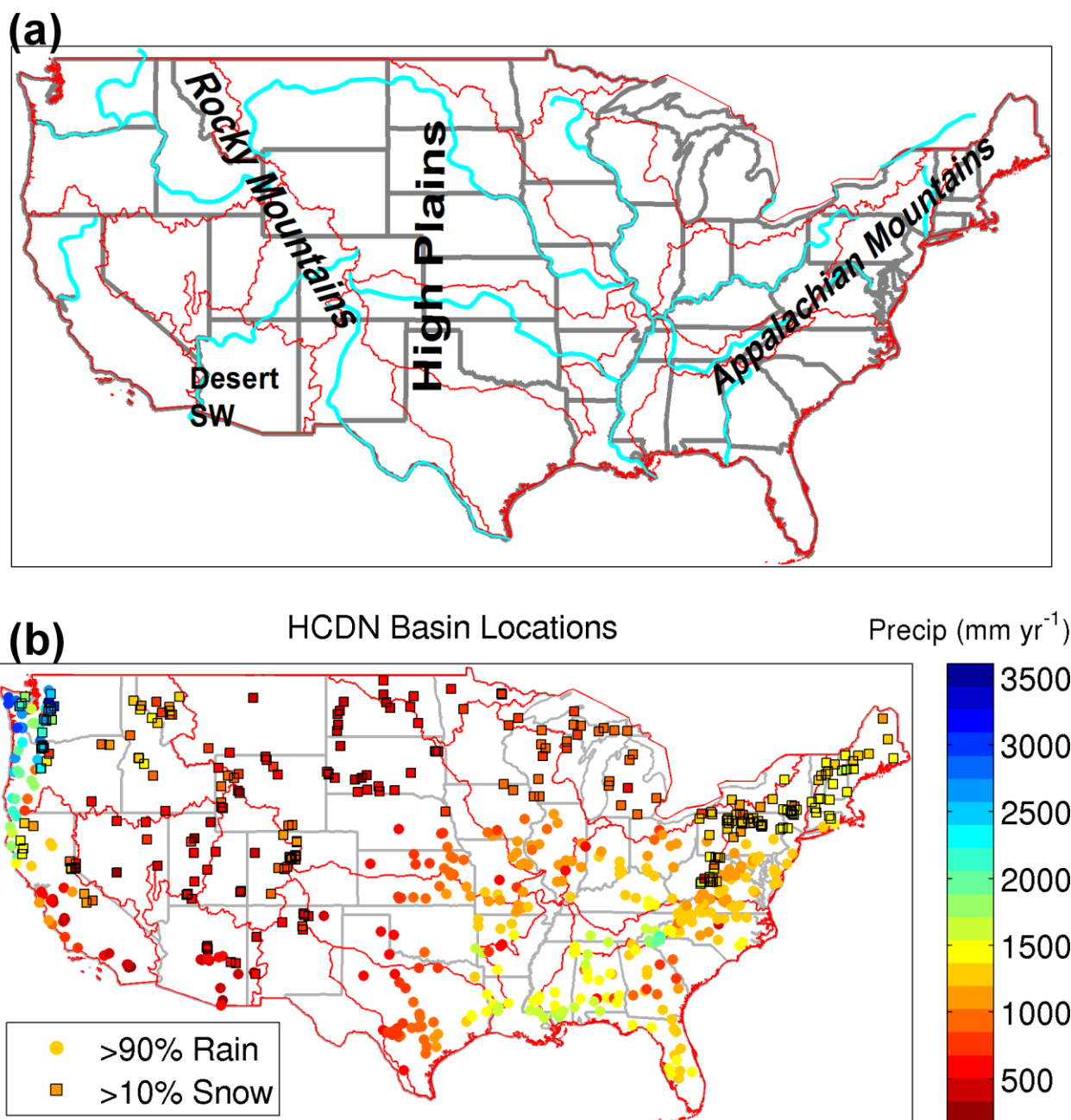
740 Table 1. Table describing all parameters calibrated and their bounds for calibration.

Parameter	Description	Units	Calibration Range
<i>Snow-17</i>			
MFMAX	Maximum melt factor	mm °C ⁻¹ 6-hr ⁻¹	0.8 – 3.0
MFMIN	Minimum melt factor	mm °C ⁻¹ 6-hr ⁻¹	0.01– 0.79
UADJ	Wind adjustment for enhanced flux during rain on snow	km 6-hr ⁻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	°C	-1.0 – 3.0
<i>SAC-SMA</i>			
UZWWM	Upper zone tension water maximum storage	mm	1.0 – 800.0
UZFWM	Upper zone free water maximum storage	mm	1.0 – 800.0
LZWWM	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	day ⁻¹	0.1 – 0.7
LZPK	Lower zone primary free water depletion rate	day ⁻¹	0.00001 – 0.025
LZSK	Lower zone secondary free water depletion rate	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
<i>Others</i>			
USHAPE	Shape of unit hydrograph	-	1.0 – 5.0
USCALE	Scale of unit hydrograph	-	0.001 – 150.0
PT	Priestly-Taylor coefficient	-	1.26 – 1.74

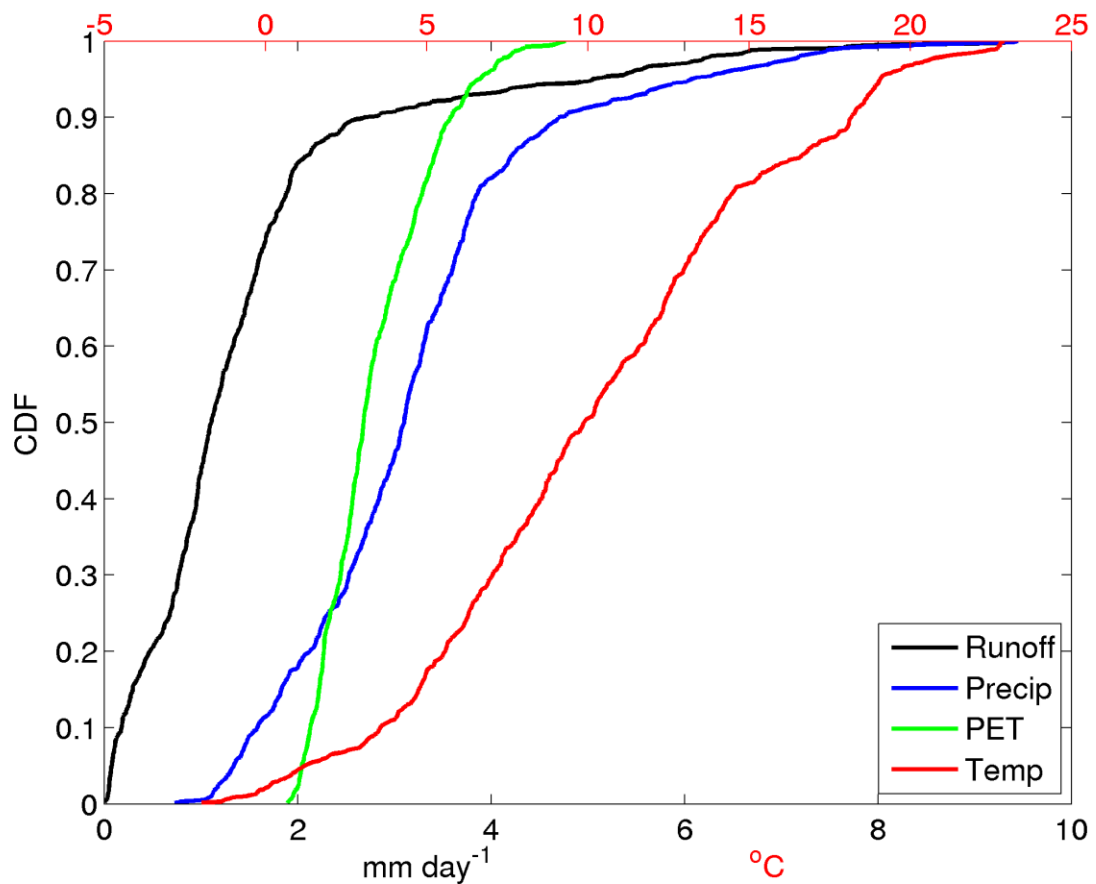
741

742

743



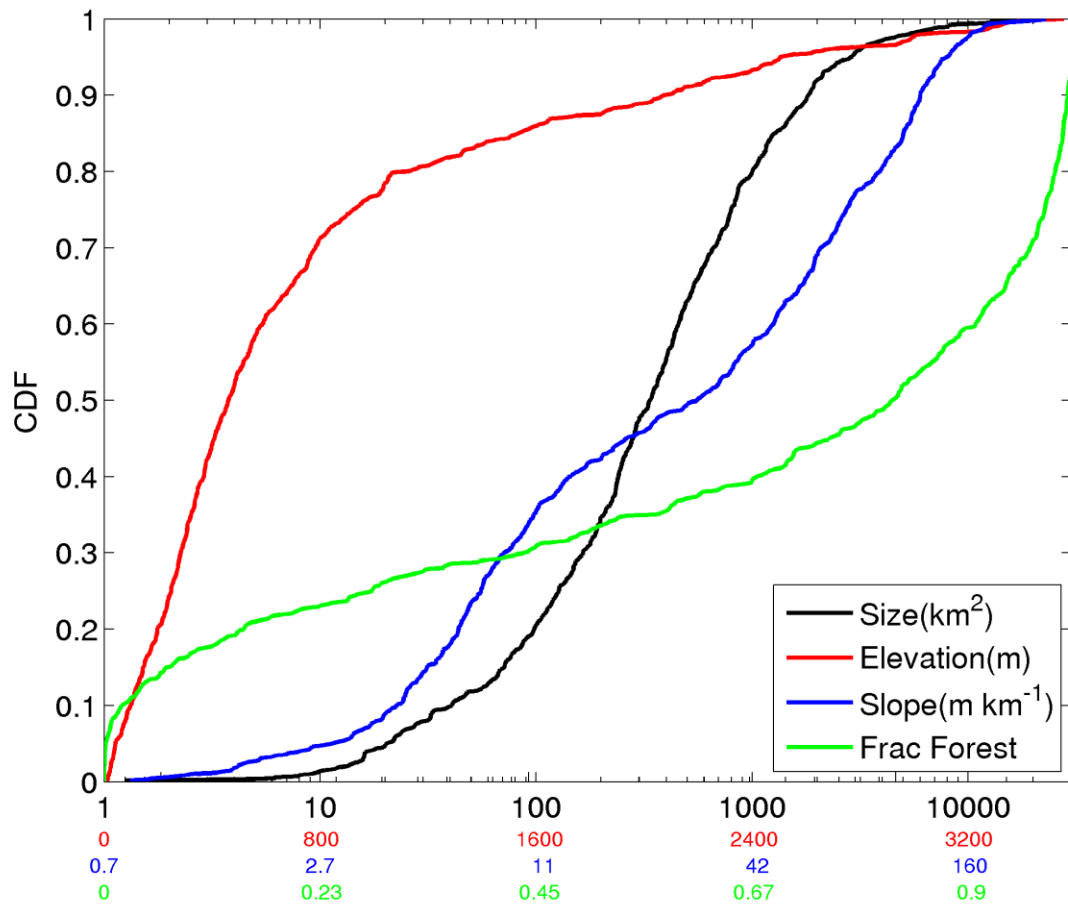
745
 746 Figure 1. (a) Contiguous United States (US) with states (gray), rivers (blue) and major
 747 hydrologic regions (red). Text indicates major geographic regions discussed in text. (b) Location
 748 of the 671 HCDN-2009 basins across the contiguous US used in the basin dataset with
 749 precipitation shaded. Circles denote basins with > 90% of their precipitation falling as rain,
 750 squares with black outlines denote basins with > 10% of their precipitation falling as snow as
 751 determined by using a 0°C daily mean Daymet temperature threshold. State outlines are in thin
 752 gray and hydrologic regions in thin red.



753

754 Figure 2. Annual cumulative density functions (CDFs) of runoff (mm day⁻¹) (black, bottom X-
 755 axis), precipitation (mm day⁻¹) (blue, bottom X-axis), potential evapotranspiration (mm day⁻¹)
 756 (green, bottom X-axis), and temperature (°C) (red, top X-axis).

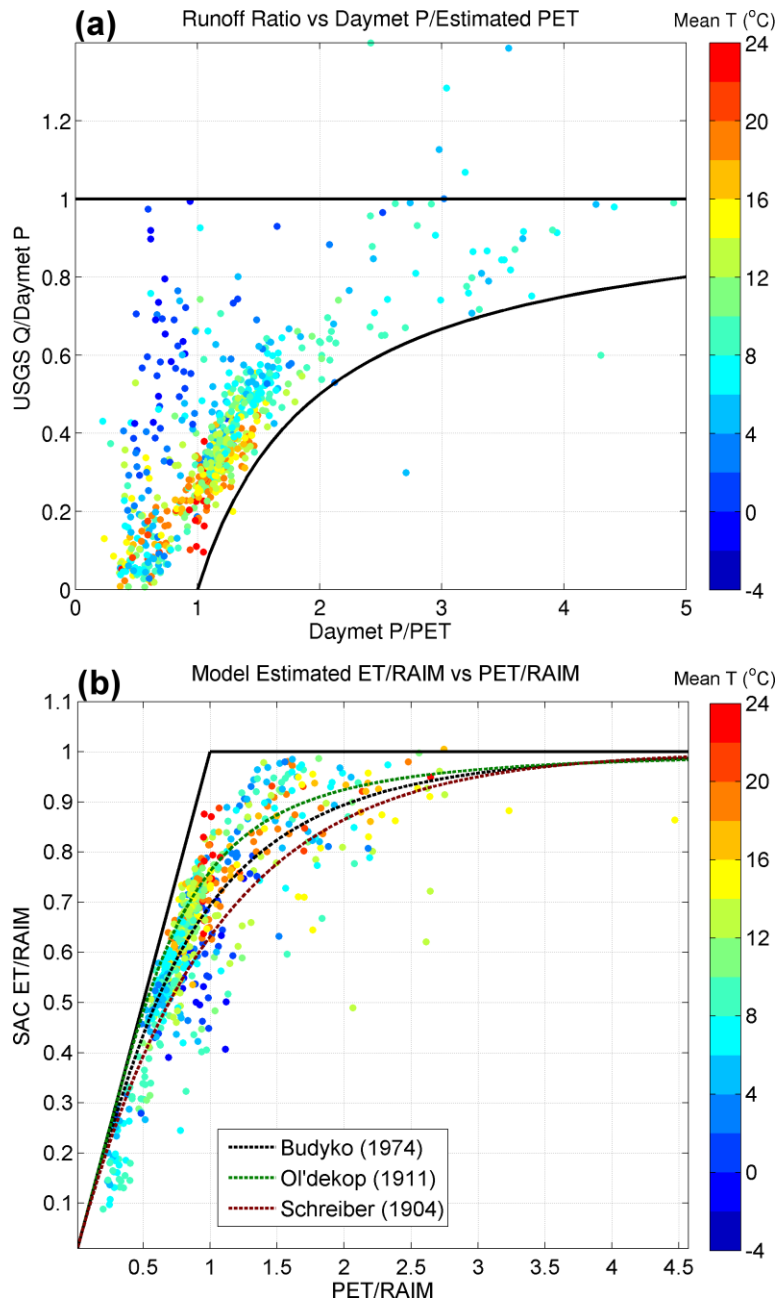
757



758

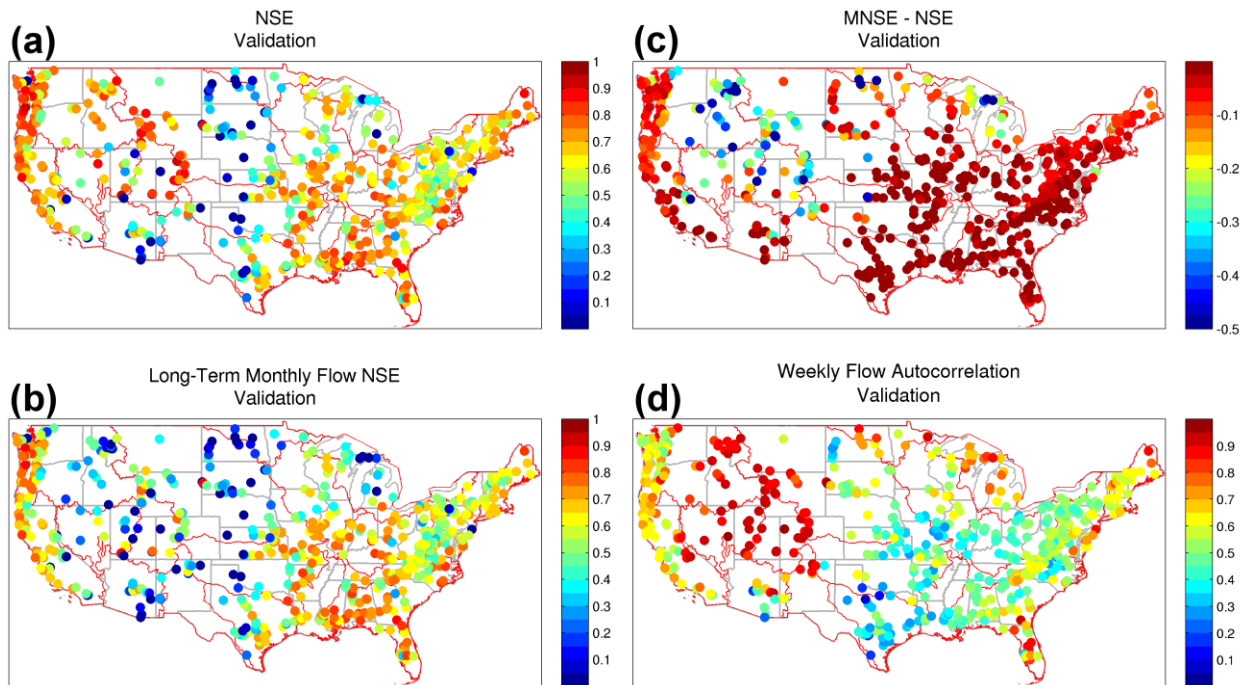
759 Figure 3. Cumulative density functions of basin size (km²) (black), basin mean elevation (m)
 760 (red), mean slope (m km⁻¹) (blue), and fractional forest cover (green) for the basin set.

761



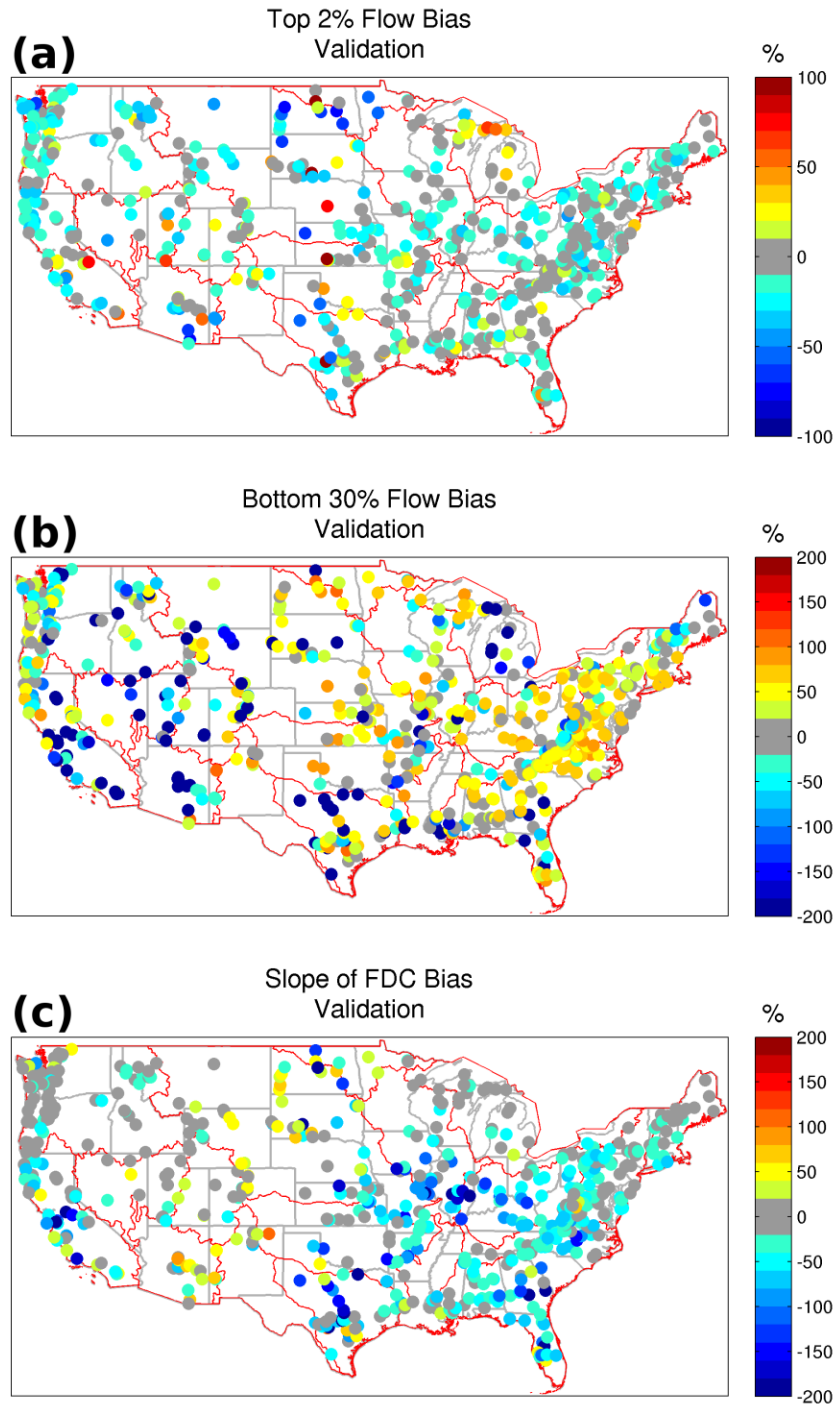
762

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



768
 769
 770
 771
 772
 773
 774

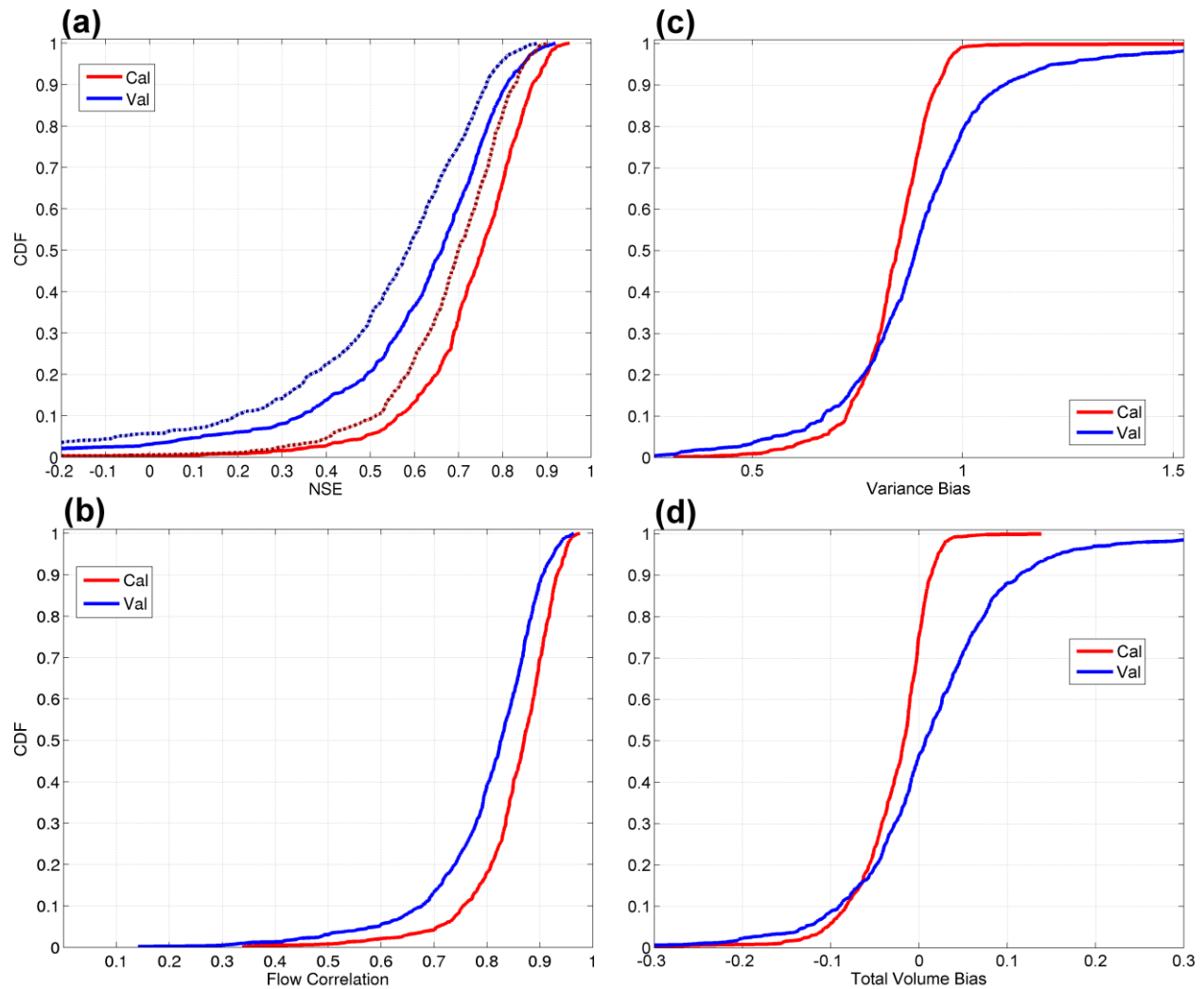
Figure 5. (a) Spatial distribution of Nash-Sutcliffe efficiency (NSE), (b) Nash-Sutcliffe efficiency using long-term monthly mean flows (MNSE) rather than the long-term mean flow, (c) MNSE – NSE for the validation period, (d) weekly flow autocorrelation.



775

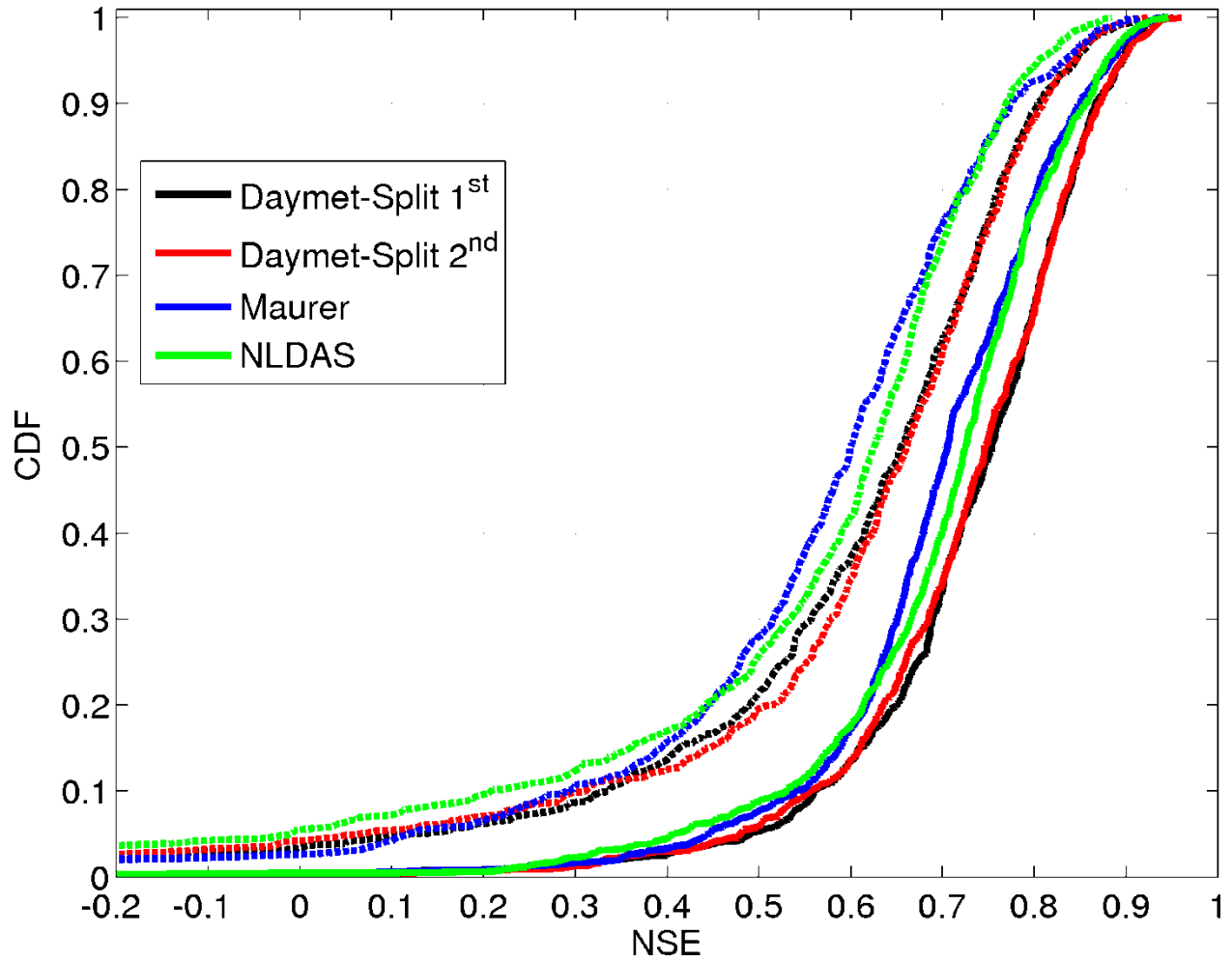
776 Figure 6. (a) Spatial distribution of the high flow bias, (b) low flow bias, (c) flow duration curve
 777 bias for the validation period.

778



779
 780 Figure 7. (a) Cumulative density functions (CDFs) for model Nash-Sutcliffe efficiency (NSE)
 781 (solid) for the calibration (red) and validation periods (blue) and NSE using the long-term
 782 monthly mean flows (MNSE, dark shaded and dashed), CDFs for (b) simulated-observed flow
 783 correlation in the decomposition of the NSE, (c) for the variance bias in the decomposition of the
 784 NSE, and (d) total volume bias in the decomposition of the NSE.

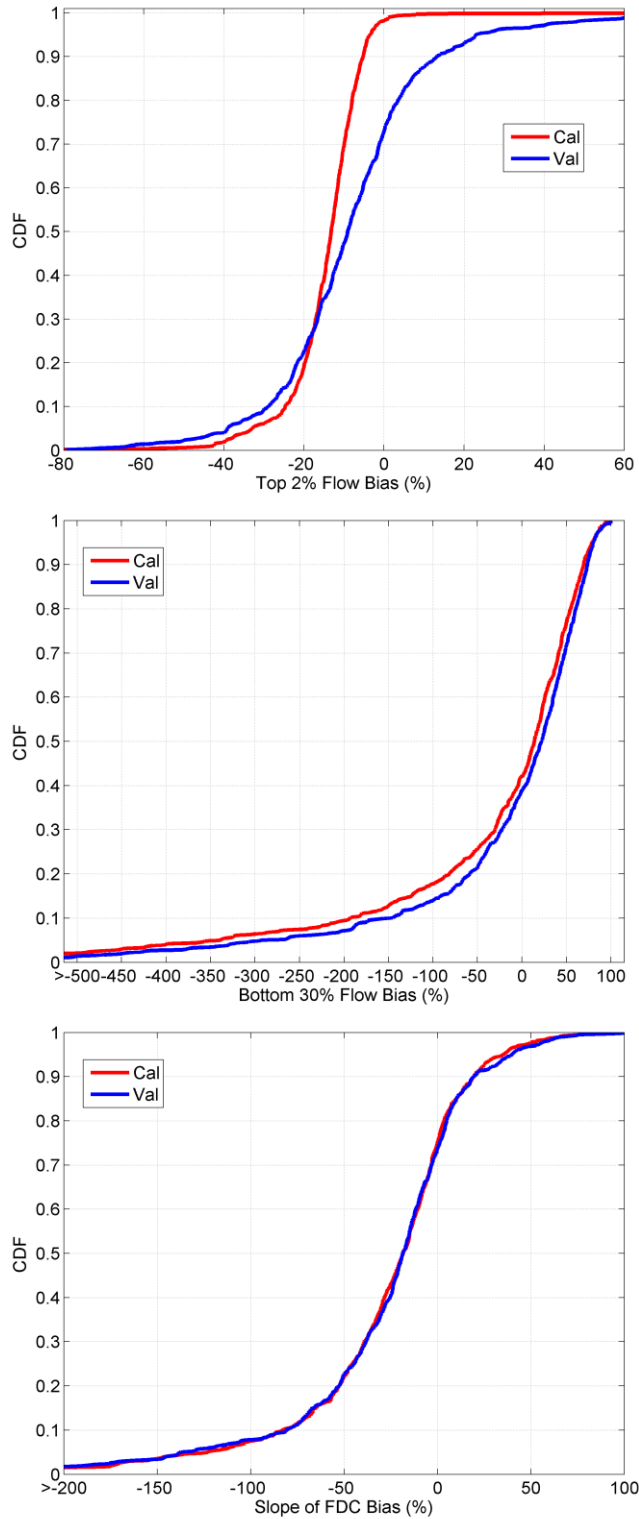
785



786

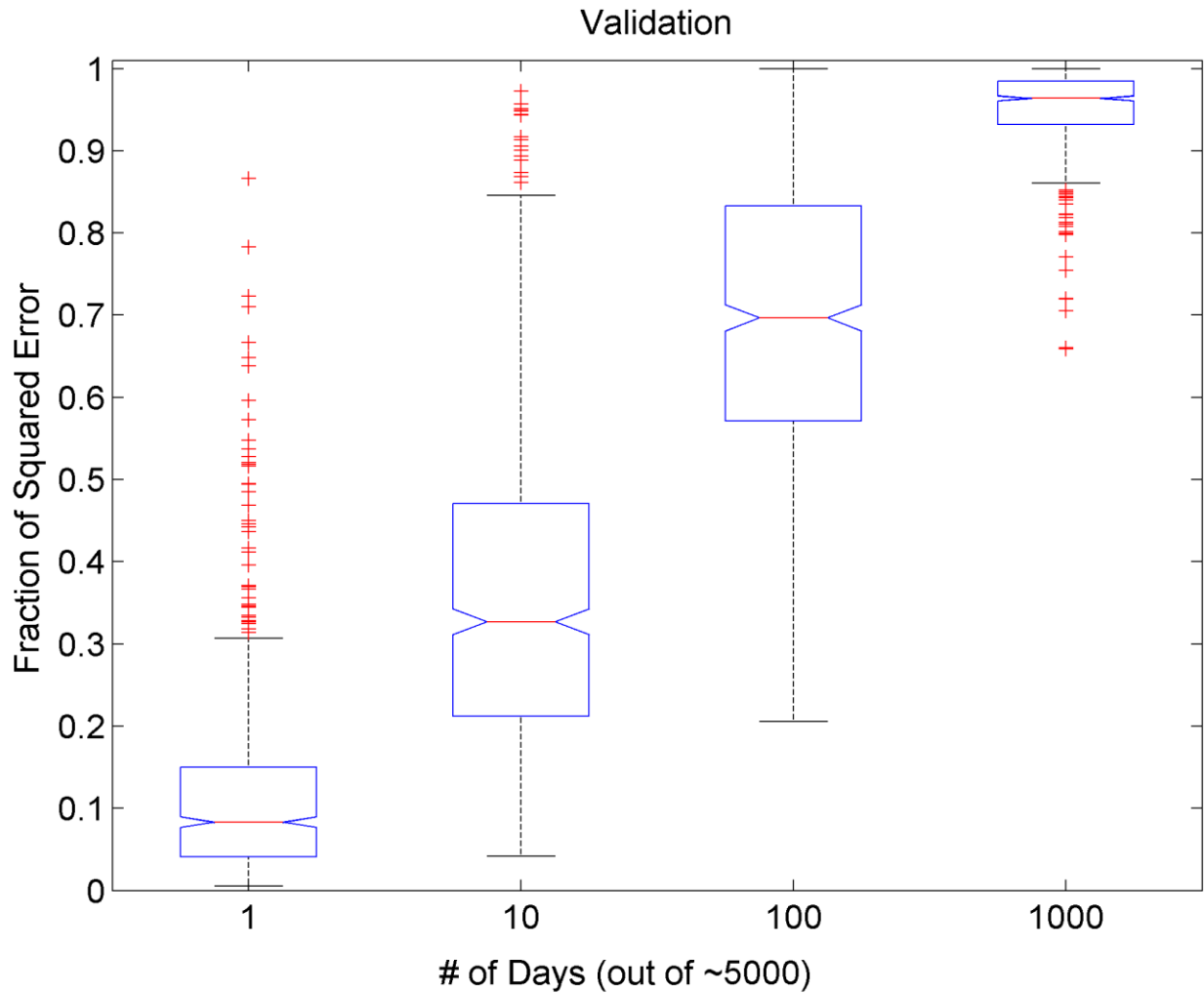
787 Figure 8. Cumulative density functions for model Nash-Sutcliffe efficiency for the calibration
 788 (solid) and validation (dashed) period using three different forcing datasets (Daymet, Maurer,
 789 NLDAS). The Daymet dataset was calibrated using the first 15 years (Split 1st) and validated
 790 against the remaining data and also calibrated using the last 15 years (Split 2nd) and validated
 791 against the initial streamflow data. Maurer and NLDAS calibrations performed using the first 15
 792 years of observed streamflow only.

793



794

795 Figure 9. (a) Cumulative density functions (CDFs) for model high flow bias for the calibration
 796 (red) and validation periods (blue), (b) model low flow bias, (c) model flow duration curve slope
 797 bias.

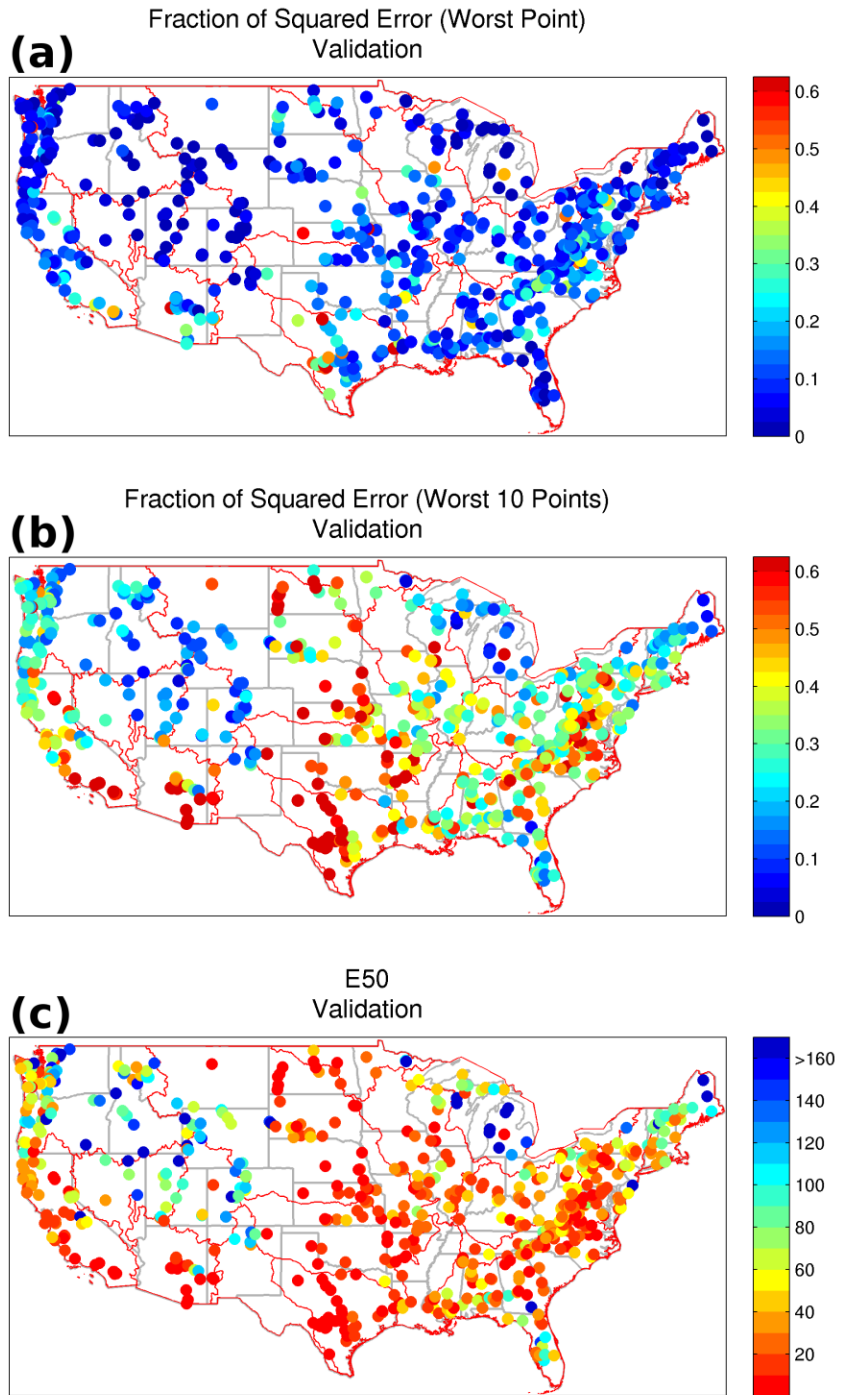


798

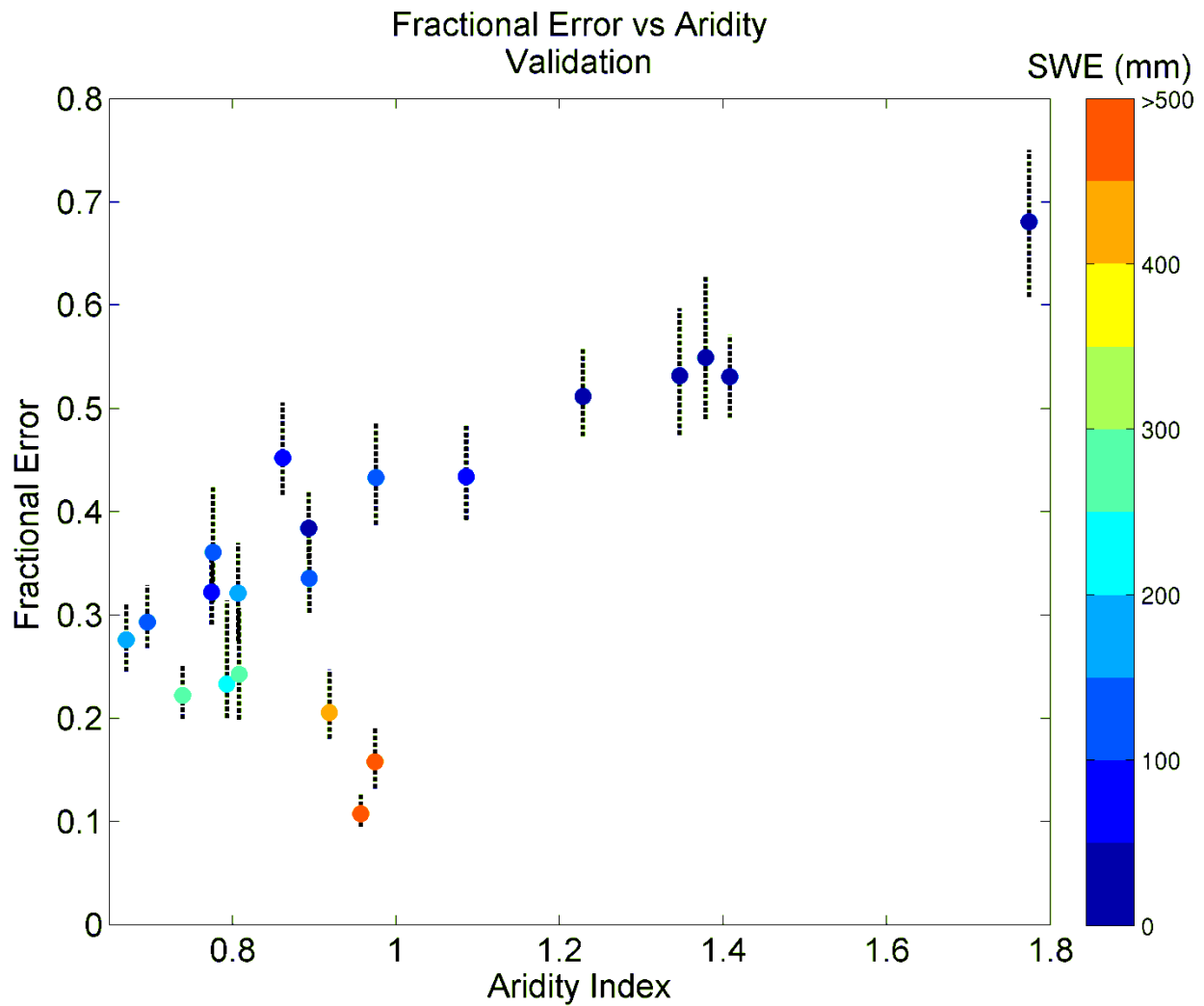
799 Figure 10. Fractional contribution of the total squared error for the 1, 10, 100, 1000 largest error
 800 days. The box plots represent the 671 basins with the blue area defining the interquartile range,
 801 the whiskers representing reasonable values and the red crosses denoting outliers. The median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

804

805



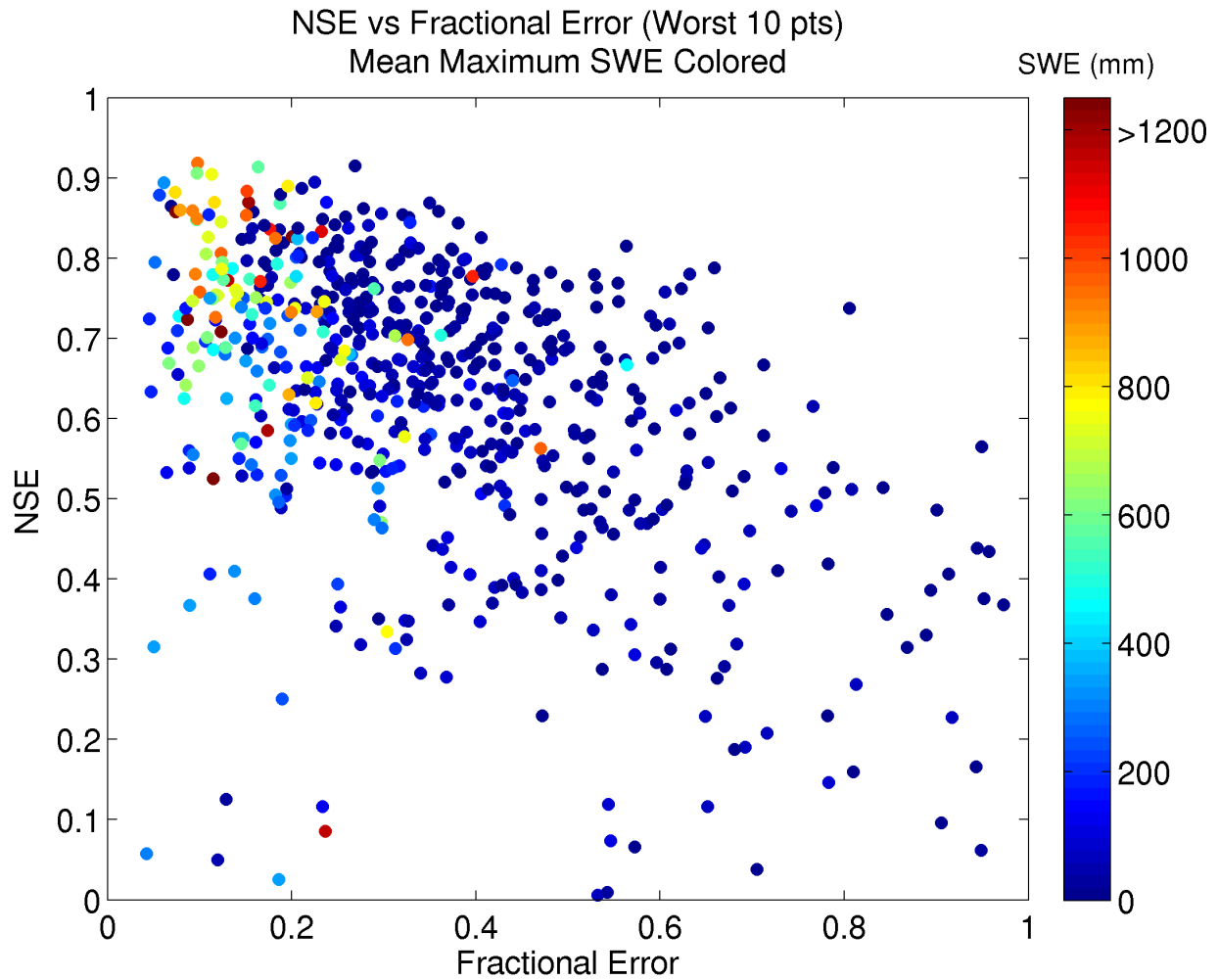
806
 807 Figure 11. (a) Spatial distribution of the fractional contribution of total squared error for the
 808 largest day during the validation period, (b) 10 largest error days, (c) the number of days
 809 contributing 50% of the total objective function error, E50.



810

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

816



817

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

820

821

822