

1 **Leveraging pre-storm soil moisture estimates for enhanced land surface model calibration**
2 **in ungauged hydrologic basins**

3
4 **Wade T. Crow¹, Jianzhi Dong^{1,2} and Rolf H. Reichle³**

5
6 ¹ USDA Hydrology and Remote Sensing Laboratory, Beltsville, MD

7 ² Massachusetts Institute of Technology, Cambridge, MA

8 ³ Global Modeling and Assimilation Office, NASA Goddard Space Flight Center, Greenbelt,
9 MD

10
11 Submitted to *Water Resources Research*

12
13 **Abstract**

14
15 Despite long-standing efforts, hydrologists still lack robust tools for calibrating land
16 surface model (LSM) streamflow estimates within ungauged basins. Using surface soil
17 moisture estimates from the Soil Moisture Active Passive Level 4 Soil Moisture (L4_SM)
18 product, precipitation observations, and streamflow gauge measurements for 617 medium-scale
19 (200-10,000 km²) basins in the contiguous United States, we measure the temporal (Spearman)
20 rank correlation between antecedent (i.e., pre-storm) surface soil moisture (ASM) and the
21 storm-scale runoff coefficient (RC; the fraction of storm-scale precipitation accumulation
22 converted into streamflow). In humid and semi-humid basins, this rank correlation is shown to
23 be sufficiently strong to allow for the substitution of storm-scale RC observations (available
24 only in basins that are both lightly regulated and gauged) with high-quality ASM values
25 (available quasi-globally from L4_SM) in streamflow calibration procedures. Using this
26 principle, we define a new, basin-wise LSM streamflow calibration approach based on L4_SM
27 alone and successfully apply it to identify LSM configurations that produce a high rank
28 correlation with observed RC. However, since the approach cannot detect RC bias, it is less
29 successful in identifying LSM configurations with low mean-absolute error.

31 **Plain Text Summary**

32 Accurately forecasting the fraction of rainfall that runs off into streams, as opposed to infiltrates
33 into the soil, is critical for flash-flood prediction, water-resource monitoring, and tracking the
34 transport of nutrients from agricultural fields into local waterways. Such forecasting is typically
35 performed by hydrologic models that attempt to represent the physical processes responsible for
36 surface runoff generation. However, to provide accurate streamflow forecasts, these models
37 typically need to be calibrated against actual streamflow observations. This is problematic
38 given the relatively poor, and declining, global availability of stream gauges. This paper
39 presents a novel model calibration strategy that uses soil moisture from remote sensing and
40 numerical modeling in place of streamflow observations during calibration. This transition has
41 significant practical advantages because, unlike streamflow observations, the soil moisture data
42 are continuously available across space. Our results demonstrate that this new approach can
43 significantly improve hydrologic models within humid and semi-humid basins lacking
44 sufficient ground-based instrumentation for traditional streamflow calibration.

45 46 **1. Introduction**

47 Despite several decades of development, land surface models (LSMs) still do not
48 generally provide adequate streamflow estimates outside of hydrologic basins in which they
49 have been directly calibrated (Xia et al., 2012; Hrachowitz et al., 2013). This is problematic due
50 to the limited, and declining, worldwide availability of streamflow gauge data (Fekete et al.,
51 2015) – as well as the proliferation of stream diversion and impoundment infrastructure that
52 degrades the quality of hydrologic information contained in streamflow observations. As a
53 result, there is a widely acknowledged need to develop effective LSM calibration strategies that
54 can be applied in the absence of reliable streamflow observations (Samaniego et al., 2017).

55 The estimation and routing of runoff is a multi-faceted process; however, one
56 fundamental aspect is the application of an LSM to estimate the fraction of storm-scale
57 precipitation converted into runoff (hereinafter, the runoff coefficient or “RC”). Storm-to-storm
58 variation in antecedent (i.e., pre-storm) soil moisture (ASM) is a well-known predictor of RC -
59 see, e.g., Song & Wang (2019) and references therein. However, past attempts to quantify the
60 impact of ASM on RC have been complicated by the presence of significant independent errors
61 in available ASM and RC estimates and the resulting attenuation bias in their sampled temporal
62 correlation (Crow et al., 2017). Attenuation bias refers to the tendency for independent random
63 errors, present in either independent or dependent variables, to spuriously decrease sampled
64 cross-correlation values (Hutcheon et al., 2010). Recent work with the Soil Moisture Active
65 Passive (SMAP) Level 4 Surface and Root-zone Soil Moisture (L4_SM) product suggests that,
66 once attenuation bias is minimized via the application of high-quality L4_SM ASM estimates,
67 storm-to-storm variations in ASM can be shown to play a dominant role in driving RC temporal
68 variability within the central and eastern United States (Crow et al., 2019). This result stands in
69 contrast with the typical representation of RC in LSMs - which generally predict a weaker role
70 for ASM in determining RC (Crow et al., 2018; 2019).

71 The apparent strength of the true coupling between ASM and RC presents an
72 opportunity for LSM streamflow calibration. Crow et al. (2019) suggest that, in certain cases,
73 the relationship between ASM estimates acquired from the L4_SM product and storm-scale RC
74 observations is sufficiently strong for the two quantities to be used inter-changeably in
75 correlation-based calibration objective functions. This is notable because the L4_SM product is
76 available globally - while meaningful RC observations are restricted to a relatively small
77 number of lightly regulated hydrologic basins with streamflow measurements available at their

78 outlet. Therefore, if L4_SM ASM and gauge-based RC values can be applied interchangeably,
79 new opportunities exist for expanding the meaningful calibration of LSMs into ungauged or
80 highly regulated basins – at least in areas like the contiguous United States (CONUS) where the
81 L4_SM product is known to provide high-quality ASM information (Crow et al., 2017).

82 Hereinafter, we will refer to the assumption of perfect temporal (i.e., storm-to-storm)
83 rank correlation between ASM and RC as the “perfect correlation” (PC) assumption. Our two
84 key objectives are to: i) evaluate the strength of the PC assumption using observations and ii)
85 investigate the potential of the PC assumption as a streamflow calibration principle for LSMs.

86 To achieve our first objective, we use ASM estimates acquired from the L4_SM product
87 (Section 2.2) and RC estimates based on streamflow and rainfall observations (Section 2.3)
88 obtained within the set of lightly regulated basins described in Section 2.1. Key steps towards
89 achieving this objective include the discrimination of individual storm events (Section 3.1),
90 modeling the impact of random errors on estimated ASM and RC values (Section 3.2), and the
91 identification of factors impacting the observed correlation between ASM and RC (Section
92 2.4). Results regarding the strength of the PC assumption are then described in Section 5.1.

93 Our second objective expands on the first by evaluating LSM calibration strategies
94 based on the PC assumption. These strategies employ an ensemble of LSM model
95 configurations (Sections 4.1-4.3) and are described in detail in Section 4.4. Evaluation metrics,
96 used to assess the performance of various calibration strategies, are introduced in Section 4.5.
97 Note that, in this context, the term “calibration” indicates the selection of an optimal LSM
98 ensemble member for an individual basin based on a particular calibration strategy. LSM
99 calibration results based on the PC assumption are then described in Sections 5.2-5.4.

100

101 **2. Domain and Data**

102 2.1. Study basins

103 Our analysis is based on an examination of ASM, precipitation, and streamflow data
104 within 617 medium-scale, lightly regulated CONUS basins during the period 1 April 2015 to 30
105 August 2020. Here, we provide background on our study basins and the datasets used to
106 examine the relationship between ASM and RC.

107 Study basins are based on a list of 1145 lightly regulated CONUS basins described in
108 Lohmann et al. (2004) and examined in Xia et al. (2012). From this original list, only basins
109 between 200 and 10,000 km² in size are considered here. Likewise, basins containing fewer
110 than 50 snow-free and frozen-soil-free storm events during our analysis period are discarded -
111 see additional details on our applied storm-event definition in Section 3.1. Finally, a small
112 number (< 10) of additional basins containing visible reservoirs (despite prior screening for
113 regulation), suffering from significant temporal gaps in daily USGS streamflow observations,
114 or providing clearly non-physical long-term streamflow statistics (e.g., mean-annual streamflow
115 exceeding mean-annual precipitation) are removed. Such screening results in the 617 selected
116 basins shown in Figure 1. These basins are generally restricted to the eastern half of CONUS –
117 along with a smaller number of basins along the west coast.

118 Daily USGS streamflow observations are acquired and processed for each basin outlet.
119 Note that, in the interest of maximizing the spatial coverage of our analysis, the rain-gauge
120 density threshold suggested by Schaake et al. (2000) is not applied. Therefore, significant
121 precipitation measurement errors are still possible. Likewise, despite our best efforts, the
122 absence of small-scale anthropogenic impoundment or diversion structures cannot be
123 guaranteed.

124

125 2.2 SMAP L4_SM product

126 ASM values are based on the area-weighted spatial averaging of 9-km resolution surface
127 (0-5 cm) soil moisture values acquired from Version 5 of the L4_SM product (Reichle et al.,
128 2019; 2020; 2021a) within each study basin. The 3-hourly L4_SM product is generated through
129 the sequential assimilation of SMAP brightness temperature data (Piepmeier et al., 2017) into
130 the NASA Catchment LSM (Reichle et al., 2017). Hourly, 0.25-degree surface meteorological
131 forcing data for the Catchment LSM is derived from the Goddard Earth Observing System
132 Forward-Processing (GEOS-FP) product (https://gmao.gsfc.nasa.gov/GMAO_products/; Lucchesi,
133 2018). Over CONUS, the GEOS-FP precipitation forcing is corrected to match daily
134 accumulations provided by the gauge-based NOAA Climate Prediction Center Unified (CPCU)
135 product at a 0.5-degree scale. Prior to the start of our analysis (1 April 2015), the Catchment
136 LSM is spun up from a cold start on 1 January 1980 using forcing data acquired from the
137 Modern-Era Retrospective Analysis for Research and Applications, Version-2 dataset (Gelaro
138 et al., 2017).

139 Daily ASM values are based on L4_SM surface soil moisture estimates within the 3-
140 hour time window (centered at 2230 UTC) closest to the end of each 0 to 24 UTC day.
141 Sampling at the end of the UTC day ensures that such ASM values are acquired as close as
142 possible to the potential start of a rainfall event on the following day.

143 Prior work has established that the L4_SM product provides a significantly better pre-
144 storm indicator of ASM than other available SM products - including a model-only version of
145 the L4_SM product that does not assimilate SMAP brightness temperature (Crow et al., 2017).
146 Likewise, L4_SM pre-storm surface (0 to 5-cm) soil moisture estimates are used because past

147 work suggests that they correlate slightly better with basin runoff response than corresponding
148 root-zone (0 to 1-m) L4_SM estimates (Crow et al., 2018).

149

150 2.3 NLDAS-2/CPCU precipitation

151 Daily (i.e., 0 to 24 UTC) precipitation totals are based on the spatial weighted averaging
152 of 0.125-degree gridded NLDAS-2 precipitation (Xia et al., 2009) estimates falling within each
153 basin. At a daily scale, these estimates are designed to match a 0.25-degree gauge-only CPCU
154 product (Chen et al., 2008) and corrected for topographic effects by the Parameter-elevation
155 Regressions on Independent Slopes Mode precipitation climatology dataset (Daly et al., 1994).
156 Note that we use the NLDAS-2 forcing data, instead of the GEOS-FP precipitation data applied
157 to force the land model in the L4_SM system, to maximize the independence of ASM and RC
158 estimates used in our analysis of the PC assumption. Nevertheless, there is considerable overlap
159 in the rain gauge data used as the basis for both. The potential impact of cross-dependency in
160 the precipitation data used to generate both ASM and RC is discussed further in Section 3.2.

161

162 2.4 Aridity index

163 Aridity index (AI) values, applied below to explain spatial trends in observed ASM
164 versus RC coupling, are taken from the Global Aridity and Potential Evaporation Dataset
165 (Zomer et al., 2007; Zomer et al., 2008). This product conforms to the AI definition offered by
166 the United Nations Environment Program (UNEP, 1992) whereby AI is the dimensionless ratio
167 between long-term, mean-annual precipitation divided by long-term, mean-annual potential
168 evapotranspiration (i.e., low AI values correspond to arid climates and high AI values to humid
169 climates).

170 AI values based on this definition are typically bounded between 0 and 1.5. Any
171 absolute labelling of AI values is somewhat subjective; however, a value of 0.4 approximates
172 the well-known wet/dry climate transition along the 100th meridian in CONUS. Note that AI
173 values are based on climatological averages sampled from long-term (1950 to 2000)
174 observations and an implied assumption of climate stationarity. As such, their temporal support
175 does not correspond directly to our 1 April 2015 to 30 August 2020 analysis period.

176

177 **3. Approach**

178 3.1 Storm events and SRCS

179 Storm-event separation is based on the approach of Crow et al. (2017) where a new
180 storm event is assumed to start on any day with a NLDAS-2/CPCU precipitation accumulation
181 exceeding P_{\min} . Unless otherwise noted, $P_{\min} = 10 \text{ mm d}^{-1}$.

182 Following a triggering daily rainfall accumulation, storm events are assumed to last for
183 an N -day period defined by rounding up the basin saturation time expression of Linsley et al.
184 (1982) to the nearest positive integer:

$$185 \quad N [\text{days}] = \text{CEIL}[(A * 2.59)^{0.2}] \quad (1)$$

186 where A is basin area [km^2] and CEIL is an upward integer rounding function. As a result, N is
187 meant to capture the period in which streamflow can be attributed to a given storm event. Here,
188 it is assumed to be independent of rainfall-event size. Derived values of N range from 4 days
189 for our smallest ($\sim 200 \text{ km}^2$) basin to 8 days for our largest ($\sim 10,000 \text{ km}^2$) basin.

190 ASM refers to the lowest end-of-day (i.e., 24 UTC or the closest available alternative),
191 basin-averaged L4_SM surface soil moisture (0-5 cm) value for the two-day period preceding
192 the start of a storm event, as will be discussed further below. Storm events interrupted by a new

193 storm (i.e., the arrival of another daily precipitation accumulation exceeding P_{\min} within an
194 earlier N -day storm event) are discarded, and a new event is assumed to begin coincident with
195 the latest triggering precipitation event.

196 Since our analysis focuses only on hydrological responses to rainfall incident on
197 unfrozen and non-snow-covered soil, only storm events where the pre-storm, basin-averaged
198 snow fractional cover is below 1% (by area) and the 24 UTC surface temperature is above 3° C
199 on the first (i.e., triggering) day of the event are considered. Surface state results from a
200 baseline LSM configuration (see Section 4 below) are used for assessing both thresholds –
201 which mask snow, rain-on-snow, and rain-on-frozen-soil storm events from our analysis.

202 As noted above, storm-scale RC is calculated as the streamflow volume during the
203 storm event divided by total precipitation accumulation volume (during the same storm-event
204 sampling period). Within basins containing more than 50 valid storm events during the study
205 period (1 April 2015 to 30 August 2020), soil moisture runoff coupling strength (SRCS) is
206 defined as the sampled Spearman rank correlation between ASM and RC values across all
207 storm events. A minimum of 50 qualifying events is required, reflecting a trade-off between the
208 competing concerns of maximizing the spatial coverage of our analysis versus adequately
209 filtering basins where SRCS values cannot be accurately sampled.

210 As discussed above, the perfect rank correlation (PC) assumption dictates that SRCS =
211 1. That is, it assumes that rank variations in RC across multiple storm events can be captured
212 perfectly given appropriate knowledge of ASM. As a correlation-based metric, SRCS is not
213 impacted by the potential presence of bias in the L4_SM product.

214 The separation of a continuous rainfall time series into discrete multi-day events
215 introduces potential ambiguities into the calculation of SRCS. One potentially problematic case

216 is when rainfall begins before 24 UTC on a given day (without exceeding the P_{\min} daily
217 accumulation threshold) and then continues into the following day – when the daily P_{\min}
218 accumulation threshold is exceeded and, as a result, a new storm event begins. In this case,
219 “end-of-day” surface soil moisture (SSM) may be enhanced on the day prior to the event,
220 because rainfall actually started before 24 UTC, and storm-scale RC will also be spuriously
221 increased, since rainfall during the previous day will be neglected in the “storm-scale”
222 calculation of RC. Such simultaneous enhancement to both ASM and RC could, conceivably,
223 inflate sampled SRCS values. To combat this, we define ASM as the minimum end-of-day
224 SSM for the two-day period prior to the start of a storm event. Therefore, in the case where
225 ASM is spuriously inflated by an event starting before 24 UTC, ASM will be defined using
226 end-of-day SSM for the *previous* day - and thereby avoid any spurious increase in ASM.

227 This approach has the benefit of not discarding any qualifying storm events. A more
228 conservative approach is to simply discard events that are preceded by more than trace amounts
229 of daily precipitation (defined here as any daily accumulation exceeding $P_{\min}/5$). While this
230 causes a significant reduction in the number of storm events available for sampling SRCS, it
231 also provides an important check that our SRCS results are not being spuriously impacted by
232 rainfall events crossing over the 24 UTC demarcation. Therefore, key results below will be re-
233 generated using this more stringent masking procedure to ensure that our main SRCS results are
234 reliable (see Section 5.1 and Figure S.1 of the Supporting Information).

235

236 3.2 Error models for ASM and RC estimates

237 As described above, a key focus of our analysis is the absolute value of SRCS sampled
238 from noisy ASM and RC estimates. As a correlation-based metric, SRCS will be spuriously

239 biased towards zero (i.e., attenuated) via the presence of independent, random error in L4_SM
240 ASM and/or USGS/CPCU RC estimates. This attenuation bias is not associated with sample
241 size limitations and will persist even in the theoretical case of an infinite sample size (Dong and
242 Crow, 2018). It is, therefore, conceptually distinct from the representation of confidence
243 intervals that vary as a function of sample size.

244 To estimate attenuation bias in sampled SRCS values, we first define random error
245 models for the L4_SM ASM analysis and the NLDAS-2/CPCU precipitation observations
246 underlying the RC estimates. L4_SM estimates are assumed to be degraded by mean-zero,
247 time-independent, additive Gaussian random error with a standard deviation of $\sigma_{L4_SM} = 0.032$
248 m^3m^{-3} . This estimate is based on L4_SM ground validation results against in situ measurements
249 at the 36-km scale (Reichle et al., 2017) and a minor adjustment to correct for random
250 uncertainty in the ground measurements themselves (Chen et al., 2019). While there is almost
251 certainly spatial and temporal variability in σ_{L4_SM} (Qiu et al., 2021), such variability is difficult
252 to assess and accurately reflect in an error model and therefore neglected here

253 Likewise, NLDAS-2/CPCU daily precipitation observations are assumed to be impacted
254 by random, multiplicative errors sampled from a log-normal distribution with unit mean and
255 standard deviation σ_{CPCU} . Here, σ_{CPCU} is estimated as a discrete function of the time-average
256 number of CPCU rain gauges (N_G) contained within each of our study basins using the relative-
257 accuracy versus gauge-density relationship summarized by Villarini et al. (2008). Specifically,
258 Figure 9 in Villarini et al. (2008) suggests that:

$$\begin{aligned} 259 \quad \sigma_{CPCU} &= 0.10 \quad \text{for } N_G \geq 9 \\ 260 \quad \sigma_{CPCU} &= 0.15 \quad \text{for } 5 \leq N_G \leq 8 \\ 261 \quad \sigma_{CPCU} &= 0.20 \quad \text{for } N_G = 4 \end{aligned} \tag{3}$$

262 $\sigma_{\text{CPCU}} = 0.25$ for $N_G = 3$

263 $\sigma_{\text{CPCU}} = 0.30$ for $N_G = 2$

264 $\sigma_{\text{CPCU}} = 0.40$ for $N_G = 1$.

265 The spatial density of the CPCU rain gauge network underlying our daily precipitation
266 estimates varies greatly in time and across CONUS. Therefore, gridded reports of station
267 densities underlying CPCU rain gauge estimates for a representative sample of days between 1
268 April 2015 to 30 August 2020 are used to estimate basin-specific values of N_G .

269 Given estimates of σ_{CPCU} and $\sigma_{\text{L4_SM}}$ for each basin, we can numerically estimate the
270 value of SRCS (i.e., SRCS_{PC}) you would expect to sample from observed ASM and RC time
271 series if the PC assumption holds and observed SRCS is degraded only by attenuation bias
272 associated with random observation error in either ASM or RC. This is done by first sorting
273 observed pre-storm L4_SM ASM and USGS/CPCU storm-scale RC values (separately) such
274 that their rank-correlation is unity and then adding random independent errors to both daily
275 precipitation values and ASM consistent with the error models described above. By re-ordering
276 the observed time series in the first step, we create ASM and RC time series that are (by
277 construction) consistent with the PC assumption. Therefore, when we subsequently add random
278 independent errors to these re-ordered time series, their rank correlation (i.e., SRCS_{PC}) reflects
279 what we would expect from our observed data if the PC assumption was valid. Based on this
280 logic, we can re-sample the Spearman rank correlation between the synthetically perturbed
281 ASM and RC values to obtain a basin-specific estimate of SRCS_{PC} .

282 This approach is repeated 5,000 times to ensure statistical convergence of the generated
283 SRCS_{PC} values. Comparisons between sampled SRCS and SRCS_{PC} values can then be used to
284 evaluate the validity of the PC assumption. If the PC assumption is valid, SRCS and SRCS_{PC}

285 will be approximately equal. However, in cases where any time-varying factor besides ASM
286 significantly impacts RC, SRCS will be less than $SRCS_{PC}$.

287 Note that our estimates of $SRCS_{PC}$ are conservative in that they likely underestimate the
288 total magnitude of attenuation bias for multiple reasons. First, USGS streamflow observations
289 are assumed to be free of random error. Second, based on the latest SMAP validation results
290 employing longer records (Colliander et al., 2021), the bias-corrected standard error of SSM
291 estimates in the L4_SM product is $0.034 \text{ m}^3\text{m}^{-3}$ and thus slightly larger than the $0.032 \text{ m}^3\text{m}^{-3}$
292 value applied here. Finally, the dual use of the gauge-based CPCU product to estimate RC
293 (from NLDAS-2/CPCU; Section 2.2) and to force the land model in the L4_SM algorithm
294 (Section 2.3) should result in negative error cross-correlation between RC and ASM estimates
295 and thus represents an additional source of SRCS attenuation bias that is missing in our
296 estimates of $SRCS_{PC}$. For example, an overestimation of precipitation in the CPCU product
297 would result in an overestimation of L4_SM ASM and an underestimation of RC values (since
298 RC represents streamflow normalized by precipitation accumulation). The sign contrast in these
299 errors would, in turn, work to spuriously degrade the otherwise positive rank correlation
300 between true ASM and true RC. Therefore, if we could properly account for ASM and RC error
301 cross-correlation in our error modeling, our derived $SRCS_{PC}$ estimates would be even lower. In
302 summary, our error modeling approach is conservative in the sense that our simplifications
303 (described above) likely result in $SRCS_{PC}$ values that underrepresent the actual impact of
304 attenuation bias. The implications of this will be discussed later in Section 5.

305

306 **4. Land Surface Modeling**

307 The methodology described above in Section 3 will be used to address our first main
308 objective - the observation-based evaluation of the PC assumption. In contrast, land surface
309 modeling, described here in Section 4, is central to addressing our second objective - the
310 assessment of the PC assumption as a viable LSM calibration strategy. Our target LSM for this
311 objective is the Noah with Multi-Parameterization options (Noah-MP) land model (Niu et al.,
312 2011).

313

314 4.1 Noah-MP set-up

315 All Noah-MP simulations are based on Version 7.2 of the NASA Land Information
316 System (LIS) (Kumar et al., 2006) and 15-minute/0.125°-resolution Noah-MP v3.6 integrations
317 between 1 April 2015 and 30 August 2020. Off-line Noah-MP forcing is based on the NLDAS-
318 2 meteorological dataset utilizing North American Regional Reanalysis variables for all fields
319 except precipitation, which instead uses the NLDAS-2/CPCU precipitation dataset described in
320 Section 2.3. Within each study basin, end-of-day (24 UTC) SSM (0-10 cm) and daily (0-24
321 UTC) runoff totals (i.e., surface runoff plus baseflow) are spatially aggregated to generate a
322 daily, basin-scale time series. Note that the 0-10 cm definition of SSM applied in Noah-MP is
323 slightly deeper than the 0-5 cm depth assumed in the L4_SM product. The impact of this
324 vertical discrepancy will be discussed below.

325

326 4.2 Noah-MP configuration ensemble

327 The Noah-MP LSM is unique in that it contains built-in options to utilize different
328 physical approaches for the representation of land surface water and energy balance processes.
329 Here, we leverage this flexibility to generate a 41-member ensemble of Noah-MP

330 configurations that reflects a range of approaches for representing the link between soil
331 moisture and both runoff and evapotranspiration. Calibration results are based on selecting a
332 single member of this ensemble, separately for each basin, that maximizes a particular
333 calibration objective function (see Section 4.4 below).

334 Given our emphasis on the representation of runoff, we start with the 16-member Noah-
335 MP configuration ensemble defined by Crow et al. (2019) via the systematic variation of Noah-
336 MP runoff processes. Within this ensemble, separate Noah-MP simulations are generated for all
337 four Noah-MP runoff-physics packages described in Niu et al. (2011): the simplified
338 groundwater (SIM GW) case, the simplified TOPMODEL (SIM TOP) case, the free-drainage
339 (FD) lower-boundary assumption, and the surface runoff parameterization taken from the
340 Biosphere Atmosphere Transfer Scheme (BATS).

341 For each of these four baseline cases, Crow et al. (2019) selected, and systematically
342 varied, one key parameter to further generate an ensemble of 16 different Noah-MP runoff
343 configurations. For the SIM GW and SIM TOP cases, we selected the TOPMODEL f
344 parameter, which describes the decay of saturated hydraulic conductivity with depth. For the
345 FD case, we selected the *REFKDT* parameter, which modulates the impact of ASM on surface
346 runoff. For the BATS case, we selected the exponential parameter (q), which links the top 2-m
347 soil moisture and surface runoff. In total, four f variations in the SIM GW case, five f variations
348 in the SIM TOP case, four *REFKDT* variations in the FD case, and three q variations in the
349 BATS case were applied to generate the entire 16-member Noah-MP runoff-configuration
350 ensemble. Table S.1 in the Supporting Information and Crow et al. (2019) provide additional
351 details on this ensemble.

352 Given the overall dominance of evapotranspiration (ET) as a soil water loss mechanism,
353 and the coupling between LSM representations of ET and runoff (Koster and Milly, 1997), the
354 16-member Noah-MP runoff configuration ensemble of Table S.1 was augmented using an
355 additional 25-member ET-configuration ensemble. The coupling relationship between soil
356 moisture and ET is impacted by a range of LSM parameters and processes. However, as shown
357 in Dong et al. (2020), SSM-ET coupling strength in Noah-MP is highly sensitive to the unitless
358 parameter λ , which controls the nonlinear relationship between soil evaporation stress and SSM.
359 Therefore, a broad range of SSM-ET coupling strengths can be captured by Noah-MP
360 configurations utilizing a corresponding range (i.e., 1, 2, 3, 5, and 10) of λ values (Dong et al.,
361 2020). Therefore, each of these five baseline λ cases was run for the five separate baseline
362 runoff-configuration cases introduced above (i.e., the baseline parameterizations for the SIM
363 GW, FD, and BATS cases, plus two separate FD cases) to generate a 25-member ET-
364 configuration ensemble. For additional details on each configuration within the ET ensemble,
365 see Table S.2 in the Supporting Information and Dong et al. (2020).

366 This new 25-member ET-configuration ensemble is combined with our earlier 16-
367 member runoff-configuration ensemble to generate a final 41-member Noah-MP ensemble. All
368 calibration results are based on the basin-wise selection of individual Noah-MP configurations
369 within this new 41-member ensemble that maximize one of the objective functions discussed
370 below in Section 4.4. Therefore, this ensemble effectively represents the parameter space for
371 our calibration analysis. Naturally, there exists an extremely wide range of approaches for
372 generating LSM configuration ensembles, and the selection of any single approach is inherently
373 subjective. However, the most important consideration is whether the selected ensemble spans a
374 sufficiently wide range of outcomes to serve as the basis for a robust calibration analysis. In this

375 regard, our 41-member ensemble appears adequate. Across all 617 study basins, the median
376 range of Spearman rank correlation against USGS RC observations within our 41-member
377 ensemble is 0.46. This range falls no lower than 0.13 for any single basin. Therefore,
378 substantial and consistent performance spread is found in all basins between the best and worst
379 Noah-MP configurations contained within our 41-member ensemble.

380

381 4.3 Noah-MP spin-up

382 All members of the final 41-member Noah-MP ensemble are spun-up individually from
383 a cold start on 1 January 2010 until the start of our analysis on 1 April 2015. To demonstrate
384 the adequacy of a five-year spin-up period, we examined Noah-MP configurations utilizing a
385 groundwater model (i.e., SIM GW runoff physics) under the assumption that they possess the
386 most stringent spin-up requirements. Noah-MP runoff and root-zone soil moisture results in
387 these (groundwater model-based) configurations generally stabilized after about five years and
388 were only marginally impacted by further increasing the model spin-up period from 5 to 15
389 years. This is consistent with Crow et al. (2019) who found that a six-year spin-up period was
390 adequate for examining the relationship between Noah-MP ASM and RC estimates.

391

392 4.4 Noah-MP calibration strategies

393 This section provides a brief description of the LSM calibration strategies applied to the
394 41-member ensemble of Noah-MP configurations described in Section 4.2 above. For our
395 purposes, the term calibration refers to the selection of a single Noah-MP configuration from
396 this ensemble, separately for each basin, by maximizing one of the three objective functions
397 described below. The result being the definition of a single “calibrated” ensemble member for

398 each individual basin and each individual calibration strategy. Note that the defining
399 characteristic of all three calibration strategies listed below is their exclusive reliance on the
400 continuous L4_SM estimates. As a result, they can all be applied directly (i.e., without
401 extrapolation) to any basin - even highly regulated ones lacking stream gauges.

402 *4.4.1 PC calibration strategy*

403 The key implication of the PC assumption is that basin-scale RC values, available only
404 in gauged/unregulated basins, can be functionally replaced by ASM estimates from the L4_SM
405 product. To examine this possibility, we define the following calibration objective function:

$$406 \quad F_1 = R_s[\mathbf{ASM}_{L4}, \mathbf{RC}_{NOAHMP}]. \quad (4)$$

407 Here, R_s is a Spearman-rank correlation operator; \mathbf{ASM}_{L4} is a set of end-of-day antecedent (see
408 Section 2.2) SSM estimates from the L4_SM product; and \mathbf{RC}_{NOAHMP} is a set of storm-scale RC
409 estimates from a particular Noah-MP configuration (Section 4.2). If the PC assumption holds,
410 selecting a single Noah-MP configuration that maximizes F_1 will simultaneously maximize the
411 correlation of Noah-MP RC estimates versus true values. Hereinafter, the application of (4) as a
412 calibration objective function will be referred to as the “PC calibration” strategy.

413 *4.4.2. SSM calibration strategy*

414 As discussed in Koster et al. (2018), an alternative calibration strategy is maximizing:

$$415 \quad F_2 = R_p[\mathbf{SSM}_{L4}, \mathbf{SSM}_{NOAHMP}] \quad (5)$$

416 where R_p is the Pearson correlation operator; \mathbf{SSM}_{L4} are end-of-day SSM (i.e., mean SSM
417 within the final 3-hour window of the UTC day – see Section 2.2) estimates from the L4_SM
418 product; and \mathbf{SSM}_{NOAHMP} are comparable end-of-day SSM values obtained from a given Noah-
419 MP configuration (Section 4.2). As a result, maximizing F_2 within a particular basin selects
420 Noah-MP configurations that maximize the temporal correlation between L4_SM and Noah-

421 MP SSM estimates. Note that the complete record of daily SSM estimates is used in (5) - not
422 just the sub-set of days when they describe ASM conditions for a new storm event. Hereinafter,
423 the application of (5) as a calibration objective function will be referred to as the “SSM
424 calibration” strategy. Since it represents an alternative, and more direct, application of L4_SM
425 SSM estimates, the SSM calibration strategy presents a useful baseline for evaluating the PC
426 calibration strategy.

427 *4.4.3 PC+SSM calibration strategy*

428 Combined calibration approaches are also possible including:

$$429 \quad F_3 = Z(F_1) + Z(F_2) \quad (6)$$

430 where Z is a normalization function defined as $Z(\mathbf{X}) = [\mathbf{X} - \text{mean}(\mathbf{X})] / \text{std}(\mathbf{X})$, and \mathbf{X} is a set of
431 (time-invariant) F_1 or F_2 values calculated across all potential Noah-MP configurations for a
432 single basin. Hereinafter, the maximization of (6) will be referred to as the “PC+SSM
433 calibration” strategy.

434

435 4.5 Calibration evaluation

436 The individual Noah-MP configurations selected via the maximization of (4)-(6) are
437 assessed using two different evaluation metrics. The first metric is the Spearman rank
438 correlation (R_s) between the storm-scale RC estimates of a particular Noah-MP configuration
439 and observed RC values; this metric assesses the skill in detecting storm-to-storm variations in
440 runoff efficiency. (Spearman rank correlation is used instead of Pearson correlation to
441 accommodate the potential for modest levels of non-linearity in the relationship between ASM
442 and RC.) Since R_s is blind to the presence of bias in Noah-MP RC estimates, the mean-
443 absolute-error (MAE) of the Noah-MP RC estimates is used as the second evaluation metric.

444 The resulting R_s and MAE RC evaluation metrics are normalized relative to the best-
445 and worst-performing Noah-MP configurations found in each basin. For each basin i , each
446 calibration function j (i.e., the PC, SSM, or PC+SSM calibration strategies), and each RC
447 evaluation metric (i.e., R_s or MAE), we individually calculate the ratio:

$$448 \quad F_{C,ij} = (C_{ij} - W_i)/(B_i - W_i) \quad (7)$$

449 where C_{ij} is the RC evaluation metric achieved in basin i through the maximization of
450 calibration function j across the Noah-MP configuration ensemble, and B_i and W_i are the best
451 (i.e., lowest MAE or highest R_s) and worst (i.e., highest MAE or lowest R_s) RC evaluation
452 metrics, respectively, across the Noah-MP configuration ensemble for the same basin i .

453 The endpoints B_i and W_i are calculated via direct comparison to observed RC values,
454 whereas C_{ij} reflects an evaluation metric obtained without access to streamflow observations -
455 and instead relies wholly on L4_SM ASM estimates. Therefore, $F_{C,ij} = 1$ indicates that, even
456 without access to observed RC, a given calibration approach j applied in basin i has accurately
457 identified the single LSM configuration within the Noah-MP ensemble that optimizes a given
458 RC evaluation metric. In the absence of meaningful calibration, all members of the ensemble
459 will be equally likely, and, sampled across all basins i , $F_{C,ij}$ will have a median value near 0.5 .
460 Therefore, any net upward shift in the spatial distribution of F_C towards one can be regarded as
461 a positive calibration outcome.

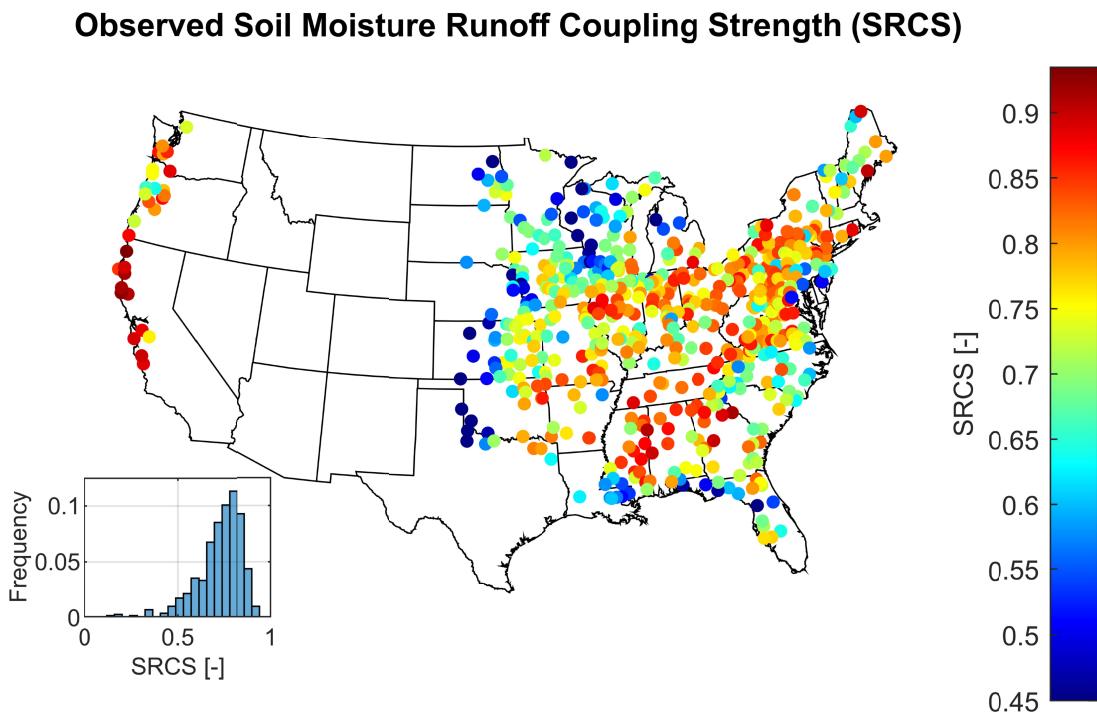
462

463 **5. Results**

464 5.1 Observation-based evaluation of the PC assumption

465 In this section, we evaluate SRCS and the PC assumption using only observed RC and
466 LS_SM data. Modeling results based on the calibration of the Noah-MP LSM are not discussed

467 until Section 5.2 below. Figure 1 plots observed SRCS values for the medium-scale CONUS
468 basins that meet the selection criteria introduced in Section 2.1. As described above, SRCS is
469 the sampled rank correlation between (pre-storm) ASM values (L4_SM) and storm-scale RC
470 (based on NLDAS-2/CPCU rainfall and USGS streamflow). Observed SRCS values are
471 relatively high with a median value of 0.75. Note that large areas of western CONUS lack
472 coverage due to our inability to sample at least 50 storm events that are not impacted by snow
473 or frozen soil during the historical SMAP period.
474



475

476 **Figure 1.** Observed (unitless) SRCS values for all 617 basins meeting the selection criteria
477 described in Section 2.1.

478

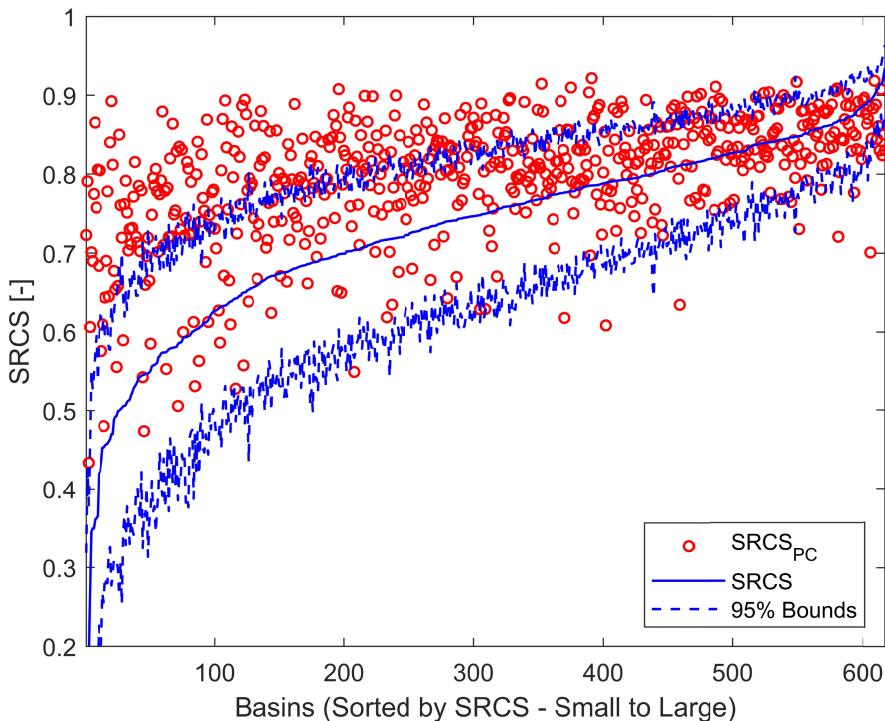
479 Assuming that random errors in ASM and RC are not positively correlated, sampled

480 SRCS values in Figure 1 will be spuriously reduced by attenuation bias. Therefore, a key

481 question is the degree to which non-unity values of SRCS in Figure 1 reflect the true sensitivity

482 of storm-scale RC variations to non-ASM factors (e.g., the time/space structure of precipitation
483 events) as opposed to the spurious impact of attenuation bias from random errors in the
484 underlying ASM and RC estimates.

485 As described above, $SRCS_{PC}$ is the hypothetical value of SRCS that we would expect to
486 sample, given attenuation bias from random errors in ASM and RC, if the PC assumption holds
487 and the true rank correlation between ASM and RC is one. As a result, sampled SRCS values
488 that are statistically indistinguishable from $SRCS_{PC}$ imply that the PC assumption is valid. In
489 contrast, SRCS values significantly less than $SRCS_{PC}$ suggest that attenuation bias alone cannot
490 explain the reduction of sampled SRCS below one. This implies that ASM variations cannot
491 explain all observed storm-to-storm variability in RC, and additional factors – presumably
492 related to dynamic land cover variations and/or the exact time/space structure of intra-storm
493 precipitation – must also be considered when modeling the storm-scale RC response. Note that
494 attenuation bias is conceptually distinct from sampling uncertainty in that it does not converge
495 towards zero with increasing sample size.



497

498 **Figure 2.** Observed SRCS (solid blue line) and $SRCS_{PC}$ (red circles) for each basin sorted by
 499 SRCS. SRCS uncertainty bounds (dashed blue lines) represent the 95% sampling confidence
 500 interval calculated based on a 5000-member bootstrapping analysis. $\Delta SRCS$ values shown in
 501 Figure 3 below are defined as $SRCS_{PC}$ minus SRCS.

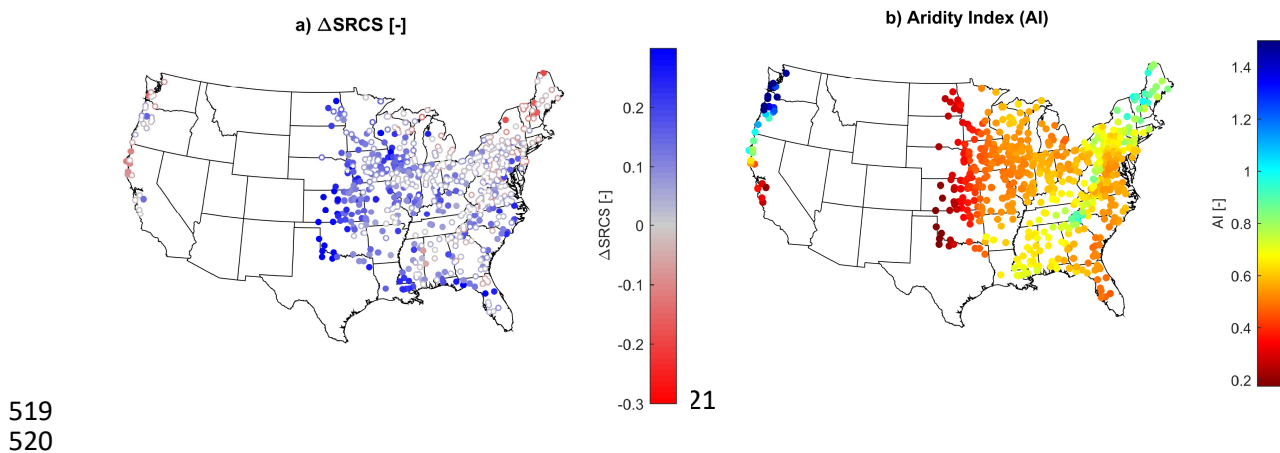
502

503 Figure 2 compares the observed SRCS with estimated values of $SRCS_{PC}$ for all basins in
 504 Figure 1 (sorted by SRCS), along with 95% confidence sampling intervals for SRCS. For 374
 505 of our 617 study basins (~61%), the estimated $SRCS_{PC}$ values fall within these intervals. In
 506 these basins, ASM appears to play a dominant role in controlling storm-to-storm variations in
 507 RC.

508 As described in Section 3.2, an analysis of potential bias in $SRCS_{PC}$ estimates suggests
 509 that Figure 2 provides a conservative assessment regarding the strength of the PC assumption
 510 (i.e., it likely under-represents the fraction of basins where $SRCS_{PC}$ values fall within the 95%
 511 uncertainty interval for sampled SRCS values). Furthermore, results in Figure S.1 (see the

512 Supporting Information) suggest that the impact of inter-day storm events, which could
513 conceivably induce a spurious positive bias into SRCS (see Section 3.2), is negligible.

514 The difference $\Delta\text{SRCS} = \text{SRCS}_{\text{PC}} - \text{SRCS}$ quantifies the impact of non-ASM factors on
515 RC temporal variability and, as such, provides a metric for evaluating the validity of the PC
516 assumption. Figure 3a plots ΔSRCS for each basin in Figure 1. As discussed above, ΔSRCS
517 values near zero support the PC assumption while $\Delta\text{SRCS} > 0$ implies that factors other than
518 ASM significantly influence RC. For reference, Figure 3b plots AI values for each basin.



522 **Figure 3.** a) ΔSRCS and b) AI for our 617 study basins. Open circles in a) denote basins where
523 ΔSRCS values are not significantly different from zero (at 95% confidence). AI is defined as
524 mean-annual precipitation divided by mean-annual potential evapotranspiration. Therefore,
525 lower AI values indicate generally drier conditions.

526
527 ΔSRCS values in Figure 3a are generally statistically insignificant in the relatively
528 humid areas of the northeastern U.S., the mid-Atlantic region, the upper Midwest, and along the
529 spine of the Appalachian Mountains. This is consistent with the expectation that, for humid
530 areas with (at least) modest levels of topographic variability, the dominant runoff mechanism is
531 so-called “saturation-excess” runoff, and the most important factor driving storm-to-storm
532 variations in RC is the fractional area of the basin saturated from below by a dynamic water
533 table (Dunne and Leopold, 1978; Zhao et al., 1980). There is a close conceptual link between

534 this spatial saturation fraction and ASM (Castillo et al., 2003) and, therefore, little observed
535 difference between SRCS and $SRCS_{PC}$.

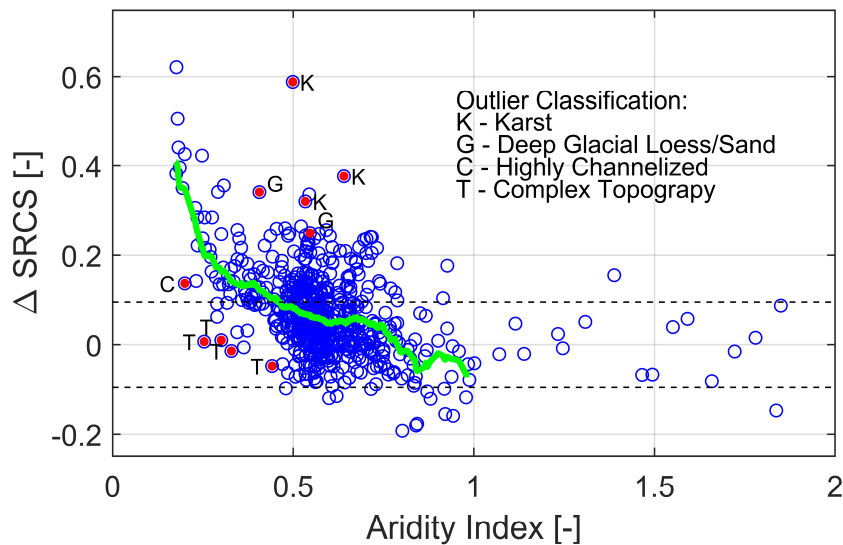
536 However, over more arid areas, one would generally expect a lower water table, reduced
537 surface saturation from below, and, therefore, relatively more emphasis placed on non-ASM
538 factors, like the time/space structure of precipitation intensity, when describing storm-to-storm
539 variations in RC (Zhang et al., 2011). This too is reflected in Figure 3a, where generally larger
540 (positive) values of $\Delta SRCS$ are noted in more arid (i.e., $AI < 0.4$) and flatter areas of central
541 CONUS.

542 Isolated negative $\Delta SRCS$ values in northeast and northwest CONUS (see Figure 3a) are
543 difficult to explain physically and likely reflect sampling errors in SRCS estimates, which tend
544 to be larger in northern CONUS basins due to the masking of storm events for snow and frozen
545 soil and/or the possible regional overestimation of $L4_SM$ or CPCU precipitation errors in the
546 statistical models we applied to estimate $SRCS_{PC}$.

547 Figure 4 summarizes the overall spatial relationship between basin $\Delta SRCS$ and AI. For
548 $AI > 0.4$, $\Delta SRCS$ values are often indistinguishable from zero and therefore roughly consistent
549 with the PC assumption. However, for $AI < 0.4$, $\Delta SRCS$ values trend significantly positive,
550 which indicates that the PC assumption is invalid for relatively arid conditions.

551 Observed outliers in the relationship between $\Delta SRCS$ and AI can often be associated
552 with the presence of unusual land surface or geological characteristics in the basin. For
553 example, multiple basins with $\Delta SRCS$ values well *above* the general $\Delta SRCS$ versus AI trend
554 line in Figure 4 contain very deep glacial soil deposits (either sand or loess) or complex karst
555 sub-surface geology (see symbol labelling in Figure 4). Both characteristics tend to dampen the
556 streamflow response to individual storm events and, therefore, reduce SRCS and increase

557 Δ SRCS. Conversely, outliers located well below the Δ SRCS versus AI trend line are frequently
 558 characterized by a high-degree of anthropogenic stream channelization and/or very-large
 559 topographic variability – characteristics that generally enhance the magnitude of streamflow
 560 response to storm events. Such enhancement tends to increase SRCS and, as a result, decrease
 561 Δ SRCS. Therefore, while AI explains a substantial fraction of observed Δ SRCS spatial
 562 variability, land surface and geologic characteristics also play an important role.



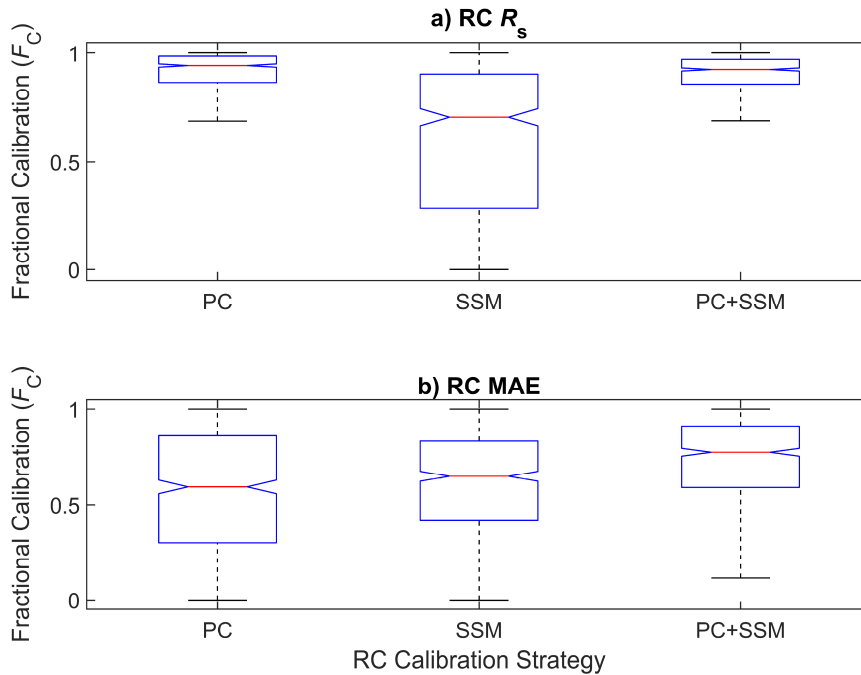
563

564 **Figure 4.** Δ SRCS versus AI values for our 617 study basins. Labelling identifies unusual land
 565 surface and geologic conditions in basins that can be visually identified as outliers. The green
 566 trend line represents median values within a 0.05-wide moving window applied along the
 567 abscissa. Dashed horizontal lines indicate 95% confidence intervals around the null hypothesis
 568 that Δ SRCS equals zero. Note that lower values of AI indicate generally drier conditions.
 569

570 5.2 Noah-MP calibration results

571 Despite considerable scatter, results in Figures 4 suggest that the PC assumption is
 572 potentially applicable for $AI > 0.4$. In such areas, it may be possible to calibrate LSMs by
 573 maximizing the temporal rank correlation between LSM-estimated RC and the L4_SM-based
 574 ASM in lieu of RC observations – see (4). Such an approach is particularly attractive in
 575 ungauged or heavily regulated basins where meaningful RC observations are not available.

576 Figure 5 examines this possibility by comparing Noah-MP RC calibration results for the
 577 PC, SSM, and PC+SSM calibration strategies introduced Section 4.4. Specifically, the figure
 578 shows box-and-whisker plots for the fractional calibration metric F_C defined in (7) across all
 579 617 study basins for both the R_s and MAE RC evaluation metrics.



580

581 **Figure 5.** Box-and whisker-plots (summarizing minimum, maximum, median, and inter-
 582 quartile spread values) for the basin-to-basin variation of Noah-MP F_C results for all three
 583 calibration strategies (i.e., PC, SSM, and PC+SSM) and the a) Spearman correlation (R_s) and b)
 584 MAE RC evaluation metrics.

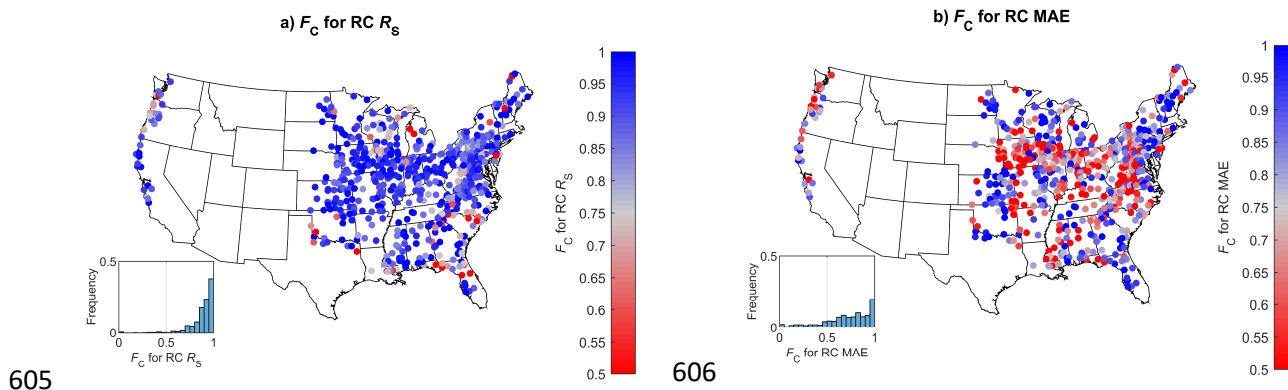
585

586 As demonstrated by the concentration of Noah-MP $R_s F_C$ values near one in Figure 5a,
 587 the PC calibration strategy consistently identifies a Noah-MP configuration for each basin that
 588 has high storm-to-storm Spearman rank correlation with observed RC values. In this regard, the
 589 PC calibration strategy is clearly superior to the direct SSM calibration strategy. However, the
 590 PC calibration strategy is markedly less successful in identifying Noah-MP configurations that
 591 exhibit low RC MAE (Figure 5b). Note that MAE F_C values for the PC strategy in Figure 5b
 592 demonstrate only a small net positive shift towards unity. Relative to the PC calibration

593 strategy, the combined PC+SSM calibration strategy provides modestly better MAE results
 594 (compare the MAE F_C results for PC and PC+SSM in Figure 5b) while simultaneously
 595 preserving its good R_s performance (Figure 5a).

596 Noah-MP calibration results in Figure 5 are based on utilizing our entire study period (1
 597 April 2015 to 30 August 2020) as both a training and a testing period. Naturally, such overlap
 598 raises concerns about overfitting. However, because all three calibration approaches, as well as
 599 the B and W values used to derive F_C in (7), are equally impacted by any potential overfitting,
 600 Noah-MP calibration results plotted in Figure 5 change very little when re-generated for the
 601 case of applying mutually exclusive training and testing periods. For reference, Figures S.2 and
 602 S.3 in the Supporting Information provide alternative, but still highly consistent, versions of
 603 Figure 5 generated for separate training and testing periods.

604



607 **Figure 6.** Maps and histograms of Noah-MP F_C results obtained for the PC+SSM calibration
 608 strategy based on the a) R_s and b) MAE RC evaluation metrics.

609

610 Spatial maps of Noah-MP RC R_s and MAE calibration results (as summarized by F_C)
 611 are shown in Figure 6 for the PC+SSM calibration strategy. With a few exceptions, Figure 6a
 612 illustrates consistently high F_C results for the Spearman rank evaluation of Noah-MP RC
 613 simulations after PC+SSM calibration. Note that some of the, relatively few, poorly performing

614 basins (identified in red in Figure 6a) were previously labeled as outliers in Figure 4 (e.g.,
615 highly karst basins in the panhandle region of Florida and basins with extremely deep glacial
616 sand deposits in Michigan). This underscores the challenge of dealing with highly unusual local
617 geology in any calibration strategy not based on local streamflow observations. Conversely, it is
618 surprising that RC R_s results for the PC+SSM strategy are not more clearly degraded in semi-
619 arid regions (e.g., west-central CONUS; Figure 6a), where the PC assumption is known to be
620 tenuous (see Figures 3-4). Given the uncertain physical basis of these results, apparently
621 successful calibration in semi-arid regions should be viewed skeptically.

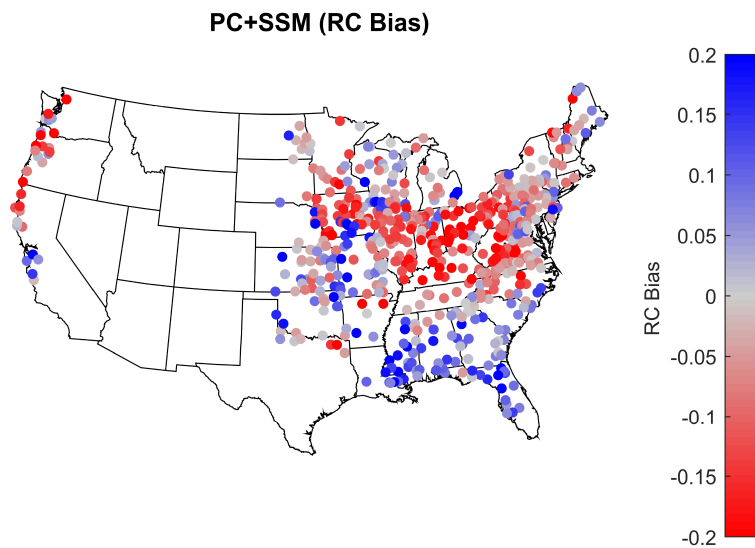
622

623 5.3 Role of RC bias

624 Despite the overall improvement noted in Figure 5b for the PC+SSM calibration
625 strategy versus the PC strategy, RC MAE values in many basins remain mediocre (grey-
626 shaded) or poor (red-shaded) following PC+SSM Noah-MP calibration (Figure 6b). MAE
627 metric values can be degraded (i.e., increased) by a variety of fit issues including additive bias,
628 multiplicative bias, and poor correlation against true values. As a result, it is worth considering
629 which component of RC MAE is driving these relatively poor results. Additional analysis (see
630 Figure S.4 in the Supporting Information) demonstrates that poor RC MAE F_C results for the
631 PC+SSM calibration strategy are strongly linked to the presence of long-term, basin-specific
632 RC biases in Noah-MP configurations identified as optimal by the PC+SSM calibration
633 strategy.

634 Figure 7 plots the spatial pattern of long-term RC bias in Noah-MP configurations
635 selected by the PC+SSM calibration strategy. A continuous area of negative RC bias is seen in
636 the northern half of CONUS while a positive RC bias dominates in southern CONUS. The

637 large-scale spatial coherence of this bias pattern suggests that it may be possible to parameterize
638 an additional RC bias-sensitive term using known regional climate or land cover characteristics
639 and add it onto the existing objective function in (6). However, the (correlation-based) PC and
640 PC+SSM calibration strategies considered here are not sensitive to bias and therefore cannot
641 meaningfully contribute to such a term. Instead, alternative, bias-sensitive analysis techniques
642 (e.g., long-term water balance calculations) will be required.



643

644 **Figure 7.** RC bias in Noah-MP configurations identified as optimal by the SSM+PC calibration
645 strategy. Note that the spatial pattern of RC bias drives the mediocre RC MAE results in Figure
646 6b.

647

648 5.4 Sensitivity to methodological changes

649 5.4.1 Storm-event threshold size

650 We re-generated Figures 1-3 after raising the storm-event threshold P_{\min} from 10 to 20
651 mm d^{-1} (not shown). Increasing P_{\min} sharply reduces the number of basins (from 617 to 296)
652 with at least 50 storm events that are not impacted by snow or frozen soil. However, within
653 these 296 basins, which tend to be concentrated in south-central and southeastern CONUS, the
654 median SRCS results for the $P_{\min} = 10 \text{ mm d}^{-1}$ and 20 mm d^{-1} cases are equal to within two

655 decimal places (0.74). As a result, there is currently no empirical indication that the PC
656 assumption weakens for larger values of P_{\min} .

657 Nevertheless, it is unreasonable to expect the PC assumption to hold indefinitely as P_{\min}
658 is raised without limit. Regardless of underlying ASM levels, storm-scale RC values will
659 eventually converge to unity for extremely intense precipitation events (e.g., P_{\min} approaching
660 or exceeding the expected annual maximum for daily precipitation in a single basin). Note that
661 the relatively short historical record of our analysis places a severe restriction on our ability to
662 sample such extreme events. Further raising P_{\min} to 30 mm d⁻¹, for example, results in only 36
663 basins with an adequate storm sample size. Therefore, results presented here cannot be
664 extrapolated to characterize basin response to major flooding events.

665 5.4.2 Vertical support of Noah-MP SSM estimates

666 As noted above, the vertical support of Noah-MP SSM estimates (0-10 cm) is deeper
667 than the corresponding 5-cm SSM estimates provided by the L4_SM product. The application
668 of a deeper (i.e., > 5 cm) surface soil layer is a common concession in LSMs to numerical
669 difficulties posed by balancing water within excessively thin soil layers. It is possible that this
670 vertical mismatch partially undermines the performance of the SSM calibration strategy posed
671 in (5). However, given the inherent difficulty of matching temporal scales of variability
672 expressed in SSM products obtained from different sources (Shellito et al., 2020), the apparent
673 robustness of the PC and PC+SSM calibration strategies to reasonable variations in surface soil
674 layer depth, and thus temporal SSM memory, is encouraging and likely necessary for their
675 practical application.

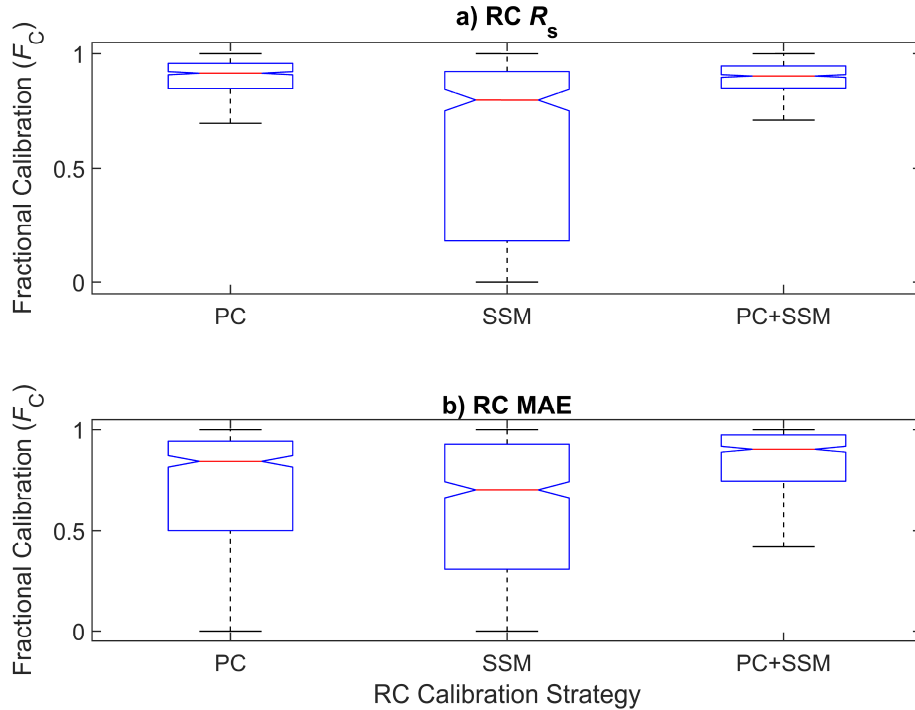
676 5.4.3 Baseflow separation

677 All calibration results presented above are based on bulk streamflow observations
678 and/or total (i.e., surface + baseflow) runoff estimates provided by Noah-MP. However, owing
679 to the particular importance of predicting fast streamflow response to rainfall, it is common in
680 hydrologic analysis to first remove the baseflow contribution to streamflow. Therefore, a
681 relevant question is whether the, generally good, RC calibration results reflected in Figure 5 for
682 the PC and PC+SSM calibration strategies remain relevant for an application focused only on
683 the fast storm response due to surface runoff.

684 To remove baseflow from the USGS streamflow observations, we applied the USGS
685 Hydrologic Separation [HYSEP; Sloto and Crouse, 1996] baseflow separation technique using
686 (non-integer) saturation time scales derived from (1) but without the application of the CEIL
687 operator. Note that Noah-MP already separately calculates both surface and baseflow runoff
688 components. Here, unlike above, we based Noah-MP storm-scale runoff estimates on only
689 surface runoff. We then re-generated Figure 5 for this surface-only runoff case, with results
690 shown in Figure 8. A comparison of the two figures reveals no apparent degradation in
691 calibrated RC results associated with considering only surface runoff contributions. In fact,
692 MAE results for the PC+SSM calibration strategy appear to be marginally improved following
693 the removal of baseflow. As a result, the PC or PC+SSM calibration strategies appear to be
694 equally effective at constraining RC estimates based on either surface-only (Figure 8) or total
695 (surface + baseflow; Figure 5) runoff.

696

697



698

699 **Figure 8.** Same as Figure 5 except for only the surface runoff response.

700

701 **6. Summary and conclusions**

702 Robust external constraints are needed for LSM-based simulation of streamflow in
 703 ungauged and/or heavily regulated basins that lack representative streamflow observations.
 704 While the potential for remote sensing to provide such constraints is widely acknowledged, the
 705 formulation of general physical principles to underlie these constraints remains elusive. Here,
 706 we hypothesize that the rank correlation between storm-to-storm variations in ASM and RC is
 707 unity (i.e., the PC assumption) and examine if this assumption can serve as a useful calibration
 708 principle for obtaining adequate LSM estimates of RC within ungauged basins.

709 Within a substantial fraction ($> 60\%$) of lightly regulated, median-scale basins in the
 710 central and eastern United States, there is indeed no significant evidence that time-varying
 711 processes other than ASM significantly impact storm-to-storm variations in RC (Figure 2). This

712 is particularly true for relatively humid ($AI > 0.4$) basins. To this end, Figures 3a and 4 suggest
713 that, at least over the eastern half of CONUS, the PC assumption is commonly valid.

714 The apparent magnitude of the rank correlation between ASM and RC suggests utilizing
715 high-quality, spatially continuous ASM estimates provided by the L4_SM product, in lieu of
716 RC observations, as a LSM calibration target in ungauged and/or heavily regulated basins that
717 lack suitable streamflow observations. This approach, summarized in (4) as the “PC calibration
718 strategy”, successfully identifies Noah-MP configurations with relatively good (i.e., high) RC
719 R_s results (Figure 5a). Such a strategy is robust for a range of storm-event rainfall thresholds
720 (i.e., 10-20 mm d^{-1}) and for hydrological applications focused on either total (i.e., surface +
721 baseflow) runoff or surface-only runoff (Figure 8). However, the PC calibration strategy is
722 generally blind to bias in RC estimates – which degrades RC MAE results for its selected
723 Noah-MP configurations (Figures 5b). RC MAE results can be marginally improved by
724 including an SSM-based correlation term in the calibration objective function – via the
725 “PC+SSM” calibration strategy summarized in (6) (Figure 5b). However, the correlation-based
726 calibration strategies we consider here cannot provide a full solution to this bias issue.

727 There are also geographic limitations to the confident application of the PC assumption.
728 However, these limitations are intuitively associated with known climate and land cover
729 characteristics (Figure 4) – raising hopes that basins where the PC assumption does not hold
730 can be readily identified and masked. We focus on medium-scale (i.e., 200 to 10,000 km^2)
731 basins for the verification of the PC assumption due to the difficulty of defining discrete storm
732 events and screening for streamflow regulation in larger basins. Nevertheless, once verified, we
733 are not aware of any barriers to applying a PC-based calibration strategy to larger basins.
734 Analogously, while the PC assumption can only be verified in lightly regulated basins, there is

735 no reason why in-channel streamflow regulation, in and of itself, should imperil its application
736 to runoff calibration. Finally, given that the integration of SMAP brightness temperature
737 information into the SMAP_L4_SM products has been shown to significantly improve SSM
738 estimation skill in data-poor regions (Dong et al., 2019; Reichle et al., 2021b), we are confident
739 that PC-based calibration strategies can be broadly applied in areas outside of CONUS.
740 However, future work will be needed to confirm this speculation.

741 Naturally, the evaluation of LSM runoff physics is a multi-objective exercise that
742 requires multiple metrics. While high RC correlation is necessary for adequate LSM
743 correlation, it is certainly not sufficient. For example, additional work is required to define a
744 bias-aware term for (6) that adequately penalizes Noah-MP configurations associated with
745 biased RC estimates (Figure 7). Potential approaches for deriving this term include the
746 application of remotely sensed ET products (obtained, for example, from thermal-infrared
747 remote sensing) or the use of improved river-stage height estimates expected from the
748 upcoming NASA Surface Water Ocean Topography mission (Biancamaria et al., 2016).
749 Likewise, a much longer historical analysis period is required to assess performance for
750 extreme rainfall events likely to cause flooding. Lacking this, care should be taken when
751 applying our proposed calibration approaches to LSMs tasked with flash flood forecasting.

752

753 **Acknowledgements**

754 This research was supported by grant 80HQTR21T0054 from the NASA SMAP mission
755 (“Using SMAP soil moisture products to improve streamflow forecasting in ungauged basins”,
756 PI: W.T. Crow). All USGS, NLDAS-2 forcing data and SMAP soil moisture data used in this
757 study are publicly available and can be accessed online (see in-text data citations). Final soil

758 moisture and runoff estimates obtained from the ensemble of Noah-MP simulations considered
759 here will be posted in an appropriate on-line data repository following acceptance of this draft
760 manuscript. USDA ARS is an equal opportunity employer.

761 **Work Cited**

762

763 Biancamaria, S., Lettenmaier, D. & Pavelsky, T. (2016). The SWOT Mission and Its
764 Capabilities for Land Hydrology. *Surveys in Geophysics*. 37(2), 307-337.
765 <https://doi.org/10.1007/s10712-015-9346-y4>.

766

767 Castillo, V.M., Gómez-Plaza, A. & Martínez-Mena, M. (2003). The role of antecedent soil
768 water content in the runoff response of semiarid catchments: A simulation approach. *J. Hydrol.*
769 284, 114–130. [https://doi.org/10.1016/S0022-1694\(03\)00264-6](https://doi.org/10.1016/S0022-1694(03)00264-6).

770

771 Chen, M., Shi, W., Xie, P., Silva, V. B. S., Kousky, V. E., Wayne Higgins, R. & Janowiak, J. E.
772 (2008). Assessing objective techniques for gauge-based analyses of global daily precipitation.
773 *J. Geophys. Res.* 113, D04110. <https://doi.org/10.1029/2007JD009132>.

774

775 Chen, F., Crow, W.T., Cosh, M.H., Colliander, A., Asanuma, J., Berg, A., Bosch, D.D.,
776 Caldwell, T.G., Collins, C.H., Jensen, K.H., Martínez-Fernández, J., McNairn, H., Starks, P.J.,
777 Su, Z. & Walker, J.P. (2019). Uncertainty of reference pixel soil moisture averages sampled at
778 SMAP core validation sites. *Journal of Hydrometeorology*. 20, 1553-1569.
779 <https://doi.org/10.1175/JHM-D-19-0049.1>.

780

781 Colliander, A. & coauthors. (2021). Validation of soil moisture data products from the NASA
782 SMAP mission. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote*
783 *Sensing*. 15, 364-392. <https://doi.org/10.1109/JSTARS.2021.3124743>.

784

785 Crow, W.T., Chen, F., Reichle, R.H. & Liu, Q. L. (2017). L-band microwave remote sensing
786 and land data assimilation improve the representation of prestorm soil moisture conditions for
787 hydrologic forecasting. *Geophysical Research Letters*. 44, 5495– 5503
788 <https://doi.org/10.1002/2017GL073642>.

789

790 Crow, W.T., Chen, F., Reichle, R. Xia, Y. & Liu, Q. (2018). Exploiting soil moisture,
791 precipitation and streamflow observations to evaluate soil moisture/runoff coupling in land
792 surface models. *Geophysical Research Letters*. 45(10), 4869-4878.
793 <https://doi.org/10.1029/2018GL077193>.

794

795 Crow, W.T., Chen, F., Reichle, R.H. & Xia, Y. (2019). Diagnosing bias in modeled soil
796 moisture/runoff coefficient correlation using the SMAP Level 4 soil moisture product. *Water*
797 *Resources Research*. 55, 7010-7026. <https://doi.org/10.1029/2019WR025245>.

798

799 Daly, C., Neilson, R.P. & Phillips, D.L. (1994). A statistical-topographic model for mapping
800 climatological precipitation over mountainous terrain. *J. Appl. Meteor.* 33, 140-158,
801 [https://doi.org/10.1175/1520-0450\(1994\)033<0140:ASTMFM>2.0.CO;2](https://doi.org/10.1175/1520-0450(1994)033<0140:ASTMFM>2.0.CO;2).
802

803 Dong, J. and Crow, W.T. (2018). The added value of assimilating remotely sensed soil moisture
804 for estimating summertime soil moisture - air temperature coupling strength. *Water Resources*
805 *Research*. 54, 6072-6084. <https://10.1029/2018WR022619>.
806

807 Dong, J., Crow, W.T., Reichle, R., Liu, Q., Lei, F. & Cosh, M. (2019) A global assessment of
808 added value in the SMAP Level 4 soil moisture product relative to its baseline land surface
809 model. *Geophysical Research Letters*. 46, 6604-6613. <https://doi.org/10.1029/2019GL083398>.
810

811 Dong, J., Dirmeyer, P. A., Lei, F., Anderson, M.C., Holmes, T.R.H., Hain, C. and Crow, W.T.
812 (2020). Soil evaporation stress determines soil moisture - evapotranspiration coupling strength
813 in land surface modeling. *Geophysical Research Letter*. 47, e2020GL090391.
814 <https://doi.org/10.1029/2020GL090391>.
815

816 Dunne, T. & Leopold, L.B. (1978) *Water in Environmental Planning*. Freeman, New York, 818
817 p.
818

819 Fekete, B.M, Robarts, R.D., Kumagai, M., Nachtnebel, H.-P., Odada, E. & Zhulidov, A.V.
820 (2015). Time for in situ renaissance. *Science*. 349(6249), 685–686.
821 <http://doi.org/10.1126/science.aac7358>.
822

823 Gelaro, R., & coauthors (2017). The Modern-Era Retrospective Analysis for Research and
824 Applications, Version-2 (MERRA-2), *J. Clim.* 30, 5419-5454,
825 <https://doi.org/10.1126/science.aac7358>.
826

827 Hrachowitz, M. & coauthors. (2013). A decade of Predictions in Ungauged Basins (PUB)—a
828 review. *Hydrological Sciences Journal*. 58(6), 1198-1255.
829 <https://doi.org/10.1080/02626667.2013.803183>.
830

831 Hutcheon, J.A., Chiolero, A. & Hanley, J.A. (2010). Random measurement error and regression
832 dilution bias. *BMJ*. 340, c2289. <https://doi.org/10.1136/bmj.c2289>.
833

834 Koster, R.D., Liu, Q., Mahanama, S.P.P. & Reichle, R.H. (2018). Improved hydrological
835 simulation using SMAP data: Relative impacts of model calibration and data assimilation.
836 *Journal of Hydrometeorology*. 19(4), 727-741. <https://doi.org/10.1175/JHM-D-17-0228.1>.
837

838 Koster, R.D. & Milly, P.C.D. (1997). The interplay between transpiration and runoff
839 formulations in land surface schemes used with atmospheric models. *J. Clim.* 10, 1578-1591.
840

841 Kumar, S.V., Peters-Lidard, C.D., Tian, Y., Houser, P.R., Geiger, J., Olden, S., Lighty, L.,
842 Eastman, J.L., Doty, B., Dirmeyer, P., Adams, J., Mitchell, K., Wood, E.G. & Sheffield, J.
843 (2006). Land Information System - An interoperable framework for high resolution land surface
844 modeling. *Environmental Modelling & Software.* 21, 1402-1415.
845 <https://doi.org/10.1016/j.envsoft.2005.07.004>.
846

847 Linsley, R.K., Kohler, M.A., & Paulhus, J.L.H. (1982). *Hydrology for Engineers*. 3rd Edition.
848 McGraw-Hill. New York. 508 p.
849

850 Lohmann, D. & coauthors (2004). Streamflow and water balance intercomparison
851 of four land surface models in the North American Land Data Assimilation System (NLDAS).
852 *J. Geophys. Res.* 109, D07S91, <https://doi.org/10.1029/2003JD003517>.
853

854 Lucchesi, R. (2018). File Specification for GEOS FP. GMAO Office Note No. 4 (Version 1.2),
855 61pp. [Available online at <http://gmao.gsfc.nasa.gov/pubs/>].
856

857 Niu, G.-Y. & coauthors (2011), The community Noah land surface model with
858 multiparameterization options (Noah-MP): 1. Model description and evaluation with local-scale
859 measurements, *J. Geophys. Res.* 116, D12109, <https://doi.org/10.1029/2010JD015139>.
860

861 Piepmeier, J. R. et al. (2017), SMAP L-Band Microwave Radiometer: Instrument Design and
862 First Year on Orbit, *IEEE Transactions on Geoscience and Remote Sensing.* 55(4), 1954–1966,
863 doi:10.1109/tgrs.2016.2631978.
864

865 Qiu, J., Dong, J., Crow, W.T., Zhang, X., Reichle, R. & De Lannoy, G. (2021). The added
866 value of brightness temperature assimilation for the SMAP Level-4 surface and root-zone soil
867 moisture analysis over mainland China. *Hydrology and Earth System Sciences.* 25, 1569-1586.
868 <https://doi.org/10.5194/hess-25-1569-2021>.
869

870 Reichle, R.H. & coauthors (2017), Global assessment of the SMAP Level-4 Surface and Root-
871 Zone Soil Moisture product using assimilation diagnostics. *Journal of Hydrometeorology.* 18,
872 3217-3237. <https://doi.org/10.1175/JHM-D-17-0130.1>.
873

874 Reichle, R.H. & coauthors. (2019). Version 4 of the SMAP Level-4 soil moisture algorithm and
875 data product. *Journal of Advances in Modeling Earth Systems.* 11, 3106-3130.
876 <https://doi.org/10.1029/2019MS001729>.
877

878 Reichle, R., De Lannoy, G., Koster, R.D., Crow, W.T., Kimball, J.S. & Liu, Q. (2020). *SMAP*
879 *L4 Global 3-hourly 9 km EASE-Grid Surface and Root Zone Soil Moisture Geophysical Data,*
880 *Version 5.* Boulder, Colorado USA. NASA National Snow and Ice Data Center Distributed
881 Active Archive Center. <https://doi.org/10.5067/9LNYIYOB5>. Accessed: 1 December
882 2020.

883

884 Reichle, R.H. & coauthors. (2021a). Soil Moisture Active Passive (SMAP) project assessment
885 report for version 5 of the L4_SM data product, R.D. Koster (Ed). *Technical Report Series on*
886 *Global Modeling and Data Assimilation.* NASA GMAO, 58.

887

888 Reichle, R. H. & coauthors (2021b). The contributions of gauge-based precipitation and SMAP
889 brightness temperature observations to the skill of the SMAP Level-4 soil moisture product.
890 *Journal of Hydrometeorology.* 22, 405-424, <https://doi.org/10.1175/JHM-D-20-0217.1>.

891

892 Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Eisner, S., Müller
893 Schmied, H., Sutanudjaja, E.H., Warrach-Sagi, K. & Attinger, S. (2017). Toward seamless
894 hydrologic predictions across spatial scales. *Hydrol. Earth Syst. Sci.* 21, 4323–4346.
895 <https://doi.org/10.5194/hess-21-4323-2017>.

896

897 Schaake, J.C., Duan, Q., Smith, M. & Koren, V. (2000). Criteria to select basins for hydrologic
898 model development and testing. Preprints in 15th Conf. On Hydrology (Long Beach,
899 California, USA, Am. Met. Soc., 10–14 January 2000), Paper P1.8.

900

901 Shellito, P.J., Kumar, S.V., Santanello, J.A., Jr., Lawston-Parker, P., Bolten, J.D., Cosh, M.H.,
902 Bosch, D.D., Holifield Collins, C.D., Livingston, S., Prueger, J., Seyfried, M. & Starks, P.J.
903 (2020). Assessing the impact of soil layer depth specification on the observability of modeled
904 soil moisture and brightness temperature. *Journal of Hydrometeorology.* 21(9), 2041-2060.
905 <https://doi.org/10.1175/JHM-D-19-0280.1>

906

907 Sloto, R.A. & Crouse. M.Y. (1996). HYSEP: A computer program for streamflow hydrograph
908 separation and analysis. *U.S. Geological Survey Water-Resources Investigations Report 1996–*
909 *4040*, 46 p., <https://pubs.er.usgs.gov/publication/wri964040>.

910

911 Song, S. & Wang, W. (2019). Impacts of antecedent soil moisture on the rainfall-runoff
912 transformation process based on high-resolution observations in soil tank experiments. *Water.*
913 *11(296).* <https://doi.org/10.3390/w11020296>.

914

915 United Nations Environment Programme. Middleton, N., & Thomas, D. S. G. (1992). *World*
916 *atlas of desertification.* London: Edward Arnold.

917

918 Villarini, G., Mandapaka, P.V., Krajewski, W.F. & Moore, R.J. (2008). Rainfall and sampling
919 uncertainties: A rain gauge perspective. *J. Geophys. Res.* 113, D11102,
920 <https://doi.org/10.1029/2007JD009214>.
921

922 Xia, Y. & coauthors. (2009). *NCEP/EMC, NLDAS Primary Forcing Data L4 Hourly 0.125 x*
923 *0.125 degree V002*, Edited by David Mocko, NASA/GSFC/HSL, Greenbelt, Maryland, USA,
924 Goddard Earth Sciences Data and Information Services Center (GES DISC), Accessed: 1
925 November 2020, <https://doi.org/10.5067/6J5LHHOHZHN4>
926

927 Xia, Y. & coauthors. (2012). Continental-scale water and energy flux analysis and validation
928 for North American Land Data Assimilation System project phase 2 (NLDAS-2): 2. Validation
929 of model-simulated streamflow. *J. Geophys. Res.*, 117, D03110,
930 <https://doi.org/10.1029/2011JD016051>.
931

932 Zhang, Y., Wei, H. & Nearing, M.A. (2011). Effects of antecedent soil moisture on runoff
933 modeling in small semiarid watersheds of southeastern Arizona. *Hydrol. Earth Syst. Sci.* 15,
934 3171–3179.
935

936 Zhao, R.-J., Zuang, Y.-L., Fang, L.-R., Liu, X.-R. & Zhang, Q.-S. (1980). The Xinanjiang
937 model. In *Proceedings of the Oxford Symposium*, Oxford, UK, 15–18 April 1980; pp. 351–356.

938 Zomer R.J., Bossio, D.A., Trabucco, A., Yuanjie, L., Gupta, D.C. & Singh, V.P. (2007). Trees
939 and water: Smallholder agroforestry on irrigated lands in northern India. *Colombo, Sri Lanka:*
940 *International Water Management Institute*. pp 45. (IWMI Research Report 122).

941 Zomer, R.J., Trabucco, A., Bossio, D.A., van Straaten, O. & Verchot, L.V. (2008). Climate
942 change mitigation: A spatial analysis of global land suitability for clean development
943 mechanism afforestation and reforestation. *Agric. Ecosystems and Envir.* 126, 67-80.

Figure 1.

Observed Soil Moisture Runoff Coupling Strength (SRCS)

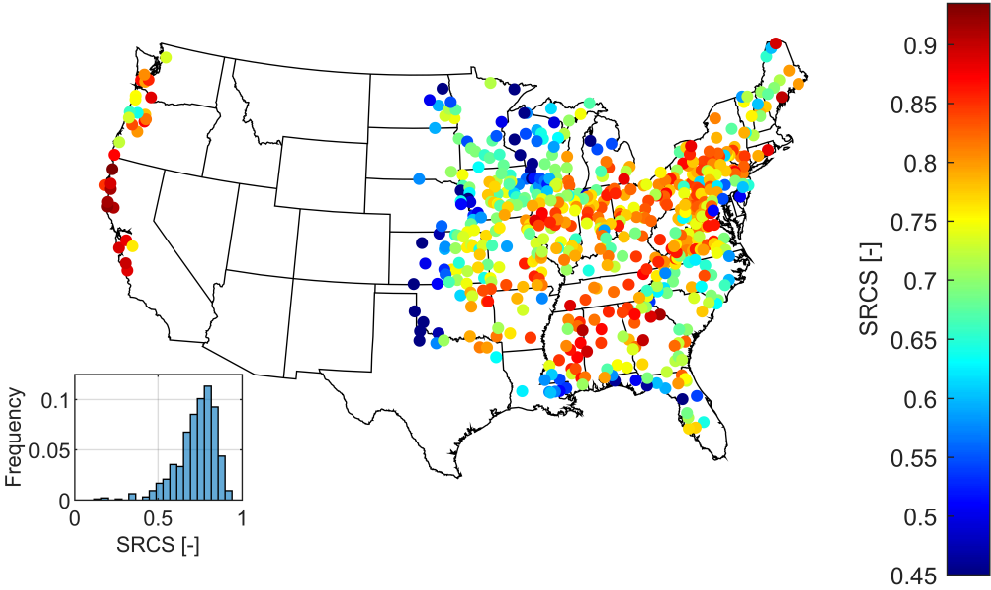


Figure 2.

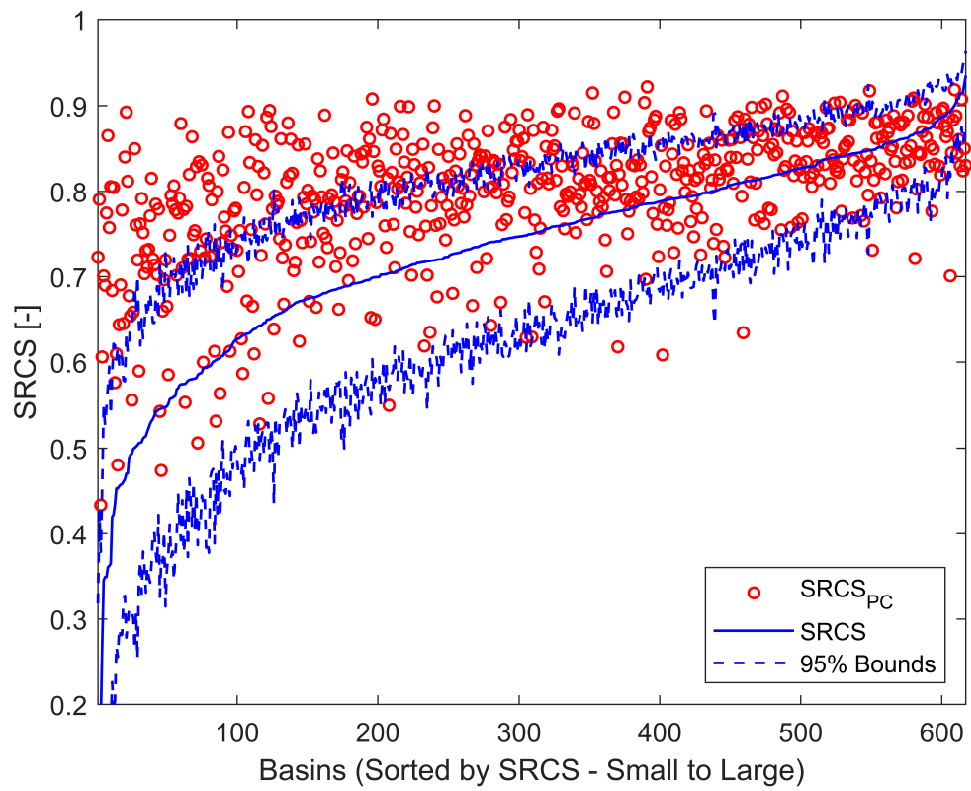


Figure 3a.

a) Δ SRCS [-]

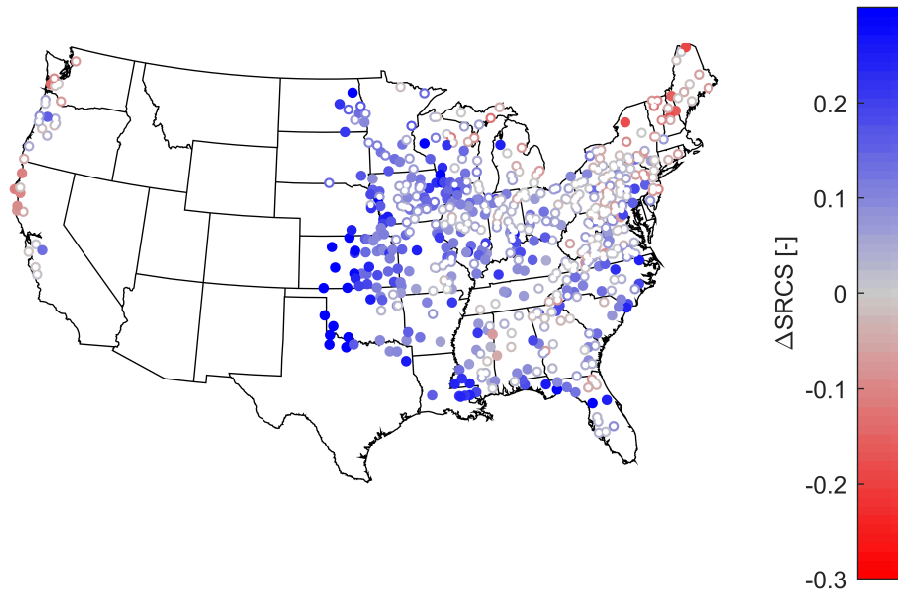


Figure 3b.

b) Aridity Index (AI)

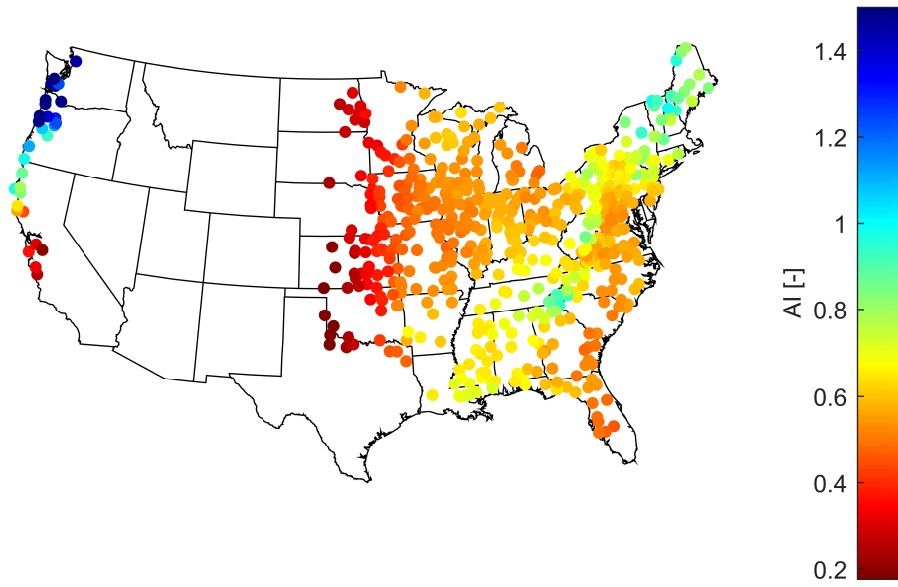


Figure 4.

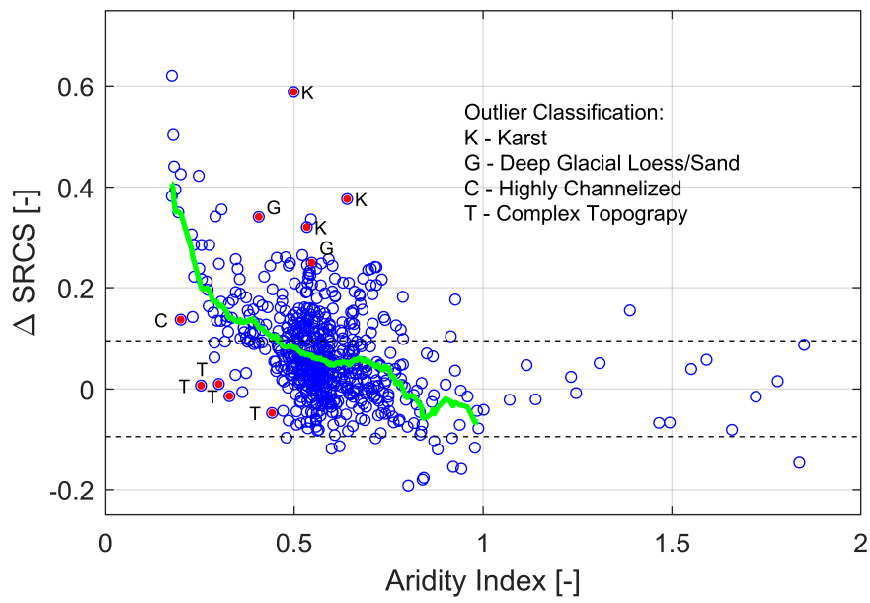


Figure 5.

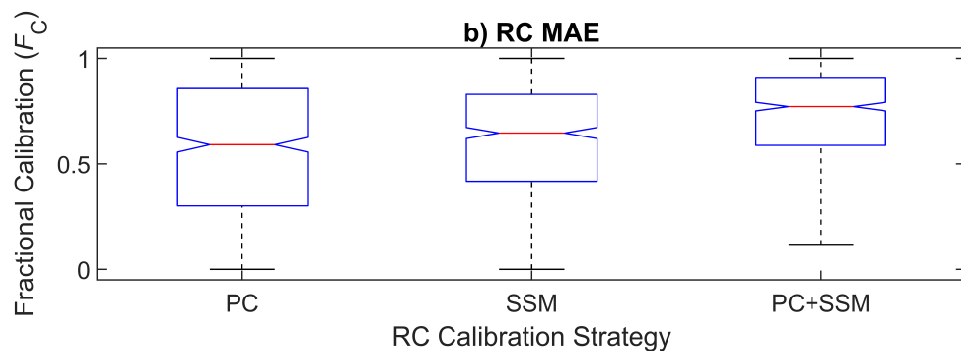
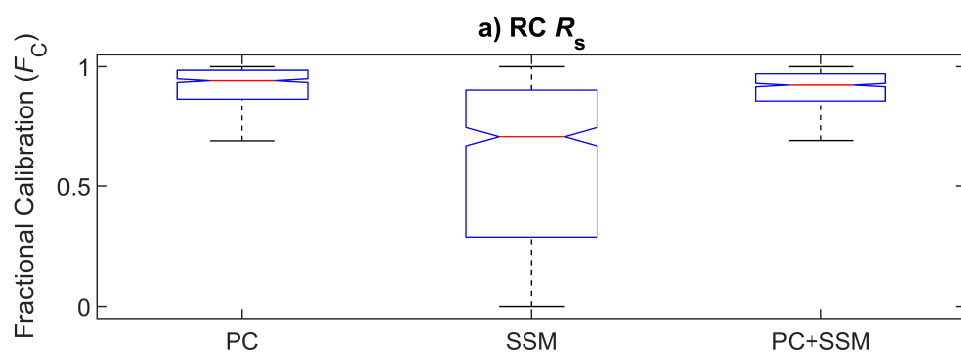


Figure 6a.

a) F_C for RC R_S

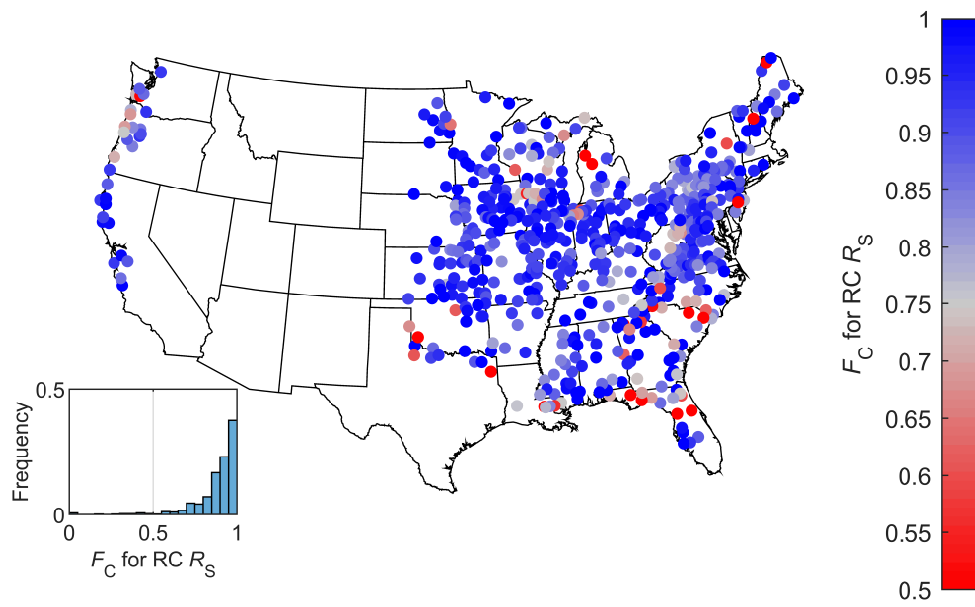


Figure 6b.

b) F_C for RC MAE

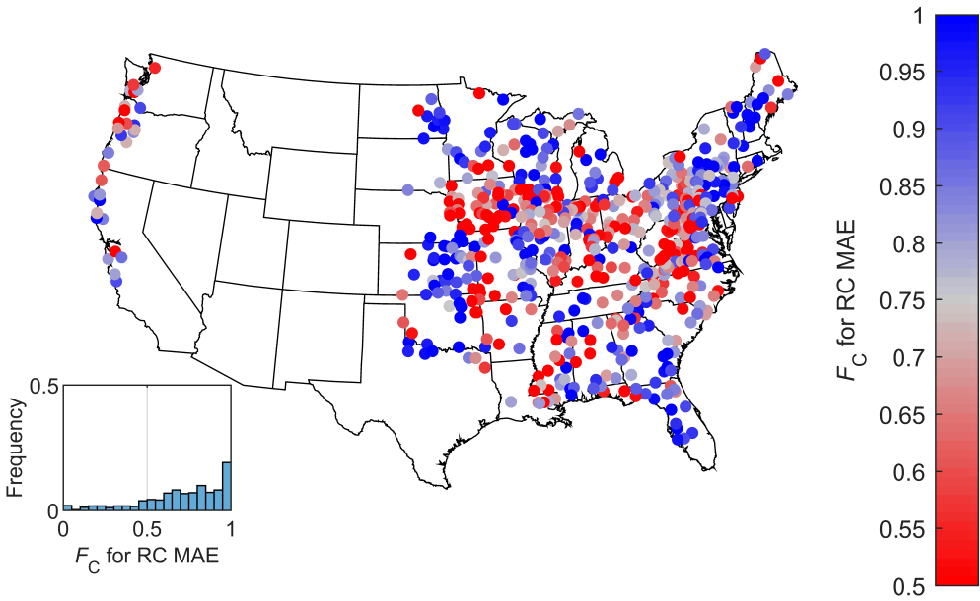


Figure 7.

PC+SSM (RC Bias)

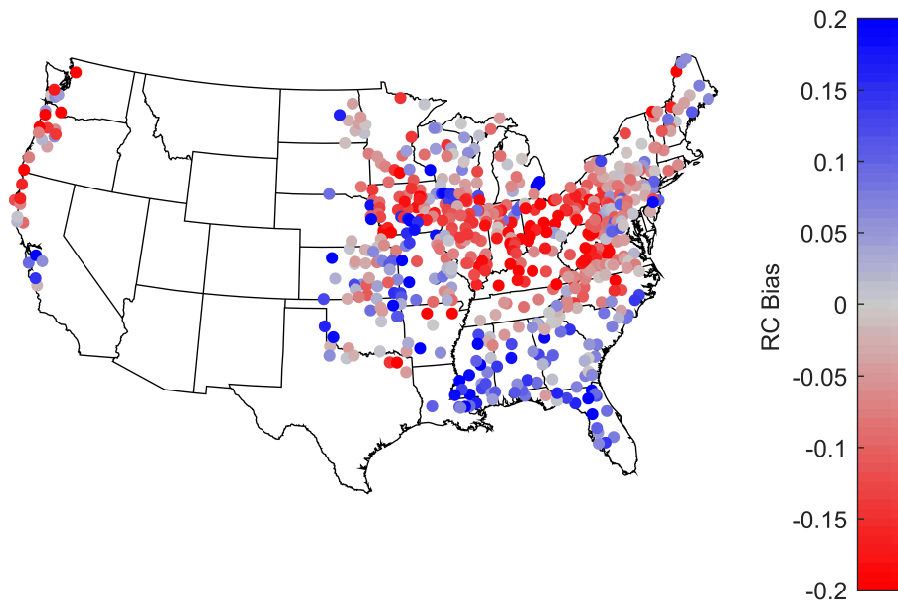


Figure 8.

