

1 Development of a “nature run” for observing system simulation experiments (OSSEs) for
2 snow mission development
3

4 Melissa L. Wrzesien^{1,2}, Sujay Kumar¹, Carrie Vuyovich¹, Ethan D. Gutmann³, Rhae Sung
5 Kim^{1,4}, Barton A. Forman⁵, Michael Durand⁶, Mark S. Raleigh⁷, Ryan Webb^{8,9}, Paul Houser¹⁰
6

7 ¹Hydrological Sciences Laboratory, NASA Goddard Space Flight Center, Greenbelt, MD

8 ²ESSIC, University of Maryland, College Park, MD

9 ³National Center for Atmospheric Research, Boulder, CO

10 ⁴GESTAR, Universities Space Research Association, Columbia, MD

11 ⁵Department of Civil and Environmental Engineering, University of Maryland, College Park,
12 MD

13 ⁶School of Earth Sciences and Byrd Polar and Climate Research Center, Ohio State
14 University, Columbus, OH

15 ⁷College of Earth, Ocean, and Atmospheric Sciences, Oregon State University, Corvallis, OR

16 ⁸Department of Civil, Construction, and Environmental Engineering, University of New
17 Mexico

18 ⁹Center for Water and the Environment, University of New Mexico, Albuquerque, NM
19 87131 USA

20 ¹⁰Department of Geography and Geoinformation Sciences, George Mason University,
21 Fairfax, VA
22

23
24 Corresponding author: Melissa Wrzesien, melissa.l.wrzesien@nasa.gov
25
26
27
28
29
30
31
32
33
34
35
36
37
38

39 **1. Abstract**

40 Snow is a fundamental component of global and regional water budgets, particularly
41 in mountainous areas and regions downstream that rely on snowmelt for water resources.
42 Land surface models (LSMs) are commonly used to develop spatially distributed estimates
43 of snow water equivalent (SWE) and runoff. However, LSMs are limited by uncertainties in
44 model physics and parameters, among other factors. In this study, we describe the use of
45 model calibration tools to improve snow simulations within the Noah-MP LSM as the first
46 step in an Observing System Simulation Experiment (OSSE). Noah-MP is calibrated against
47 the University of Arizona (UA) SWE product over a Western Colorado domain. With
48 spatially varying calibrated parameters, we run calibrated and default Noah-MP
49 simulations for water years 2010-2020. By evaluating both simulations against the UA
50 dataset, we show that calibration decreases domain averaged temporal RMSE and bias for
51 snow depth from 0.15 to 0.13 m and from -0.036 to -0.0023 m, respectively, and improves
52 the timing of snow ablation. Increased snow simulation performance also improves
53 estimates of model-simulated runoff in four of six study basins, though only one has
54 statistically significant improvement. Spatially distributed Noah-MP snow parameters
55 perform better than default uniform values. We demonstrate that calibrating variables
56 related to snow albedo calculations and rain-snow partitioning, among other processes, is a
57 necessary step for creating a nature run that reasonably approximates true snow
58 conditions for the OSSEs. Additionally, the inclusion of a snowfall scaling term can address
59 biases in precipitation from meteorological forcing datasets, further improving the utility of
60 LSMs for generating reliable spatiotemporal estimates of snow.

61

62 **2. Introduction**

63 Snow is a critical part of global and local water budgets, particularly in watersheds
64 with headwaters in mountainous regions (Viviroli et al., 2007; Immerzeel et al., 2020).
65 Millions of people around the world rely on snowmelt-derived runoff (Barnett et al., 2005;
66 Li et al., 2017), especially in semi-arid regions. Despite being an integral component of
67 global and regional water balances, estimating mountain snow accumulation remains one
68 of the largest challenges of snow hydrology (Bormann et al., 2018; Dozier et al., 2016).
69 While some mountain ranges have relatively dense in situ networks, other areas lack
70 observations (Dozier et al., 2016), limiting techniques for interpreting point observations
71 across a larger scale. Beyond in situ observations, remote sensing offers the ability to
72 observe snow extent from space (Hall et al., 2002), but estimating snow water equivalent
73 (SWE) to understand the water content of the snowpack remains a significant challenge,
74 particularly in the mountains (Lettenmaier et al., 2015; Nolin, 2010; Takala et al., 2011;
75 Vuyovich et al., 2014).

76 Due to limited in situ networks and uncertainty in remotely sensed observations,
77 models are a practical alternative for developing spatiotemporal estimates of snow depth
78 and SWE across large regions. Model intercomparison efforts have helped to identify
79 important processes to improve simulating snow (Essery et al., 2009; Etchevers et al.,
80 2004; van den Hurk et al., 2016; Krinner et al., 2018; Rutter et al., 2009), such as multi-
81 layer snowpack. While snow models often have complex physics and parameterizations,
82 resulting in accurate simulations of snow compared to in situ observations (Dutra et al.,
83 2012; Etchevers et al., 2004), such processes are often too computationally complex for
84 land surface models (LSMs) designed to run over large geographical areas. Additionally,

85 snow models are typically focused only on modeling the snowpack processes whereas
86 LSMs also enable the linkages to the water, energy, and carbon cycle processes. Though
87 LSMs allow for simulations across a range of spatial and temporal scales in a
88 computationally efficient manner, the relatively simple nature of their conceptual
89 formulations and model parameterizations, as compared to complex process models,
90 increases the uncertainties of their predictions. Further, biases in model forcing data,
91 particularly precipitation, are a major driver of model error (Raleigh et al., 2015; Schmucki
92 et al., 2014; Henn et al., 2018), and studies suggest that reanalyses, which are often used for
93 model meteorological forcing, underestimate precipitation in mountainous areas (Henn et
94 al., 2018; Enzinger et al., 2019; He et al., 2019). Such limitations are well documented in
95 the literature, where it has been suggested that common LSMs, such as the Noah LSM with
96 multiple parameterization options (Noah-MP; Niu et al., 2011), underestimate snow mass
97 (Chen, Liu et al., 2014; Kumar et al., 2019; Xia et al., 2017; Chen, Barlage, et al., 2014).
98 Despite these issues, LSMs are an essential tool for producing multi-year estimates of snow
99 accumulation over continental or global study domains.

100 To reduce biases, models are often calibrated against reliable observation-based
101 datasets (e.g. Ahl et al., 2008; Franz & Karsten, 2013; Henn et al., 2016; Rutter et al., 2009).
102 Calibration has a long history in operational snow modeling (e.g. Turcotte et al., 2017;
103 Franz et al., 2008) and previous intercomparison projects explicitly considered the
104 performance of calibrated vs. non-calibrated models (Rutter et al., 2009; Essery et al.,
105 2009). Often in snow and hydrological modeling, simulations are calibrated against
106 discharge for improving model performance (Franz and Karsten, 2013; Hay et al., 2006; Ahl
107 et al., 2008; Turcotte et al., 2017;). More recently, efforts have aimed to improve snow

108 estimation by calibrating against SWE (Chen et al., 2017; Franz et al., 2010), snow-covered
109 area (Franz and Karsten, 2013; Parajka and Bloschl, 2008), or multi-objective strategies
110 that include two or more calibration variables (Nemri and Kinnard, 2020; Parajka et al.,
111 2007; Chen et al., 2017; Franz and Karsten, 2013).

112 The performance of a calibrated model, however, will depend upon parameter
113 selection for use during calibration, and complex LSMs such as Noah-MP have hundreds of
114 parameters throughout the model code, some that are hard-coded to spatially uniform
115 values. Cuntz et al. (2016) examined over 100 Noah-MP parameters, dozens of which are
116 hard-coded into the LSM, and showed that simulated surface runoff is sensitive to almost
117 all selected snow parameters; the authors conclude that it is necessary to expose some of
118 the hard-coded parameters during calibration in order to improve model performance.
119 Similarly, Mendoza et al. (2015) discussed that hard-coding parameters diminishes model
120 agility; they identify several important hard-coded snow parameters that are treated as
121 spatially uniform constants but in actuality likely vary through both time and space.

122 Here we calibrate Noah-MP against SWE estimates from the University of Arizona
123 gridded observation-based snow data product (here referred to as UA; Zeng et al., 2018) in
124 an effort to address dry biases in Noah-MP and improve snow estimation. We evaluate the
125 impact of calibration on simulation of snow mass in a mountainous region. Since
126 calibration will have implications beyond snow-related variables, we also examine impacts
127 to other hydrologic processes, including runoff. The overarching motivation for the
128 calibration is to produce a Noah-MP simulation that better approximates snow conditions
129 through improvements to snow depth and SWE.

130 We aim for the calibrated simulation to be used as the “nature run” (NR) in a
131 forthcoming snow-focused Observing System Simulation Experiment (OSSE). OSSEs are
132 data assimilation experiments, performed to evaluate the type and impact of data to be
133 collected from proposed missions and to enable the assessment of the utility from
134 competing mission designs and design configurations (Garnaud et al., 2019; Crow et al.,
135 2001, 2005; Wang et al., 2008; Nearing et al., 2012). Further, these experiments help to
136 quantify the utility of observations beyond the immediate variable of interest (e.g., the
137 impact of assimilating snow information on other aspects of the water budget, such as
138 streamflow). OSSEs are useful in developing assessments of proposed observational
139 methods and can be performed in addition to field work, such as the extensive NASA
140 SnowEx campaigns, for evaluating proposed sensors.

141 A NR is the foundational step of an OSSE, upon which the data assimilation
142 experiments are built (see Figure S1 for general steps to an OSSE). Within an OSSE, the NR
143 simulation is considered the “true” state of the variable of interest. Therefore, NRs are
144 developed using a high-quality model and meteorological inputs and should not have large
145 uncertainty. Synthetic observations are then generated from the NR, after accounting for
146 sources of errors and uncertainty associated with the anticipated sensor. The synthetic
147 observations are assimilated into an open loop model simulation, and the assimilated
148 result is compared back to the original NR to understand how well the proposed sensor
149 captures the “true” conditions. The quality of the NR, therefore, significantly impacts the
150 conclusions made from the OSSE. Since previous studies highlight biases in LSMs related to
151 snow depth and SWE estimation, it is critical to reduce LSM bias and uncertainty to assess
152 how proposed technologies perform in a variety of environments. If the NR and resulting

153 synthetic observations are biased low, for example, it will be difficult to understand how a
154 proposed sensor observes deep snowpacks. While a NR is not expected to be a perfect
155 simulation, if it has a known systematic negative bias for SWE and snow depth, the
156 assimilation experiments may not provide much information for how a sensor performs in
157 regions where models have larger uncertainty, such as deep snow and forested regions
158 (Kim et al., 2021). The calibration procedure described below is the first and an essential
159 step in an OSSE designed to test potential configurations for a snow mission.

160 In addition to producing an improved NR for the OSSE, we aim to address three
161 research questions: 1.) Can calibration address known dry biases in LSMs that cause
162 underestimation of snow accumulation? 2.) How does calibration impact streamflow,
163 beyond the targeted snow variables? and 3.) Can calibration suggest areas of model
164 configuration that need improvement, such as meteorological forcing data for use as model
165 boundary conditions? We test whether Noah-MP with calibration pre-processing yields
166 similar snow estimates as a higher resolution, computationally expensive and complex
167 snow physics model (SnowModel). We introduce the study area and calibration procedure
168 in Section 3 below. In Section 4, we report results from the calibration experiments, and in
169 Section 5, we discuss implications and provide thoughts for future studies.

170

171 **3. Data and Methods**

172 **3.1. Model Setup**

173 We use the NASA Land Information System (LIS; Kumar et al., 2006; Peters-Lidard
174 et al., 2007) for simulations over a western Colorado domain. The domain is selected to
175 include sites from previous NASA SnowEx field campaign locations, including Grand Mesa

176 and Senator Beck (Figure 1). LIS is a land surface modeling framework designed to be
177 highly flexible, offering users choice of LSM, meteorological forcing, and assimilation of in
178 situ and remotely sensed observations, among other options. Created to be
179 computationally efficient, LIS can perform simulations over large regional and global
180 domains. The central component of the LIS framework is the LSM selection; LIS offers
181 several community-supported LSMs relevant to operations and research. Here we use
182 Noah-MP version 4.0.1. Recent work demonstrates that Noah-MP has superior
183 performance to the original Noah LSM for simulating snow (Chen, Barlage et al., 2014;
184 Chen, Liu et al., 2014; Kim et al., 2021; Minder et al., 2016; Wrzesien et al., 2015) due to
185 model physics updates, including a multilayer (three layer) snowpack. Table S1 lists the
186 physics options selected for the Noah-MP simulation.

187 In the LIS framework, Noah-MP simulates both surface water and energy fluxes as
188 they respond to meteorological boundary conditions supplied by LIS. Simulations are from
189 September 2009 through July 2020 at 0.01° spatial resolution (~ 1 km) and use hourly
190 meteorological forcing data from the North American Land Data Assimilation System phase
191 2 (NLDAS-2; Xia et al., 2012). LIS includes statistical downscaling procedures for matching
192 meteorological data to the specified spatial resolution of the LSM. The $1/8^\circ$ spatial
193 resolution NLDAS-2 forcing data are downscaled to ~ 1 km through a bilinear spatial
194 interpolation approach. The model was first spun up for 72 years beginning in January
195 1979 and running through January 2020 twice until the simulation begins in September
196 2009. We also simulate the same time period using the default parameters to understand
197 how calibration impacts the Noah-MP results. We distinguish between the two simulations

198 as Noah-MP-Cal and Noah-MP-Def to represent the calibrated and default configurations,
199 respectively.

200

201 3.2. Noah-MP Parameter Calibration

202 Previous studies suggest that LSMs underestimate snow accumulation, particularly
203 in mountains (Broxton, Zeng, et al., 2016; Wrzesien et al., 2017, 2018). A recent model
204 intercomparison using an ensemble of LSM simulations from LIS highlighted the model
205 disagreement and uncertainty of snow estimation over North America, including mountain
206 areas (Kim et al., 2021). To improve Noah-MP simulations, we select 24 parameters for
207 calibration (Table 1), based on previous sensitivity studies (Cuntz et al., 2016; Mendoza et
208 al., 2015) and their relationship to modeled snow processes. In Noah-MP-Def, these
209 parameters are either hard-coded, often to a single spatially uniform value, or provided in
210 lookup tables that vary based on land or soil properties. In contrast, the results from
211 calibration are spatially distributed parameters that can vary across the domain (Figure 2).
212 In addition to 23 existing parameters within Noah-MP, we include a snowfall scale factor in
213 the calibration. Precipitation underestimation will impact the snow simulation and lead to
214 biases throughout the snow season. The inclusion of a snowfall scale factor allows us to
215 target the uncertainty resulting from biases in precipitation forcing. All 24 parameters are
216 explored in point scale and full domain tests, though only the parameters that are sensitive
217 enough to warrant calibration are described in Section 4.1.

218 Noah-MP is calibrated against SWE estimates from the University of Arizona dataset
219 (UA; Zeng et al., 2018) in an optimization approach. The UA data product provides SWE at 4
220 km spatial resolution over the conterminous United States (Zeng et al., 2018). Estimates

221 are provided daily between 1981 and 2020. UA is based on the assimilation of in situ
222 measurements of both SWE and snow depth (Broxton, Dawson, et al., 2016) and
223 precipitation and temperature values from the PRISM dataset (Daly et al., 2000). UA has
224 been evaluated against multiple datasets (Dawson et al., 2018), including airborne lidar
225 measurements of snow depth. We note that any biases in UA SWE will likely be reflected in
226 the calibrated parameters and the resulting simulations; however, such biases, especially in
227 gridded observation-based data products like UA, are unavoidable.

228 We calibrate over water years 2007-2009. This period was selected by examining
229 domain averaged SWE from water years 1982-2020 from the UA record. From
230 comparisons of domain-wide average maximum SWE and depth, this period included
231 average (2009), high (2008), and low (2007) snow conditions for the study region.
232 Domain-wide average maximum SWE (snow depth) for water years 2007, 2008, and 2009
233 is 135.5 mm, 231.0 mm, and 162.4 mm, respectively, (524.0 mm, 850.1 mm, and 536.9 mm,
234 respectively) vs. the long-term mean of 163.5 mm (618.5 mm).

235 For calibration, we use a genetic algorithm (GA), which is part of the LIS-
236 Optimization and Uncertainty subsystem (Kumar et al., 2012). The GA is a common
237 stochastic tool used in hydrology model optimization (Duethmann et al., 2014; Isenstein et
238 al., 2015; Shafii & De Smedt, 2009; Wang, 1991; Yapo et al., 1998) and is designed to mimic
239 biological evolution where the fittest of the population (i.e., parameter sets), as determined
240 through comparison to an observational dataset, survive and move to the next generation.
241 Within each generation, crossover and mutation operators are used to produce new
242 parameter estimates and to introduce diversity in the parameter set. To ensure good
243 solutions are not lost between generations due to either crossover or mutation operators,

244 an elitism strategy is used, where the best solution is carried over to the next generation.
245 Over many generations, the average fitness, which reflects the quality of the solution, tends
246 to increase due to the selection of individuals that compare favorably to observations.

247 GAs aim to prevent overfitting through an ensemble approach and by introducing
248 poor performing solutions through mutation operators. Since they do not rely on gradient
249 information, GAs can handle local optima and discontinuities in the search space, unlike
250 gradient search. Since GAs require an ensemble that must be run over several generations,
251 they are computationally expensive. Running 50 generations of the GA with 30 ensemble
252 members for three water years over the study domain requires a total running time over
253 480 hours, or over 20 days of continuous simulation, with 532 processors.

254 Within LIS, the GA does not provide estimates of parameter uncertainty. For
255 estimating parameter uncertainty, variants of Markov Chain Monte Carlo methods such as
256 Differential Evolution Monte Carlo (te Braak et al., 2008) would be required; however,
257 algorithms such as these have a high computational cost, with run times an order of
258 magnitude higher than GA (Harrison et al., 2012), making their implementation over a
259 domain size such as ours difficult. Since the primary objective of this study is to produce a
260 better snow simulation, a thorough investigation into the parameter uncertainty is omitted.
261 More detail on GAs within the LIS framework is discussed by Kumar et al. (2012).

262 The GA results in calibrated values for a set of parameters that allow for the best
263 match with observations. The range in parameter values for calibration (see Table 1) are
264 either taken from the literature or allowed to vary +/- 20% of the default value, following
265 Cuntz et al. (2016). As an objective function, we consider the squared difference between
266 the observation and the model:

267
$$J_i = (d_i^o - d_i^m)^2 \quad (1)$$

268 where d_i^o is snow depth from the observations (UA) for grid cell i and d_i^m is snow depth
269 from the model (Noah-MP) for grid cell i . We minimize J_i for each grid cell i independently
270 in the calibration, resulting in parameters that vary spatially (Figure 2). In contrast, Noah-
271 MP-Def has spatially uniform parameters. UA, produced at 4 km, is rescaled to match the
272 Noah-MP resolution through bilinear interpolation during calibration.

273

274 3.3. Evaluation Datasets

275 In addition to comparing Noah-MP estimates to UA, we evaluate snow simulations
276 against a suite of independent datasets using the Land surface Verification Toolkit (LVT;
277 Kumar et al., 2012). First, we compare snow depth across the full domain to the Snow Data
278 Assimilation System (SNODAS; Carroll et al., 2001), which is an operational dataset
279 available over the contiguous United States at approximately 1 km spatial resolution. Both
280 Noah-MP simulations are evaluated against UA and SNODAS for the full analysis period of
281 water years 2010-2020. UA and SNODAS are both reprocessed in LVT to match the spatial
282 resolution of Noah-MP.

283 We also compare to snow depth measurements from the Global Historical
284 Climatology Network (GHCN; Menne et al., 2012); the western Colorado domain includes
285 79 GHCN stations with snow depth observations. Stations within the domain include a
286 range of elevations (1467-3422 m) with an average station elevation of 2349 m. This
287 compares to the full Noah-MP domain with elevations ranging from 1399-4185 m and an
288 average elevation of 2639 m; approximately 9% of GHCN stations within the domain have
289 elevations > 3000 m, compared to 26% of the full domain. While GHCN stations

290 undersample higher elevations within the western Colorado domain, they provide an
291 additional evaluation dataset for snow depth. GHCN data are available for water years
292 2010-2016.

293 We also compare Noah-MP to datasets collected from the 2017 NASA SnowEx field
294 campaign in Colorado. First, we evaluate Noah-MP against snow pit observations of snow
295 depth and SWE from SnowEx (Elder et al., 2018) at Grand Mesa and Senator Beck, which
296 were collected between February 6-25, 2017. For a spatial comparison, we evaluate Noah-
297 MP snow depth against Airborne Snow Observatory (ASO) lidar observations of snow
298 depth, which are produced at 3 m spatial resolution (Painter, 2018). Here we use ASO
299 flights over Grand Mesa from February 8 and February 16; though other flights are
300 available for the 2017 field campaign, other days either included artefacts from the lidar
301 collection or excluded portions of the mesa.

302 In addition to observations from SnowEx, Noah-MP is evaluated against a
303 SnowModel simulation over Grand Mesa for the 2017 campaign, as described in Webb et al.
304 (2020). SnowModel is a widely used snow model that simulates distributed snow
305 properties in space and time and can be configured to simulate a single or multi-layer
306 snowpack (Liston and Elder, 2006a; Liston and Sturm, 1998). SnowModel is designed to
307 include four interconnected models: MicroMet for processing and downscaling
308 meteorological forcing data (Liston and Elder, 2006b), EnBal for calculating the energy
309 balance of the snowpack, SnowPack for simulating the snowpack in space and time, and
310 SnowTran-3D for computing redistribution of snow due to wind (Liston and Sturm, 1998;
311 Liston et al., 2007). Webb et al. (2020) configure SnowModel to simulate a single layer
312 snowpack over Grand Mesa for the 2016-2017 water year to coincide with the SnowEx

313 field campaign in February 2017. They use station observations as meteorological forcing
 314 data, including data from the Grand Mesa Study Plot (Skiles, 2018), four SnowEx campaign
 315 weather stations, and three nearby Snow Telemetry (SNOTEL) sites. SNOTEL sites provide
 316 temperature and precipitation observations, and all other stations provide temperature,
 317 wind speed/direction, humidity, and radiation. No adjustment of precipitation or other
 318 forcing data were made, and SnowModel simulations were independent of any snow
 319 observations. Elevation data were from the 1/3 arc-second USGS National Elevation
 320 Dataset, while vegetation data were taken from 30m USGS LANDFIRE v.1.4 Existing
 321 Vegetation Type data (Rollins, 2009) and reclassified to SnowModel vegetation types.
 322 Webb et al. (2020) ran SnowModel at multiple spatial resolutions, but here we consider
 323 SWE and snow depth outputs from their 30 m simulation. Webb et al. (2020) provide
 324 additional information on the SnowModel configuration and evaluation.

325 For spatial evaluations against both ASO and SnowModel, we calculate the Spatial
 326 Efficiency (SPAEF; Koch et al., 2018; Demirel et al., 2018), which combines histogram
 327 matching, spatial correlation coefficient, and spatial variability error to evaluate spatial
 328 patterns. SPAEF is defined as:

$$329 \quad SPAEF = 1 - \sqrt{(\alpha - 1)^2 + (\beta - 1)^2 + (\gamma - 1)^2} \quad (2)$$

330 where $\alpha = \rho(obs, mod)$, $\beta = \frac{\frac{\sigma_{mod}}{\mu_{mod}}}{\frac{\sigma_{obs}}{\mu_{obs}}}$, and $\gamma = \frac{\sum_{j=1}^n \min(K_j, L_j)}{\sum_{j=1}^n K_j}$. Here α is the Pearson correlation
 331 coefficient between the observation (ASO lidar or SnowModel simulation) and the model
 332 (Noah-MP), β is the fraction of the coefficient of variation, which represents spatial
 333 variability, and γ is the histogram intersection for the histogram of the observation, K , and

334 the histogram from the model, L (Swain and Ballard, 1991). SPAEF has an optimal value of
335 1.

336 For streamflow, we compare to natural flow estimates for four basins in the Upper
337 Colorado River Basin (UCRB) that lie completely within the model domain (see Table 6).
338 Natural flow estimates are from the Bureau of Reclamation and are available monthly
339 between 1901 and 2018 (Prairie and Callejo, 2005). We also compare to daily, unregulated
340 streamflow for two basins from the Catchment Attributes and Meteorology for Large
341 Sample Studies (CAMELS; Newman et al., 2015; Newman et al., 2014) dataset. We only use
342 streamflow observations between 2009 and 2014 for the two CAMELS basins, and the daily
343 streamflow has been processed into monthly averages. Since Noah-MP does not include
344 human management on streamflow networks, we cannot compare model-simulated runoff
345 to streamgauge observations, due to water diversions, dams, and other water management
346 practices. Instead, we compare monthly grid-cell generated runoff – the summation of
347 surface runoff and subsurface runoff – to monthly observations over small unmanaged
348 basins and to estimated natural flow (i.e., runoff in the absence of human management) in
349 larger basins. Using total runoff at monthly scales as a proxy to streamflow is a valid
350 assumption (Chow, 1964) and a strategy used in other studies (e.g., Koster et al., 2010). We
351 evaluate monthly streamflow with Nash-Sutcliffe Efficiency metrics (NSE; Nash & Sutcliffe,
352 1970), where a perfect fit with observations has $NSE = 1$, and $NSE > 0$ indicates the model
353 has better predictive skill than the mean of the observations.

354

355 **4. Results**

356 4.1. Calibration

357 We initially run point-scale calibration tests with 23 selected parameters from the
358 snow modules within Noah-MP (Table 1). Noah-MP-Cal generally improved the snow
359 ablation timing in spring months relative to Noah-MP-Def. However, maximum snow
360 conditions remained largely underestimated, particularly for sites with deep snowpack
361 (not shown). After implementation of a snowfall scaling factor, described below in equation
362 5, as an additional calibration parameter, test simulations resulted in snow depths in better
363 agreement with UA estimates. Therefore, for calibration over the full domain, we include
364 24 spatially variable parameters: 23 from Noah-MP and an additional snowfall scale term
365 (Figure 2).

366 Though we include 24 parameters in the GA procedure, only 11 were sensitive to
367 calibration. We determine that 13 are not sensitive because they do not demonstrate any
368 noticeable spatial patterns such as those reported in Figure 2 and instead calibrated values
369 have noisy spatial patterns (see Supplemental Figure S2). Some of the 11 selected
370 parameters have regions of noisy artificial patterns in regions of the domain that were
371 insensitive to calibration, often in portions of the domain where less snow accumulates
372 (Figure 2). Despite these regions, we look further into the 11 sensitive parameters. The first
373 four parameters are used within the CLASS snow albedo scheme (Verseghy, 1991) and
374 include minimum snow albedo (*MNSNALB*), maximum snow albedo (*MXSNALB*), the
375 exponent in the snow albedo decay relationship (*SNDECAYEXP*), and the new snow mass
376 required to cover old snow (*SWEMX*). These parameters are used in each time step to
377 calculate snow albedo. First, the albedo of the snow cover for the new time step is
378 determined as:

379
$$\alpha_s(t) = MNSNALB + [\alpha_s(t - 1) - MNSNALB] * \exp\left[-\frac{SNDECAYEXP * \Delta t}{3600}\right] \quad (3)$$

380 where α_s is snow albedo at time step t or $t - 1$ and Δt is the model time step. If new snow
381 has fallen in an amount larger than $SWEMX$, snow albedo is refreshed to a value of
382 $MXSNALB$.

383 The next group of calibration parameters relates to the rain-snow partitioning
384 scheme used here, i.e., the Jordan (1991) scheme from the SNTHERM model. In this
385 method, if air temperature is above the upper temperature limit ($T_U LIMIT$), all
386 precipitation is rainfall. At air temperatures below the lower temperature limit ($T_L LIMIT$),
387 all precipitation is snowfall. For temperatures between $T_L LIMIT$ and a middle threshold
388 ($T_M LIMIT$), the fraction of precipitation that is frozen is a function of air temperature. At
389 temperatures between $T_M LIMIT$ and $T_U LIMIT$, the fraction of precipitation that is frozen is
390 set to 0.6. In the calibration procedure, $T_L LIMIT < T_M LIMIT < T_U LIMIT$.

391 The remaining four parameters are from different schemes throughout Noah-MP,
392 and three were highlighted by Mendoza et al. (2015) as key parameters for model
393 sensitivity. These include the exponent used in the snow depletion curve ($MFSNO$), liquid
394 water holding capacity (SSI), and snow surface roughness length ($Z0SNO$). $MFSNO$ is used
395 within Noah-MP to calculate the fractional portion of the grid cell that is snow covered, as
396 shown in equation 4 below from Niu and Yang (2007):

$$397 \quad f_{sno} = \tanh \frac{h_{sno}}{2.5z_{0g} \left(\frac{\rho_{sno}}{\rho_{new}} \right)^{MFSNO}} \quad (4)$$

398 where h_{sno} is snow depth, z_{0g} is the bare soil roughness length, ρ_{sno} is bulk density of
399 snow, and ρ_{new} is the density of new snow, which is set to 100 kg/m³. f_{sno} is used
400 throughout Noah-MP to scale grid-cell calculations into snow-covered and non-snow-
401 covered fractions, including within surface radiation calculations.

402 *SSI* and *ZOSNO* are each used only once in the Noah-MP code. *SSI* is included in the
403 calculation of snow layer liquid water, which determines the rate of exfiltration of
404 snowmelt release from the bottom of the snowpack. *ZOSNO* is used to calculate the surface
405 roughness length for turbulent flux calculations over snow covered ground.

406 The final calibration parameter is the snowfall scaling term, *SNOWF_SCALEF*,
407 which was included to address uncertainty in precipitation forcing data. *SNOWF_SCALEF*
408 is described as:

$$409 \quad S = P * f_{ice} * SNOWF_SCALEF \quad (5)$$

410 where S is snowfall, P is total precipitation, and f_{ice} is the fraction of the precipitation that
411 is frozen. The snowfall scale factor is applied to frozen precipitation to reduce the bias
412 introduced from NLDAS-2. Other studies introduce a similar precipitation scaling factor in
413 optimization or assimilation experiments. Smyth et al. (2020), who also used NLDAS-2 for
414 model forcing data, use a snowfall correction factor to scale precipitation at their SNOTEL
415 study sites across the western United States. In their work, the average snowfall correction
416 factor is 1.64, indicating NLDAS-2 underestimates mountain snowfall by more than 50%. In
417 Smyth et al. (2020) and here, NLDAS-2 snowfall is too low and must be scaled to larger
418 values to produce realistic snow accumulation. Other studies have also included a
419 correction factor to address biases in snowfall from meteorological data (Magnusson et al.,
420 2017; He et al., 2011; Franz and Karsten, 2013). Errors in forcing data, particularly
421 precipitation, have a large impact on snow modeling performance (Raleigh et al., 2015;
422 Schmucki et al., 2014; Henn et al., 2018), and including a snowfall scaling term in the
423 calibration procedure can help address this bias.

424

425 4.2. SWE and Snow Depth Evaluation

426 4.2.1. UA and SNODAS Comparisons

427 4.2.1.1. Full Domain Comparison

428 Figure 3 shows the time series of average SWE and average snow depth across the
429 domain for Noah-MP-Cal, Noah-MP-Def, and the UA dataset. In nearly all cases, calibration
430 results in more snow and later snowmelt. Occasionally, Noah-MP-Cal produces more snow
431 accumulation than the UA dataset, such as in 2015 and 2017 (Figure 3). Over the 11-year
432 simulation, Noah-MP-Cal has larger magnitudes of snow depth and SWE; average maximum
433 SWE (depth) from Noah-MP-Cal is 166.7 mm (0.61 m), while average maximum SWE
434 (depth) from Noah-MP-Def is 131.8 mm (0.52 m).

435 Spatially, Noah-MP-Cal produces greater April 1 SWE at higher elevations across the
436 domain, averaged over the water year 2010-2020 simulation period (Figure 4). Estimates
437 from Noah-MP-Def have similar domain-wide averages as Noah-MP-Cal (Figure 3), but the
438 snow is less spatially variable. This is contrasted with Noah-MP-Cal where snow
439 accumulation more closely follows local topography. We also compare Noah-MP-Cal and
440 Noah-MP-Def to UA and SNODAS at six evaluation points throughout the domain that
441 correspond to SnowEx field campaign sites (Table 2). At these points, Noah-MP-Cal
442 generally has smaller biases and RMSE than Noah-MP-Def for the UA comparison (Table 3).
443 Noah-MP-Cal also tends to perform better than Noah-MP-Def when evaluated against
444 SNODAS (Table 3). Noah-MP-Cal has smaller bias and RMSE at all evaluation points except
445 Fool Creek and Senator Beck, the two highest elevations stations. For a similar comparison
446 but for SWE, see Table S2.

447 To compare Noah-MP-Cal and Noah-MP-Def against UA and SNODAS over the full
448 domain, we first calculate the SPAEF (Equation 2) to evaluate spatial performance.
449 Compared to UA, Noah-MP-Cal has a SPAEF of 0.799 and Noah-MP-Def has a SPAEF of
450 0.508. For SNODAS, Noah-MP-Cal also has a higher SPAEF metric: 0.722 vs 0.460 for Noah-
451 MP-Def. For RMSE (Figure 5), higher elevations tend to have larger RMSE values,
452 particularly for Noah-MP-Def compared to both SNODAS and UA. Noah-MP-Cal has high
453 RMSE values in the central northern portion of the study domain. This area has much larger
454 values of snow depth in Noah-MP-Cal than Noah-MP-Def, and the snowfall scale factor from
455 calibration is high in the area (up to 2.5-3, compared to the domain average of 1.16),
456 leading to increased precipitation and higher snow accumulations (discussed in Section 5).
457 Aside from this anomalous region and an area in the southern portion of the domain, Noah-
458 MP-Cal generally reduces the UA snow depth RMSE (Figure 5c), particularly at higher
459 elevations. Averaged over the domain, Noah-MP-Cal has a slightly lower RMSE (0.13 m)
460 than Noah-MP-Def (0.15 m) compared to UA (Table 3). Performance between Noah-MP-Cal
461 and Noah-MP-Def is similar for SNODAS as for the UA comparison. Averaged over the full
462 domain, Noah-MP-Def is in better agreement with SNODAS (RMSE of 0.18 m) than Noah-
463 MP-Cal (RMSE of 0.19 m), though results are generally similar.

464 Similar to RMSE, we also compare temporal bias over the full domain (Figure 6).
465 Noah-MP-Def has a negative bias for higher elevation grid cells compared to both UA and
466 SNODAS. This suggests that Noah-MP-Def is underestimating snow accumulation in the
467 mountains, highlighting the known dry bias of LSMs (e.g., Chen, Liu et al., 2014; Holtzman
468 et al., 2020; Kumar et al., 2019; Wang et al., 2019; Xia et al., 2017). Noah-MP-Cal bias spatial
469 patterns are similar between both UA and SNODAS, with a large positive bias in the central

470 northern portion of the domain due to anomalously high values of snow depth. Averaged
471 over the full domain, Noah-MP-Cal vs. UA has a bias of nearly zero (-0.0023 m), compared
472 to Noah-MP-Def of -0.036 m (Table 3). For both UA and SNODAS comparisons, Noah-MP-
473 Cal has more instances of positive bias at higher elevations (>3500 m), while these same
474 grid cells in Noah-MP-Def tend to have negative biases. Noah-MP-Def underestimates snow
475 accumulation at high elevations and calibration somewhat addresses these biases, though
476 can result in too much snow in some regions.

477 *4.2.1.2. Seasonal Comparison*

478 During the accumulation season (December through February), calibration
479 increases the domain averaged snow depth by almost 18%, from a -14.0% difference with
480 Noah-MP-Def to a +1.4% difference with Noah-MP-Cal, relative to UA. RMSE also improves
481 slightly from 0.162 m to 0.142 m. Similarly, for the peak snow season (March and April),
482 calibration results in an improvement of snow depth percent difference from -24.8%
483 (Noah-MP-Def) to -5.1% (Noah-MP-Cal). RMSE decreases from 0.269 m with Noah-MP-Def
484 to 0.215 m with Noah-MP-Cal, a 20% improvement. Calibration results in large
485 improvements for the ablation season (May through July), increasing the domain averaged
486 snow depth by 45.4%. Noah-MP-Def mean snow depth is 31.6% less than the UA estimate,
487 while Noah-MP-Cal is comparable to UA, only -0.5% smaller. RMSE decreases by over 12%,
488 from 0.0981 m with Noah-MP-Def to 0.0863 m with Noah-MP-Cal. Across the full domain,
489 calibration addresses the underestimation of snow throughout the full water year, though
490 with slightly too much snow during the peak snow season.

491 At the grid cell scale, Noah-MP-Cal generally has more snow accumulation and a
492 later end to the snow season than Noah-MP-Def, as shown in Figure 7 for Senator Beck.

493 Point scale evaluations have a better agreement between UA and Noah-MP-Cal, with RMSE
494 declining by 4.23 cm for peak season. During the accumulation and ablation seasons,
495 results are different, where Noah-MP-Def has smaller bias and RMSE. Noah-MP-Cal
496 overestimates UA in the spring for several years (Figure 7a), with snow lingering longer
497 than observed in UA for water years 2015, 2017, and 2019. Performance is similarly mixed
498 at other study points (Table 4), where calibration may improve performance during all
499 seasons (Cameron Pass, Niwot Ridge, Skyway/Grand Mesa) or may degrade performance,
500 depending on the season (accumulation and peak for Fool Creek, ablation for Rock Creek,
501 and accumulation and ablation for Senator Beck). Comparing SWE bias and RMSE over
502 different seasons has similar results (Table S3).

503 *4.2.1.3. Comparison over Vegetation Class*

504 We next aggregate the 20 LIS land cover classifications into five broader groups –
505 forest, shrubland, grassland, cropland, and barren (see inset in Figure 1 and see Table S4
506 for statistics) – and compare Noah-MP-Cal vs. Noah-MP-Def against both UA and SNODAS
507 (Figure 8a,c). For average snow depth bias, Noah-MP-Cal performs better than Noah-MP-
508 Def across land covers. Most comparisons have a negative bias, indicating that Noah-MP-
509 Cal and Noah-MP-Def have less snow than either UA or SNODAS, though magnitude of the
510 bias is generally smaller than 0.05 m. The exception is for the barren land cover class,
511 which is the category with the fewest grid cells (1564 or 1.4% of the domain) and the land
512 class with the highest average elevation (3178 m vs. a domain average of 2639 m).
513 Comparing with SNODAS, Noah-MP-Def has smaller RMSE than Noah-MP-Cal. In all classes
514 except cropland, Noah-MP-Cal has a lower RMSE when compared to UA.

515 Modeling in forested regions is often challenging due to uncertainty in snow-canopy
516 interactions (Essery et al., 2009; Krinner et al., 2018). Therefore, we further subdivide the
517 forest class into elevation bands to single out the impact of elevation on a land cover class
518 with higher uncertainty (Figure 8b,d and see Table S4 for number of grid cells within each
519 category). Results are similar to the full land cover comparison, where Noah-MP biases are
520 negative, and Noah-MP-Cal has smaller bias and RMSE than Noah-MP-Def. Higher
521 elevations have larger biases and RMSEs. At forested elevations below 3000 m, Noah-MP-
522 Cal and Noah-MP-Def have similar values of RMSE. Calibration decreases errors in the
523 higher elevation grid cells, which is often where more snow accumulates due to colder
524 temperatures coupled with orographic lifting. We also calculate the ratio of RMSE to mean
525 snow depth (not shown), and for Noah-MP-Cal, this metric decreases with elevation, while
526 for Noah-MP-Def, it increases above 2500 m. Much of the increase in Noah-MP-Cal RMSE is
527 due to deeper snowpacks at higher elevation. SNODAS and UA are both based on
528 observational datasets, which likely have larger uncertainty in forests. Noah-MP-Cal is in
529 better agreement with the observation-based gridded data products than Noah-MP-Def,
530 but the “true” accuracy in forested environments is limited by a lack of observations in
531 forests.

532

533 *4.2.2 GHCN Comparisons*

534 Across 79 GHCN stations, Noah-MP-Cal is less biased (0.0049 m) and has a lower
535 RMSE (0.15 m) than Noah-MP-Def (bias of -0.04 m and RMSE of 0.20 m). Noah-MP-Cal
536 generally reduces the snow depth bias in Noah-MP-Def in the Front Range and broadly
537 reduces RMSE across the full domain (Figure 9). While results are generally similar

538 between Noah-MP-Cal and Noah-MP-Def, the evaluation with GHCN demonstrates an
539 additional independent check that calibration improves the performance of modeled snow
540 depth.

541

542 *4.2.3 SnowEx Comparisons*

543 Finally, we also evaluate snow depth and SWE against 264 snow pit observations
544 from the NASA SnowEx 2017 field campaigns at Grand Mesa and Senator Beck (Figure 10).
545 Here we include SnowModel simulations in the comparison to consider a snow process
546 model. For both Noah-MP and SnowModel, we select the grid that contains each snow pit
547 for the comparison. SnowModel is kept at its native 30 m resolution, though we also tested
548 average SnowModel grid cells to the Noah-MP resolution and results were similar. The
549 majority of pit observations ($n = 224$) are from Grand Mesa, where there is better
550 agreement after calibration for snow depth (Table 5): mean bias decreases from -48.2 cm
551 to -12.1 cm (mean percent absolute difference decreases from 32.2% to 20.0%) and RMSE
552 decreases from 54.4 cm to 34.9 cm. Similar for SWE, Noah-MP-Cal has a smaller SWE mean
553 bias at Grand Mesa than Noah-MP-Def (-23.0 mm vs. -160.6 mm) and a smaller RMSE
554 (132.9 mm vs 185.4 mm). For SWE, SnowModel has better agreement with snow pits than
555 either Noah-MP simulation, though the performance of SnowModel and Noah-MP-Cal are
556 comparable for snow depth, with Noah-MP-Cal having smaller MAE and RMSE. Similar
557 performance for snow depth and SWE disagreements may be due to different density
558 estimates in SnowModel and Noah-MP-Cal. At Senator Beck ($n = 40$ pits), where we do not
559 have SnowModel simulations, Noah-MP-Cal greatly improves upon Noah-MP-Def
560 evaluation metrics for both snow depth and SWE: for snow depth (SWE), RMSE increases

561 from 49.7 cm to 102.5 cm (167.6 mm to 413.0). This highlights the uneven performance
562 across the domain after calibration.

563 Spatially, Noah-MP-Def has much lower values of snow depth than measured in the
564 snow pits on a single day (Figure 11). Noah-MP-Cal, on the other hand, has spatial patterns
565 that better match the snow pits observations throughout the Grand Mesa study site,
566 capturing the overall east-west gradient seen in the snow pit observations and in the
567 SnowModel simulation. Calibrated Noah-MP at a 1-km resolution has similar error metrics
568 to an uncalibrated snow process model at a 30-m resolution, but this evaluation is only
569 possible over a small portion of the full domain.

570 Finally, we evaluate Noah-MP simulations against ASO lidar snow depth
571 observations from SnowEx flights on February 8 and 16 (Painter, 2018). Spatially,
572 estimates from ASO, Noah-MP-Cal, and Noah-MP-Def have somewhat similar patterns on
573 each flight day, with snow depth tending to increase toward the eastern portion of the
574 domain (Figure 12). ASO and Noah-MP-Cal also show that snow depth increases from the
575 north to the south across the domain; Noah-MP-Def, on the other hand, has lower
576 variability across the domain. Note the deeper band of snow in the ASO observations along
577 the northern portion of the domain. The deeper snow here is likely due to snow
578 accumulating at the base of the cliff. Snow persistence maps (Figure S3) show that snow
579 historically lingers longer along the base of the cliff, suggesting that the deeper snow
580 depths in ASO are plausible. Noah-MP, with grid cells orders of magnitude coarser than
581 ASO, cannot capture this fine scale spatial pattern. For both flight days, Noah-MP-Cal has
582 higher values of SPAEF, which indicates better spatial agreement with ASO observations:
583 on February 8, Noah-MP-Cal has a SPAEF value of 0.408 and Noah-MP-Def has a value of

584 0.253; on February 16, Noah-MP-Cal has a SPAEF of 0.516 compared to 0.195 from Noah-
585 MP-Def.

586 Originally collected at 3 m spatial resolution, ASO snow depth observations are
587 aggregated to 0.01° resolution to match the Noah-MP simulations by averaging together
588 over 100,000 ASO 3 m grid cells. In evaluations of Noah-MP grid cells against aggregated
589 ASO depth observation, Noah-MP-Def underestimates ASO, and Noah-MP-Cal
590 overestimates for snow depths above 1.5 m (Figure 12). For each flight day, Noah-MP-Cal
591 has smaller RMSE, MAE, and bias magnitude than Noah-MP-Def. From this comparison,
592 calibration may lead to overestimates of snow depth in some regions, but calibration
593 introduces more realistic spatial patterns of snow depth, as compared to ASO observations.

594

595 4.3. Streamflow Evaluation

596 Beyond impacts on snow depth and SWE, calibration will impact LSM simulation of
597 other hydrological variables. For six basins within the Colorado domain with little-to-no
598 human management, calibration can improve streamflow estimation (Figure 13). Of the six
599 basins, four have higher NSE values for Noah-MP-Cal than Noah-MP-Def. After calibration,
600 however, four of the six basins still have negative NSE values, though the streamflow bias
601 may not all be due to snow. For the two basins with NSE>0 (9072500 and 9081600),
602 calibration improves performance, though only 9072500 has a statistically significant
603 difference in monthly streamflow between Noah-MP-Cal and Noah-MP-Def. In two basins,
604 both on the Gunnison River (9124700 and 9127800), calibration leads to a larger
605 overestimation in streamflow for some evaluation years. For the Colorado River at
606 Glenwood Springs (9072500), Noah-MP-Def largely underestimates streamflow, and

607 calibration addresses this bias through increased runoff. In most years for most basins,
608 Noah-MP-Cal has later peak streamflow, in agreement with the observations, which is also
609 noted in Figure 13. In 9107000, where Noah-MP-Def overestimates observations, Noah-
610 MP-Cal decreases the magnitude of the bias, though Noah-MP-Cal still overestimates
611 slightly; in 9081600, Noah-MP-Def underestimates observed streamflow, and the
612 calibrated runoff value is a better match for the observations. This demonstrates that
613 calibration does not increase snow and runoff in one direction, but rather calibration can
614 improve upon both positive and negative biases. Results similar to the small basin analysis
615 are seen across the full model domain, including higher springtime streamflow in Noah-
616 MP-Cal compared to Noah-MP-Def (Figure S4a,b,c). Peak streamflow in Noah-MP-Cal also
617 generally occurs later in the year than Noah-MP-Def (Figure S4f), in agreement with later
618 snowmelt in Noah-MP-Cal (Figures 3 and 8).

619

620 **5. Discussion**

621 **5.1. Summary of results**

622 Here we investigate the impact of model calibration on simulations of snow depth,
623 SWE, and streamflow. From this calibration exercise, we aim to answer the three research
624 questions posed in the introduction. First, calibration can address dry biases in LSMs,
625 which often result in underestimation of snow. We show improvements to not only
626 simulated SWE and snow depth magnitude but also to timing, of both accumulation and
627 ablation periods. Calibration also results in Noah-MP-Cal performing about as well as an
628 uncalibrated, high-resolution snow process model (Table 5, Figure 11). Though evaluations
629 of Noah-MP and SnowModel are limited to Grand Mesa, the spatial variability of snow

630 depth across the mesa are similar in SnowModel and Noah-MP-Cal, though SnowModel
631 simulations produce more detail with the finer spatial resolution. When comparing both
632 models to snow pit measurements, Noah-MP-Cal actually has better performance for snow
633 depth, though error Noah-MP-Cal metrics are larger for SWE. Results are similar for high
634 resolution ASO lidar, where Noah-MP-Cal captures the realistic spatial variability in ASO
635 estimates, suggesting that, over Grand Mesa at least, the calibration procedure largely
636 improves the model simulation.

637 Second, impacts from calibration are observed beyond snow variables. For
638 streamflow, Noah-MP-Cal improves estimates for four of the six study basins. Here we are
639 limited to small unmanaged basins or reconstructed estimates of natural streamflow since
640 Noah-MP does not include human management. We note that we are not evaluating routed
641 streamflow here, but instead, we consider grid cell estimates of surface and subsurface
642 runoff. Future work should consider dynamically routed streamflow in order to account for
643 time lags between the upper reaches of the watershed and the evaluation point with the
644 stream gage. Even with those considerations, improved NSE metrics suggest that the
645 increased snowpack in Noah-MP-Cal results in streamflow magnitude and timing that
646 better matches observations.

647 Finally, the calibration highlights potential avenues for improving both model
648 configuration and meteorological forcing data, though calibrated parameters may be
649 reflective of the choice of forcing dataset (Elsner et al., 2014). The genetic algorithm
650 procedure produces spatially varying model parameters, as compared to the spatially
651 uniform parameters used in the default Noah-MP configuration. In particular, we highlight
652 ten parameters within Noah-MP that are likely candidates for further investigation. We

653 show that allowing these parameters to vary in space results in improved model
654 performance compared to the default, spatially uniform values. Some of the parameters,
655 such as *SNOWF_SCALE_*, appear to have a relationship with elevation (compare Figure 2f
656 with Figure 1), while other parameters, such as *MXSNALB* and *ZOSNO*, appear to be more
657 related to land class category (compare Figure 2b,j with Figure 1). Future efforts should
658 determine new estimates for these parameters, perhaps through investigation of
659 relationships with landscape characteristics, such as elevation, vegetation class, and soil
660 type.

661 Global maps of the sensitive parameters could likely improve simulation of snow
662 without the need for a computationally expensive calibration procedure. In addition to
663 investigating relationships for creating spatially varying parameters, work should consider
664 whether parameters should also vary in time. Creating new estimates of spatially and
665 temporally varying parameters could improve snow modeling without the data
666 requirement of calibration, which would have implications for our ability to estimate global
667 snow, regardless of data availability. Efforts to scale snow parameters examined here to
668 larger domains are under development, resulting in spatially varying parameter estimates
669 for all of CONUS.

670 In addition to the ten parameters from Noah-MP discussed above, results from
671 calibration demonstrate that introducing the snowfall scaling term has a large impact on
672 the snow accumulation magnitude. This points to the need for better meteorological
673 forcing data, particularly for precipitation at high elevations. There is often high variability
674 between precipitation estimates from differing models and reanalyses (Decker et al., 2012;
675 Essou et al., 2016; Henn et al., 2018; Hughes et al., 2017; Wrzesien et al., 2019), and

676 previous studies have suggested that NLDAS-2 precipitation is too low in mountain regions
677 (Enzinger et al., 2019; He et al., 2019; Henn et al., 2018; Smyth et al., 2020); such
678 uncertainty will be propagated into the LSM. However, improving large scale precipitation
679 estimates is not trivial, and model-based precipitation estimates often outperform
680 observation-based estimates in mountain areas (Lundquist et al., 2019), despite known
681 model biases. If we cannot improve estimates of precipitation and snowfall in the forcing
682 datasets, informing modeled snowpack estimates with observations of SWE and snow
683 depth is likely the best option. This calibration procedure highlights a method for
684 addressing biases in both meteorological forcing and the LSM itself and results in improved
685 simulations of snow in a topographically complex region.

686

687 5.2. Implications for Snow OSSE

688 As discussed, the Noah-MP-Cal simulation presented here will be used as the nature
689 run (NR) in a snow-focused Observing System Simulation Experiment (OSSE), where the
690 NR is designed to approximate the “truth”, i.e., actual snow conditions. Though calibration
691 is not a panacea for reducing all model uncertainty, the improved performance from Noah-
692 MP-Def to Noah-MP-Cal provides compelling support for Noah-MP-Cal to be the NR for the
693 OSSE. Of particular concern when designing the OSSE was whether the NR could address
694 the common underestimation of snow at higher elevations, which is necessary for
695 understanding how proposed sensors will observe realistic ranges of snow conditions.
696 Calibrating Noah-MP against UA SWE estimates reduces the negative bias for SWE and
697 snow depth and results in snow spatial heterogeneity that better matches both UA and
698 SNODAS.

699 While Noah-MP-Cal is not without error, a NR is not expected to perfectly replicate
700 actual conditions, and no true observations are used in an OSSE. Therefore, the spatial and
701 temporal variability in Noah-MP-Cal is adequate for approximating realistic snow
702 conditions for the western Colorado domain. The main drawback of Noah-MP as the NR –
703 whether the default or calibrated configuration – is that Noah-MP does not provide
704 estimates of snow grain size. Understanding how satellite observations are impacted by
705 snow grain size and metamorphism is a fundamentally important question (Durand et al,
706 2018; Nolin, 2010; Foster et al., 2005). However, no models within the current LIS
707 framework provide estimates of snow grain size, though work is ongoing to implement
708 new snow models into LIS. While Noah-MP-Cal will be used in the OSSE described here,
709 future work will consider a follow on OSSE that incorporates a model that does include the
710 simulation of grain size.

711

712 5.3. Challenges with Calibration

713 With a calibration exercise such as this one, there are a few notable challenges.
714 While calibration can lead to domain-averaged improvements in the targeted variable, as
715 presented here for SWE and snow depth, it can cause degraded performance in individual
716 regions across the domain. We see this in the northern portion of the domain (Figure 4a) to
717 the west of the Cameron Pass evaluation site (Figure 1). After the genetic algorithm
718 optimization, the snowfall scale term is high in this region (Figure 2k), resulting in snow
719 depths and SWE values that much larger than either UA or SNODAS. In the calibration
720 period, SWE estimates from UA were particularly large, where the 2007-2009 average peak
721 SWE value for this area from UA is higher than the average peak SWE value for 2010-2020,

722 causing the calibration to be trained on higher-than average SWE. Anomalies such as this
723 from calibration are often unavoidable.

724 Another challenge with our calibration setup is that the parameters are constant in
725 time. Therefore, even if Noah-MP-Def performs well compared to UA, the calibrated
726 parameters will still be applied. For example, in water years 2017, 2019, and 2020,
727 domain-averaged SWE and snow depth from Noah-MP-Def is similar to UA (Figure 3).
728 Applying the calibrated parameters generally results in increased snow values, and as a
729 result, Noah-MP-Cal overestimates SWE in these years. Calibration improves performance
730 over the full study period (Table 3), but it does not always result in better performance for
731 an individual year or season. As discussed above with the spatial anomalies, calibration will
732 not result in uniformly improved performance.

733 For all calibration procedures, such as the genetic algorithm used here, a “truth”
734 dataset is required to calibrate against, and data availability is limited in many regions,
735 especially in high elevations and high latitudes where much of the global snow
736 accumulates. Therefore, while the calibration procedure presented here is a critical step for
737 the ongoing OSSE and for improving the representation of the truth, calibrating over a well-
738 observed Colorado domain may not necessarily improve the model performance of global
739 snow. Results presented here may not reflect other regions with differing snow conditions,
740 such as maritime snow in the Pacific Northwest or tundra snow in the high latitudes (e.g.,
741 Kim et al., 2021). Future work will investigate similar calibration methods in other regions.
742 We hypothesize that in regions with high precipitation uncertainty, such as mountainous
743 regions, the snowfall scaling term will have similar impacts on snow magnitude as
744 presented here.

745 Since we only calibrate against SWE and do not include additional constraints in the
746 objective function, such as for streamflow, the calibration cannot directly address biases in
747 other model processes. In the streamflow analyses, we see that Noah-MP-Def does not have
748 good agreement with the observed runoff (Figure 13 and Table 6). However, further
749 observational constraints or model improvements (possibly unrelated to snow processes)
750 are required to address runoff biases that we show here. In operational modeling, it is
751 standard to calibrate snowmelt rates to runoff (e.g. Hay et al., 2006; Franz and Karsten,
752 2013; Turcotte et al., 2017), in order to constrain snow ablation. Here, though, we do not
753 calibrate against runoff. Degradation in unconstrained variables, such as runoff, are not
754 uncommon during calibration efforts (e.g. Franz and Karsten, 2013; Nemri and Kinnard).
755 Future efforts could consider multi-criteria objective functions to reduce biases in both
756 snow variables and streamflow.

757 Beyond the calibration, evaluating gridded data with point observations presents
758 additional challenges. There are significant differences in what an ~ 1 m observation, such
759 as a snow pit or a GHCN station, measures and what a ~ 1000 m model grid cell simulates.
760 Since snow depth and SWE measurements are typically point observations, this imperfect
761 comparison is often necessary for evaluating models. However, during extensive field
762 campaigns, such as SnowEx 2017, numerous observations are made in a small domain over
763 a short period of time. While the result is still point-to-grid comparisons, the high density of
764 observation allows for a more complete evaluation, if over a limited domain. Though we do
765 acknowledge the uncertainty from scale differences, we aim to provide as thorough an
766 evaluation as we can for demonstrating evidence of the improved performance of Noah-
767 MP-Cal over Noah-MP-Def through the numerous independent comparison datasets.

768

769 **6. Conclusion**

770 The Noah-MP-Cal and Noah-MP-Def evaluation demonstrates that calibrating a land
771 surface model against an observation-based SWE dataset (e.g., the University of Arizona
772 dataset), improves model performance of snow, though not uniformly across the domain.
773 The calibration procedure was motivated by an ongoing Observing System Simulation
774 Experiment (OSSE) to evaluate the utility of proposed snow satellite sensors, and Noah-
775 MP-Cal will be used as the nature run for the OSSE. That is, the improved Noah-MP
776 simulation will act as the “truth” for the OSSE, upon which synthetic observations will be
777 created. However, results presented here have important implications beyond the OSSE.
778 We demonstrate a method for improving spatiotemporal estimates of snow, and we show
779 that spatially uniform values of key model parameters result in worse performance.
780 Allowing parameters to vary spatially, as we do in the Noah-MP-Cal simulation after a
781 genetic algorithm optimization procedure, results in improved model performance of both
782 snow depth and SWE. Future model development could consider implementing distributed
783 values of sensitive parameters, which might improve LSM simulations without the need for
784 an initial calibration step.

785

786 **7. Acknowledgements**

787 This research was supported by Grants from the National Aeronautics and Space
788 Administration (#AIST18-0041, #AIST18-0045, and #NNH16ZDA001N). Computing was
789 supported by the resources at the NASA Center for Climate Simulations. We thank three
790 anonymous reviewers for feedback that strengthened the overall manuscript.

791

792 **8. Data Availability Statement**

793 All datasets described here for modeling forcing and evaluation are available through the
794 provided citations within the text. The University of Arizona data are available for
795 download from the National Snow and Ice Data Center (NSIDC, Broxton et al., 2019). The
796 Noah-MP model simulations upon which this study is based are too large to be publicly
797 archived with available resources, though all model output are stored on the NASA
798 Discover supercomputer system through the NASA Center for Climate Simulations. To
799 replicate the simulation, interested users can access the NASA Land Information System at
800 <https://github.com/NASA-LIS/LISF>.

801

802 9. References

- 803 Aguado, E. (1985). Radiation Balances of Melting Snow Covers at an Open Site in the
804 Central Sierra Nevada, California. *Water Resources Research*, 21(11), 1649–1654.
805 <https://doi.org/10.1029/WR021i011p01649>
- 806 Ahl, R. S., Woods, S. W., & Zuuring, H. R. (2008). Hydrologic Calibration and Validation of
807 SWAT in a Snow-Dominated Rocky Mountain Watershed, Montana, U.S.A.1. *JAWRA*
808 *Journal of the American Water Resources Association*, 44(6), 1411–1430.
809 <https://doi.org/10.1111/j.1752-1688.2008.00233.x>
- 810 Amorocho, J., & Espildora, B. (1966). *Mathematical Simulation of the Snow Melting Process*.
811 Department of Water Science and Engineering, University of California, Davis.
- 812 Anderson, E. (1973), National Weather Service River Forecast system—Snow accumulation
813 and ablation model, NOAA Tech. Memo. NWS HYDRO-17, 217 pp., U.S. Dep. of
814 Commer., Silver Spring, Md.
- 815 Barnett, T. P., Adam, J. C., & Lettenmaier, D. P. (2005). Potential impacts of a warming
816 climate on water availability in snow-dominated regions. *Nature*, 438(7066), 303–
817 309. <https://doi.org/10.1038/nature04141>
- 818 Barrett, A. P. (2003), National Operational Hydrologic Remote Sensing Center Snow Data
819 Assimilation System (SNODAS) Products at NSIDC, 19 pp., Natl. Snow and Ice Data
820 Cent., Coop. Inst. for Res. in Environ. Sci., Boulder, Colo.
- 821 Bormann, K. J., Brown, R. D., Derksen, C., & Painter, T. H. (2018). Estimating snow-cover
822 trends from space. *Nature Climate Change*, 8(11), 924–928.
823 <https://doi.org/10.1038/s41558-018-0318-3>
- 824 Broxton, P. D., Dawson, N., & Zeng, X. (2016). Linking snowfall and snow accumulation to
825 generate spatial maps of SWE and snow depth. *Earth and Space Science*, 3(6), 246–
826 256. <https://doi.org/10.1002/2016EA000174>
- 827 Broxton, P. D., Zeng, X., & Dawson, N. (2016). Why Do Global Reanalyses and Land Data
828 Assimilation Products Underestimate Snow Water Equivalent? *Journal of*
829 *Hydrometeorology*, 17(11), 2743–2761. <https://doi.org/10.1175/JHM-D-16-0056.1>
- 830 Broxton, P., X. Zeng, and N. Dawson. 2019. *Daily 4 km Gridded SWE and Snow Depth from*
831 *Assimilated In-Situ and Modeled Data over the Conterminous US, Version 1*. [Indicate
832 subset used]. Boulder, Colorado USA. NASA National Snow and Ice Data Center
833 Distributed Active Archive Center. <https://doi.org/10.5067/OGGPB220EX6A>.

834 Carroll, T., Cline, D., Fall, G., Nilsson, A., Li, L., & Rost, A. 2001. NOHRSC operations and the
835 simulation of snow cover properties for the coterminous US. In Proc. 69th Annual
836 Meeting of the Western Snow Conf (pp. 1-14).

837 Chen, Fei, Barlage, M., Tewari, M., Rasmussen, R., Jin, J., Lettenmaier, D., et al. (2014).
838 Modeling seasonal snowpack evolution in the complex terrain and forested
839 Colorado Headwaters region: A model intercomparison study. *Journal of Geophysical
840 Research: Atmospheres*, 119(24), 13,795-13,819.
841 <https://doi.org/10.1002/2014JD022167>

842 Chen, Feng, Liu, C., Dudhia, J., & Chen, M. (2014). A sensitivity study of high-resolution
843 regional climate simulations to three land surface models over the western United
844 States. *Journal of Geophysical Research: Atmospheres*, 119(12), 7271–7291.
845 <https://doi.org/10.1002/2014JD021827>

846 Cho, E., & Jacobs, J. M. (2020). Extreme value snow water equivalent and snowmelt for
847 infrastructure design over the contiguous United States. *Water Resources Research*,
848 56, e2020WR028126. <https://doi.org/10.1029/2020WR028126>

849 Cho, E., Jacobs, J. M., & Vuyovich, C. M. (2020). The Value of Long-Term (40 years) Airborne
850 Gamma Radiation SWE Record for Evaluating Three Observation-Based Gridded
851 SWE Data Sets by Seasonal Snow and Land Cover Classifications. *Water Resources
852 Research*, 56(1), e2019WR025813. <https://doi.org/10.1029/2019WR025813>

853 Chow, V.T. (1964) Handbook of Applied Hydrology. McGraw-Hill Book Company, New York.

854 Clow, D. W., Nanus, L., Verdin, K. L., & Schmidt, J. (2012). Evaluation of SNODAS snow depth
855 and snow water equivalent estimates for the Colorado Rocky Mountains, USA.
856 *Hydrological Processes*, 26(17), 2583–2591. <https://doi.org/10.1002/hyp.9385>

857 Crow, W. T., Drusch, M., & Wood, E. F. (2001). An observation system simulation
858 experiment for the impact of land surface heterogeneity on AMSR-E soil moisture
859 retrieval. *IEEE Transactions on Geoscience and Remote Sensing*, 39(8), 1622-1631,
860 <https://doi.org/10.1109/36.942540>.

861 Crow, W. T. et al. (2005). An observing system simulation experiment for hydros
862 radiometer-only soil moisture products. *IEEE Transactions on Geoscience and
863 Remote Sensing*, 43(6), 1289-1303. <https://doi.org/10.1109/TGRS.2005.845645>.

864 Cuntz, M., Mai, J., Samaniego, L., Clark, M., Wulfmeyer, V., Branch, O., et al. (2016). The
865 impact of standard and hard-coded parameters on the hydrologic fluxes in the
866 Noah-MP land surface model. *Journal of Geophysical Research: Atmospheres*, 121(18),
867 10,676-10,700. <https://doi.org/10.1002/2016JD025097>

868 Daly, C., G. H. Taylor, W. P. Gibson, T. W. Parzybok, G. L. Johnson, & P. A. Pasteris. (2000).
869 HIGH-QUALITY SPATIAL CLIMATE DATA SETS FOR THE UNITED STATES AND
870 BEYOND. *Transactions of the ASAE*, 43(6), 1957–1962.
871 <https://doi.org/10.13031/2013.3101>

872 Dawson, N., Broxton, P., & Zeng, X. (2017). A New Snow Density Parameterization for Land
873 Data Initialization. *Journal of Hydrometeorology*, 18(1), 197–207.
874 <https://doi.org/10.1175/JHM-D-16-0166.1>

875 Dawson, N., Broxton, P., & Zeng, X. (2018). Evaluation of Remotely Sensed Snow Water
876 Equivalent and Snow Cover Extent over the Contiguous United States. *Journal of
877 Hydrometeorology*, 19(11), 1777–1791. <https://doi.org/10.1175/JHM-D-18-0007.1>

878 Decker, M., Brunke, M. A., Wang, Z., Sakaguchi, K., Zeng, X., & Bosilovich, M. G. (2012).
879 Evaluation of the Reanalysis Products from GSFC, NCEP, and ECMWF Using Flux
880 Tower Observations. *Journal of Climate*, 25(6), 1916–1944.
881 <https://doi.org/10.1175/JCLI-D-11-00004.1>

882 Demirel, M. C., Mai, J., Mendiguren, G., Koch, J., Samaniego, L., & Stisen, S. (2018). Combining
883 satellite data and appropriate objective functions for improved spatial pattern
884 performance of a distributed hydrologic model. *Hydrology and Earth System
885 Sciences*, 22(2), 1299-1315. <https://doi.org/10.5194/hess-22-1299-2018>

886 Dirmhirn, I., & Eaton, F. D. (1975). Some Characteristics of the Albedo of Snow. *Journal of
887 Applied Meteorology*, 14(3), 375–379. [https://doi.org/10.1175/1520-
888 0450\(1975\)014<0375:SCOTAO>2.0.CO;2](https://doi.org/10.1175/1520-0450(1975)014<0375:SCOTAO>2.0.CO;2)

889 Dozier, J., Bair, E. H., & Davis, R. E. (2016). Estimating the spatial distribution of snow water
890 equivalent in the world's mountains: Spatial distribution of snow in the mountains.
891 *Wiley Interdisciplinary Reviews: Water*, 3(3), 461–474.
892 <https://doi.org/10.1002/wat2.1140>

893 Duethmann, D., Peters, J., Blume, T., Vorogushyn, S., & Güntner, A. (2014). The value of
894 satellite-derived snow cover images for calibrating a hydrological model in snow-
895 dominated catchments in Central Asia. *Water Resources Research*, 50(3), 2002–2021.
896 <https://doi.org/10.1002/2013WR014382>

897 Durand, M., Gatebe, C., Kim, E., Molotch, N., Painter, T.H., Raleigh, M., Sandells, M. &
898 Vuyovich, C. (2018). NASA SnowEx Science Plan: Assessing approaches for measuring
899 water in Earth's Seasonal Snow, version 1.6.

900 Dutra, E., Viterbo, P., Miranda, P. M. A., & Balsamo, G. (2012). Complexity of Snow Schemes
901 in a Climate Model and Its Impact on Surface Energy and Hydrology. *Journal of
902 Hydrometeorology*, 13(2), 521–538. <https://doi.org/10.1175/JHM-D-11-072.1>

903 Elder, K., L. Brucker, C. Hiemstra, and H. Marshall. 2018. *SnowEx17 Community Snow Pit
904 Measurements, Version 1*. [Indicate subset used]. Boulder, Colorado USA. NASA
905 National Snow and Ice Data Center Distributed Active Archive Center. doi:
906 <https://doi.org/10.5067/Q0310G1XULZS>.

907 Elsner, M. M., Gangopadhyay, S., Pruitt, T., Brekke, L. D., Mizukami, N., & Clark, M. P. (2014).
908 How Does the Choice of Distributed Meteorological Data Affect Hydrologic Model
909 Calibration and Streamflow Simulations? *Journal of Hydrometeorology*, 15(4), 1384–
910 1403. <https://doi.org/10.1175/JHM-D-13-083.1>

911 Enzinger, T. L., Small, E. E., & Borsa, A. A. (2019). Subsurface Water Dominates Sierra
912 Nevada Seasonal Hydrologic Storage. *Geophysical Research Letters*, 46(21), 11993–
913 12001. <https://doi.org/10.1029/2019GL084589>

914 Essery, R., Rutter, N., Pomeroy, J., Baxter, R., Stähli, M., Gustafsson, D., et al. (2009).
915 SNOWMIP2: An Evaluation of Forest Snow Process Simulations. *Bulletin of the
916 American Meteorological Society*, 90(8), 1120–1136.
917 <https://doi.org/10.1175/2009BAMS2629.1>

918 Essou, G. R. C., Sabarly, F., Lucas-Picher, P., Brissette, F., & Poulin, A. (2016). Can
919 Precipitation and Temperature from Meteorological Reanalyses Be Used for
920 Hydrological Modeling? *Journal of Hydrometeorology*, 17(7), 1929–1950.
921 <https://doi.org/10.1175/JHM-D-15-0138.1>

922 Etchevers, P., Martin, E., Brown, R., Fierz, C., Lejeune, Y., Bazile, E., et al. (2004). Validation
923 of the energy budget of an alpine snowpack simulated by several snow models

924 (Snow MIP project). *Annals of Glaciology*, 38, 150–158.
 925 <https://doi.org/10.3189/172756404781814825>

926 Foster, J. L., Sun, C., Walker, J. P., Kelly, R., Chang, A., Dong, J., & Powell, H. (2005).
 927 Quantifying the uncertainty in passive microwave snow water equivalent
 928 observations. *Remote Sensing of environment*, 94(2), 187-203.

929 Franz, K. J., Hogue, T. S., & Sorooshian, S. (2008). Operational snow modeling: Addressing
 930 the challenges of an energy balance model for National Weather Service forecasts.
 931 *Journal of Hydrology*, 360(1), 48–66. <https://doi.org/10.1016/j.jhydrol.2008.07.013>

932 Franz, K. J., Butcher, P., & Ajami, N. K. (2010). Addressing snow model uncertainty for
 933 hydrologic prediction. *Advances in Water Resources*, 33(8), 820–832.
 934 <https://doi.org/10.1016/j.advwatres.2010.05.004>

935 Franz, K. J., & Karsten, L. R. (2013). Calibration of a distributed snow model using MODIS
 936 snow covered area data. *Journal of Hydrology*, 494, 160–175.
 937 <https://doi.org/10.1016/j.jhydrol.2013.04.026>

938 Garnaud, C., Bélair, S., Carrera, M. L., Derksen, C., Bilodeau, B., Abrahamowicz, M., Gauthier,
 939 N., & Vionnet, V. (2019). Quantifying Snow Mass Mission Concept Trade-Offs Using
 940 an Observing System Simulation Experiment, *Journal of Hydrometeorology*, 20(1),
 941 155-173.

942 Hall, D. K. and G. A. Riggs. 2016. *MODIS/Terra Snow Cover Daily L3 Global 500m SIN Grid,*
 943 *Version 6.* Boulder, Colorado USA. NASA National Snow and Ice Data Center
 944 Distributed Active Archive Center. <https://doi.org/10.5067/MODIS/MOD10A1.006>.

945 Hall, D. K., Riggs, G. A., Salomonson, V. V., DiGirolamo, N. E., & Bayr, K. J. (2002). MODIS
 946 snow-cover products. *Remote Sensing of Environment*, 83(1–2), 181–194.
 947 [https://doi.org/10.1016/S0034-4257\(02\)00095-0](https://doi.org/10.1016/S0034-4257(02)00095-0)

948 Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S. A., Tyukavina, A., et al.
 949 (2013). High-resolution global maps of 21st-century forest cover change. *Science*,
 950 342(6160), 850-853. <https://doi.org/10.1126/science.1244693>

951 Harrison, K. W., Kumar, S. V., Peters-Lidard, C. D., and Santanello, J. A. (2012), Quantifying
 952 the change in soil moisture modeling uncertainty from remote sensing observations
 953 using Bayesian inference techniques, *Water Resour. Res.*, 48, W11514,
 954 [doi:10.1029/2012WR012337](https://doi.org/10.1029/2012WR012337).

955 He, C., Chen, F., Barlage, M., Liu, C., Newman, A., Tang, W., et al. (2019). Can Convection-
 956 Permitting Modeling Provide Decent Precipitation for Offline High-Resolution
 957 Snowpack Simulations Over Mountains? *Journal of Geophysical Research:*
 958 *Atmospheres*, 124(23), 12631–12654. <https://doi.org/10.1029/2019JD030823>

959 He, M., Hogue, T. S., Franz, K. J., Margulis, S. A., & Vrugt, J. A. (2011). Characterizing
 960 parameter sensitivity and uncertainty for a snow model across hydroclimatic
 961 regimes. *Advances in Water Resources*, 34(1), 114–127.
 962 <https://doi.org/10.1016/j.advwatres.2010.10.002>

963 Henn, B., Clark, M. P., Kavetski, D., McGurk, B., Painter, T. H., & Lundquist, J. D. (2016).
 964 Combining snow, streamflow, and precipitation gauge observations to infer basin-
 965 mean precipitation. *Water Resources Research*, 52(11), 8700–8723.
 966 <https://doi.org/10.1002/2015WR018564>

967 Henn, B., Newman, A. J., Livneh, B., Daly, C., & Lundquist, J. D. (2018). An assessment of
 968 differences in gridded precipitation datasets in complex terrain. *Journal of*
 969 *Hydrology*, 556, 1205–1219. <https://doi.org/10.1016/j.jhydrol.2017.03.008>

970 Holtzman, N. M., Pavelsky, T. M., Cohen, J. S., Wrzesien, M. L., & Herman, J. D. (2020).
971 Tailoring WRF and Noah-MP to Improve Process Representation of Sierra Nevada
972 Runoff: Diagnostic Evaluation and Applications. *Journal of Advances in Modeling*
973 *Earth Systems*, 12(3), e2019MS001832. <https://doi.org/10.1029/2019MS001832>
974 Hughes, M., Lundquist, J. D., & Henn, B. (2017). Dynamical downscaling improves upon
975 gridded precipitation products in the Sierra Nevada, California. *Climate Dynamics*.
976 <https://doi.org/10.1007/s00382-017-3631-z>
977 Immerzeel, W.W., Lutz, A.F., Andrade, M. *et al.* Importance and vulnerability of the world's
978 water towers. *Nature* **577**, 364–369 (2020). [https://doi.org/10.1038/s41586-019-](https://doi.org/10.1038/s41586-019-1822-y)
979 [1822-y](https://doi.org/10.1038/s41586-019-1822-y)
980 Isenstein, E. M., Wi, S., & Ethan Yang, Y. C. (2015). Calibration of a Distributed Hydrologic
981 Model Using Streamflow and Remote Sensing Snow Data. In *World Environmental*
982 *and Water Resources Congress 2015* (pp. 973–982). Austin, TX: American Society of
983 Civil Engineers. <https://doi.org/10.1061/9780784479162.093>
984 Jordan, R. (1991). A one-dimensional temperature model for a snow cover: technical
985 documentation for SNTherm.89. CRREL Spec. Rep. 91-16.
986 Kim, R. S., Kumar, S., Vuyovich, C., Houser, P., Lundquist, J., Mudryk, L., et al. (2021). Snow
987 Ensemble Uncertainty Project (SEUP): quantification of snow water equivalent
988 uncertainty across North America via ensemble land surface modeling, *The*
989 *Cryosphere*, 15, 771–791, <https://doi.org/10.5194/tc-15-771-2021>
990 Koch, J., Demirel, M. C., & Stisen, S. (2018). The SPAtial EFFiciency metric (SPAEF): multiple-
991 component evaluation of spatial patterns for optimization of hydrological models.
992 *Geoscientific Model Development*, 11(5), 1873-1886. [https://doi.org/10.5194/gmd-](https://doi.org/10.5194/gmd-11-1873-2018)
993 [11-1873-2018](https://doi.org/10.5194/gmd-11-1873-2018)
994 Koster, R. D., Mahanama, S. P., Livneh, B., Lettenmaier, D. P., & Reichle, R. H. (2010). Skill in
995 streamflow forecasts derived from large-scale estimates of soil moisture and snow.
996 *Nature Geoscience*, 3(9), 613-616. Krinner, G., Derksen, C., Essery, R., Flanner, M.,
997 Hagemann, S., Clark, M., et al. (2018). ESM-SnowMIP: assessing snow models and
998 quantifying snow-related climate feedbacks. *Geoscientific Model Development*,
999 11(12), 5027–5049. <https://doi.org/10.5194/gmd-11-5027-2018>
1000 Kumar, S. V., Peters-Lidard, C. D., Tian, Y., Houser, P. R., Geiger, J., Olden, S., et al. (2006).
1001 Land information system: An interoperable framework for high resolution land
1002 surface modeling. *Environmental Modelling & Software*, 21(10), 1402–1415.
1003 <https://doi.org/10.1016/j.envsoft.2005.07.004>
1004 Kumar, S. V., Peters-Lidard, C. D., Santanello, J., Harrison, K., Liu, Y., & Shaw, M. (2012). Land
1005 surface Verification Toolkit (LVT) – a generalized framework for land surface model
1006 evaluation. *Geoscientific Model Development*, 5(3), 869–886.
1007 <https://doi.org/10.5194/gmd-5-869-2012>
1008 Kumar, Sujay V., M. Mocko, D., Wang, S., Peters-Lidard, C. D., & Borak, J. (2019). Assimilation
1009 of Remotely Sensed Leaf Area Index into the Noah-MP Land Surface Model: Impacts
1010 on Water and Carbon Fluxes and States over the Continental United States. *Journal*
1011 *of Hydrometeorology*, 20(7), 1359–1377. <https://doi.org/10.1175/JHM-D-18-0237.1>
1012 Kumar, S. V., Reichle, R. H., Harrison, K. W., Peters-Lidard, C. D., Yatheendradas, S., &
1013 Santanello, J. A. (2012). A comparison of methods for a priori bias correction in soil
1014 moisture data assimilation. *Water Resources Research*, 48(3).
1015 <https://doi.org/10.1029/2010WR010261>

1016 Lettenmaier, D. P., Alsdorf, D., Dozier, J., Huffman, G. J., Pan, M., & Wood, E. F. (2015).
1017 Inroads of remote sensing into hydrologic science during the WRR era: REMOTE
1018 SENSING. *Water Resources Research*, 51(9), 7309–7342.
1019 <https://doi.org/10.1002/2015WR017616>

1020 Li, D., Wrzesien, M. L., Durand, M., Adam, J., & Lettenmaier, D. P. (2017). How much runoff
1021 originates as snow in the western United States, and how will that change in the
1022 future?: Western U.S. Snowmelt-Derived Runoff. *Geophysical Research Letters*,
1023 44(12), 6163–6172. <https://doi.org/10.1002/2017GL073551>

1024 Liston, G. E., & Elder, K. (2006). A distributed snow-evolution modeling system
1025 (SnowModel). *Journal of Hydrometeorology*, 7(6), 1259–1276.
1026 <https://doi.org/10.1175/JHM548.1>

1027 Liston, G. E., & Elder, K. (2006b). A meteorological distribution system for high-resolution
1028 terrestrial modeling (MicroMet). *Journal of Hydrometeorology*, 7(2), 217–234.
1029 <https://doi.org/10.1175/JHM486.1>

1030 Liston, G. E., Haehnel, R. B., Sturm, M., Hiemstra, C. A., Berezovskaya, S., & Tabler, R. D.
1031 (2007). Simulating complex snow distributions in windy environments using
1032 SnowTran-3D. *Journal of Glaciology*, 53(181), 241–256.
1033 <https://doi.org/10.3189/172756507782202865>

1034 Liston, G. E., & Sturm, M. (1998). A snow-transport model for complex terrain. *Journal of*
1035 *Glaciology*, 44(148), 498–516. <https://doi.org/10.1017/S0022143000002021>

1036 Lundquist, J., Hughes, M., Gutmann, E., & Kapnick, S. (2019). Our skill in modeling mountain
1037 rain and snow is bypassing the skill of our observational networks. *Bulletin of the*
1038 *American Meteorological Society*, BAMS-D-19-0001.1.
1039 <https://doi.org/10.1175/BAMS-D-19-0001.1>

1040 Magnusson, J., Winstral, A., Stordal, A. S., Essery, R., & Jonas, T. (2017). Improving physically
1041 based snow simulations by assimilating snow depths using the particle filter. *Water*
1042 *Resources Research*, 53(2), 1125–1143. <https://doi.org/10.1002/2016WR019092>

1043 Marks, D., & Dozier, J. (1992). Climate and energy exchange at the snow surface in the
1044 Alpine Region of the Sierra Nevada: 2. Snow cover energy balance. *Water Resources*
1045 *Research*, 28(11), 3043–3054. <https://doi.org/10.1029/92WR01483>

1046 Mendoza, P. A., Clark, M. P., Barlage, M., Rajagopalan, B., Samaniego, L., Abramowitz, G., &
1047 Gupta, H. (2015). Are we unnecessarily constraining the agility of complex process-
1048 based models? *Water Resources Research*, 51(1), 716–728.
1049 <https://doi.org/10.1002/2014WR015820>

1050 Menne, M. J., Durre, I., Vose, R. S., Gleason, B. E., & Houston, T. G. (2012). An Overview of the
1051 Global Historical Climatology Network-Daily Database. *Journal of Atmospheric and*
1052 *Oceanic Technology*, 29(7), 897–910. <https://doi.org/10.1175/JTECH-D-11-00103.1>

1053 Minder, J. R., Letcher, T. W., & Skiles, S. M. (2016). An evaluation of high-resolution regional
1054 climate model simulations of snow cover and albedo over the Rocky Mountains,
1055 with implications for the simulated snow-albedo feedback: Evaluating mountain
1056 snow cover in RCMs. *Journal of Geophysical Research: Atmospheres*, 121(15), 9069–
1057 9088. <https://doi.org/10.1002/2016JD024995>

1058 Nash, J. E., & Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I
1059 — A discussion of principles. *Journal of Hydrology*, 10(3), 282–290.
1060 [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6)

1061 Nearing, G. S., Crow, W. T., Thorp, K. R., Moran, M. S., Reichle, R. H., and Gupta, H. V. (2012),
1062 Assimilating remote sensing observations of leaf area index and soil moisture for
1063 wheat yield estimates: An observing system simulation experiment, *Water*
1064 *Resources Research*, 48, W05525. <https://doi.org/10.1029/2011WR011420>.

1065 Nemri, S., & Kinnard, C. (2020). Comparing calibration strategies of a conceptual snow
1066 hydrology model and their impact on model performance and parameter
1067 identifiability. *Journal of Hydrology*, 582, 124474.

1068 Newman, A. J., Clark, M. P., Sampson, K., Wood, A., Hay, L. E., Bock, A., et al. (2015).
1069 Development of a large-sample watershed-scale hydrometeorological data set for
1070 the contiguous USA: data set characteristics and assessment of regional variability
1071 in hydrologic model performance. *Hydrology and Earth System Sciences*, 19(1), 209–
1072 223. <https://doi.org/10.5194/hess-19-209-2015>

1073 Newman, A., K. Sampson, M.P. Clark, A. Bock, and R.J. Viger, and D. Blodgett, 2014: A large-
1074 sample watershed-scale hydrometeorological dataset for the contiguous USA.
1075 Boulder, CO: UCAR/NCAR. doi:[10.5065/D6MW2F4D](https://doi.org/10.5065/D6MW2F4D)

1076 Niu, G.-Y., & Yang, Z.-L. (2007). An observation-based formulation of snow cover fraction
1077 and its evaluation over large North American river basins. *Journal of Geophysical*
1078 *Research*, 112(D21), D21101. <https://doi.org/10.1029/2007JD008674>

1079 Niu, G.-Y., Yang, Z.-L., Mitchell, K. E., Chen, F., Ek, M. B., Barlage, M., et al. (2011). The
1080 community Noah land surface model with multiparameterization options (Noah-
1081 MP): 1. Model description and evaluation with local-scale measurements. *Journal of*
1082 *Geophysical Research*, 116(D12), D12109. <https://doi.org/10.1029/2010JD015139>

1083 Nolin, A. W. (2010). Recent advances in remote sensing of seasonal snow. *Journal of*
1084 *Glaciology*, 56(200), 1141–1150. <https://doi.org/10.3189/002214311796406077>

1085 Painter, T. (2018). ASO L4 LiDAR snow depth 3m UTM grid, version 1. Grand Mesa
1086 Colorado. NASA National Snow and Ice Data Center Distributed Active Archive
1087 Center. <https://doi.org/10.5067/KIE9QNVG7HP0>

1088 Parajka, J., & Blöschl, G. (2008). The value of MODIS snow cover data in validating and
1089 calibrating conceptual hydrologic models. *Journal of Hydrology*, 358(3), 240–258.
1090 <https://doi.org/10.1016/j.jhydrol.2008.06.006>

1091 Parajka, J., Merz, R., & Blöschl, G. (2007). Uncertainty and multiple objective calibration in
1092 regional water balance modelling: case study in 320 Austrian catchments.
1093 *Hydrological Processes*, 21(4), 435–446. <https://doi.org/10.1002/hyp.6253>

1094 Peters-Lidard, C. D., Houser, P. R., Tian, Y., Kumar, S. V., Geiger, J., Olden, S., et al. (2007).
1095 High-performance Earth system modeling with NASA/GSFC's Land Information
1096 System. *Innovations in Systems and Software Engineering*, 3(3), 157–165.
1097 <https://doi.org/10.1007/s11334-007-0028-x>

1098 Pomeroy, J. W., Gray, D. M., Shook, K. R., Toth, B., Essery, R. L. H., Pietroniro, A., & Hedstrom,
1099 N. (1998). An evaluation of snow accumulation and ablation processes for land
1100 surface modelling. *Hydrological Processes*, 12(15), 2339–2367.
1101 [https://doi.org/10.1002/\(SICI\)1099-1085\(199812\)12:15<2339::AID-](https://doi.org/10.1002/(SICI)1099-1085(199812)12:15<2339::AID-HYP800>3.0.CO;2-L)
1102 [HYP800>3.0.CO;2-L](https://doi.org/10.1002/(SICI)1099-1085(199812)12:15<2339::AID-HYP800>3.0.CO;2-L)

1103 Raleigh, M. S., Lundquist, J. D., & Clark, M. P. (2015). Exploring the impact of forcing error
1104 characteristics on physically based snow simulations within a global sensitivity
1105 analysis framework. *Hydrology and Earth System Sciences*, 19(7), 3153–3179.
1106 <https://doi.org/10.5194/hess-19-3153-2015>

1107 Reba, M. L., Marks, D., Link, T. E., Pomeroy, J., & Winstral, A. (2014). Sensitivity of model
1108 parameterizations for simulated latent heat flux at the snow surface for complex
1109 mountain sites. *Hydrological Processes*, 28(3), 868–881.
1110 <https://doi.org/10.1002/hyp.9619>

1111 Rittger, K., Painter, T. H., & Dozier, J. (2013). Assessment of methods for mapping snow
1112 cover from MODIS. *Advances in Water Resources*, 51, 367–380.
1113 <https://doi.org/10.1016/j.advwatres.2012.03.002>

1114 Rollins, M. G.: LANDFIRE: a nationally consistent vegetation, wildland fire, and fuel
1115 assessment, *Int. J. Wildl. Fire*, 18(3), 235–249. <https://doi.org/10.1071/WF08088>,
1116 2009.

1117 Rutter, N., Essery, R., Pomeroy, J., Altimir, N., Andreadis, K., Baker, I., et al. (2009).
1118 Evaluation of forest snow processes models (SnowMIP2). *Journal of Geophysical
1119 Research*, 114(D6), D06111. <https://doi.org/10.1029/2008JD011063>

1120 Schmucki, E., Marty, C., Fierz, C., & Lehning, M. (2014). Evaluation of modelled snow depth
1121 and snow water equivalent at three contrasting sites in Switzerland using
1122 SNOWPACK simulations driven by different meteorological data input. *Cold Regions
1123 Science and Technology*, 99, 27–37.
1124 <https://doi.org/10.1016/j.coldregions.2013.12.004>

1125 Shafii, M., & De Smedt, F. (2009). Multi-objective calibration of a distributed hydrological
1126 model (WetSpa) using a genetic algorithm. *Hydrology and Earth System Sciences*,
1127 13(11), 2137–2149. <https://doi.org/10.5194/hess-13-2137-2009>

1128 Skiles, S. M. (2018). Grand Mesa study plot (version 1), Zenodo.
1129 <https://doi.org/10.5281/zenodo.1479859>

1130 Smyth, E. J., Raleigh, M. S., & Small, E. E. (2020). Improving SWE Estimation With Data
1131 Assimilation: The Influence of Snow Depth Observation Timing and Uncertainty.
1132 *Water Resources Research*, 56(5), e2019WR026853.
1133 <https://doi.org/10.1029/2019WR026853>

1134 Swain, M. J. and Ballard, D. H.: Color indexing, *Int. J. Comput. Vis.*, 7, 11–32,
1135 <https://doi.org/10.1007/BF00130487>, 1991.

1136 Takala, M., Luojus, K., Pulliainen, J., Derksen, C., Lemmetyinen, J., Kärnä, J.-P., et al. (2011).
1137 Estimating northern hemisphere snow water equivalent for climate research
1138 through assimilation of space-borne radiometer data and ground-based
1139 measurements. *Remote Sensing of Environment*, 115(12), 3517–3529.
1140 <https://doi.org/10.1016/j.rse.2011.08.014>

1141 ter Braak, C. J., & Vrugt, J. A. (2008). Differential evolution Markov chain with snooker
1142 updater and fewer chains. *Statistics and Computing*, 18(4), 435–446.

1143 Viviroli, D., Dürr, H. H., Messerli, B., Meybeck, M. & Weingartner, R. Mountains of the world,
1144 water towers for humanity: typology, mapping, and global significance. *Water
1145 Resources Research*, 43, 1–13 (2007). <https://doi.org/10.1029/2006WR005653>

1146 Versegny, D.L. (1991), Class—A Canadian land surface scheme for GCMs. I. Soil model. *Int. J.
1147 Climatol.*, 11: 111–133. <https://doi.org/10.1002/joc.3370110202>

1148 Vuyovich, C. M., Jacobs, J. M., & Daly, S. F. (2014). Comparison of passive microwave and
1149 modeled estimates of total watershed SWE in the continental United States. *Water
1150 Resources Research*, 50(11), 9088–9102. <https://doi.org/10.1002/2013WR014734>

1151 Wang, Q. J. (1991). The Genetic Algorithm and Its Application to Calibrating Conceptual
1152 Rainfall-Runoff Models. *Water Resources Research*, 27(9), 2467–2471.
1153 <https://doi.org/10.1029/91WR01305>

1154 Wang, X., Barker, D. M., Snyder, C., & Hamill, T. M. (2008). A Hybrid ETKF–3DVAR Data
1155 Assimilation Scheme for the WRF Model. Part I: Observing System Simulation
1156 Experiment, *Monthly Weather Review*, 136(12), 5116–5131.

1157 Wang, Y.-H., Broxton, P., Fang, Y., Behrangi, A., Barlage, M., Zeng, X., & Niu, G.-Y. (2019). A
1158 Wet-Bulb Temperature-Based Rain-Snow Partitioning Scheme Improves Snowpack
1159 Prediction Over the Drier Western United States. *Geophysical Research Letters*,
1160 46(23), 13825–13835. <https://doi.org/10.1029/2019GL085722>

1161 Webb, R. W., Raleigh, M. S., McGrath, D., Molotch, N. P., Elder, K., & Hiemstra, C., et al. (2020).
1162 Within-stand boundary effects on snow water equivalent distribution in forested
1163 areas. *Water Resources Research*, 56, e2019WR024905.
1164 <https://doi.org/10.1029/2019WR024905>

1165 Wrzesien, M. L., Pavelsky, T. M., Kapnick, S. B., Durand, M. T., & Painter, T. H. (2015).
1166 Evaluation of snow cover fraction for regional climate simulations in the Sierra
1167 Nevada. *International Journal of Climatology*, 35(9), 2472–2484.
1168 <https://doi.org/10.1002/joc.4136>

1169 Wrzesien, M. L., Durand, M. T., Pavelsky, T. M., Howat, I. M., Margulis, S. A., & Huning, L. S.
1170 (2017). Comparison of Methods to Estimate Snow Water Equivalent at the Mountain
1171 Range Scale: A Case Study of the California Sierra Nevada. *Journal of*
1172 *Hydrometeorology*, 18(4), 1101–1119. <https://doi.org/10.1175/JHM-D-16-0246.1>

1173 Wrzesien, M. L., Durand, M. T., Pavelsky, T. M., Kapnick, S. B., Zhang, Y., Guo, J., & Shum, C. K.
1174 (2018). A New Estimate of North American Mountain Snow Accumulation From
1175 Regional Climate Model Simulations. *Geophysical Research Letters*, 45(3), 1423–
1176 1432. <https://doi.org/10.1002/2017GL076664>

1177 Wrzesien, M. L., Durand, M. T., & Pavelsky, T. M. (2019). A Reassessment of North American
1178 River Basin Cool-Season Precipitation: Developments From a New Mountain
1179 Climatology Data Set. *Water Resources Research*, 55(4), 3502–3519.
1180 <https://doi.org/10.1029/2018WR024106>

1181 Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., et al. (2012). Continental-
1182 scale water and energy flux analysis and validation for the North American Land
1183 Data Assimilation System project phase 2 (NLDAS-2): 1. Intercomparison and
1184 application of model products: water and energy flux analysis. *Journal of Geophysical*
1185 *Research: Atmospheres*, 117(D3). <https://doi.org/10.1029/2011JD016048>

1186 Xia, Y., Mocko, D., Huang, M., Li, B., Rodell, M., Mitchell, K. E., et al. (2017). Comparison and
1187 Assessment of Three Advanced Land Surface Models in Simulating Terrestrial Water
1188 Storage Components over the United States. *Journal of Hydrometeorology*, 18(3),
1189 625–649. <https://doi.org/10.1175/JHM-D-16-0112.1>

1190 Yapo, P. O., Gupta, H. V., & Sorooshian, S. (1998). Multi-objective global optimization for
1191 hydrologic models. *Journal of Hydrology*, 204(1), 83–97.
1192 [https://doi.org/10.1016/S0022-1694\(97\)00107-8](https://doi.org/10.1016/S0022-1694(97)00107-8)

1193 Zeng, X., Broxton, P., & Dawson, N. (2018). Snowpack Change From 1982 to 2016 Over
1194 Conterminous United States. *Geophysical Research Letters*, 45(23), 12,940–12,947.
1195 <https://doi.org/10.1029/2018GL079621>

1196

1197 **10. Tables**

1198

1199 Table 1. Calibration parameters including default values, calibration range, and average
 1200 calibrated value. The calibration range reference, when applicable, is noted. Otherwise, the
 1201 calibration range is +/- 20% of the default value.
 1202

Parameter	Description [units]	Default Value	Calibration Range	Calibrated Average Value	Reference
ALBDRY1	Dry soil albedos (VIS) [-]	0.10-0.27	0.08-0.32	0.200	--
ALBDRY2	Dry soil albedos (NIR) [-]	0.20-0.54	0.16-0.65	0.404	--
ALBICE1	Albedo land ice (VIS) [-]	0.55	0.44-0.66	0.552	--
ALBICE2	Albedo land ice (NIR) [-]	0.80	0.64-0.96	0.795	--
ALBSAT1	Saturated soil albedos (VIS) [-]	0.05-0.15	0.04-0.18	0.110	--
ALBSAT2	Saturated soil albedos (NIR) [-]	0.10-0.30	0.08-0.36	0.224	--
BETADS	Two stream parameter β_d for snow [-]	0.5	0.4-0.6	0.498	--
BETAIS	Two stream parameter β_i for snow [-]	0.5	0.4-0.6	0.501	--
EG1	Emissivity soil surface (soil) [-]	0.97	0.78-1.0	0.895	--
EG2	Emissivity soil surface (lake) [-]	0.98	0.78-1.0	0.891	--
MFSNO	Snowmelt parameter [-]	2.5	0.5-3.0	1.158	Niu & Yang (2007)
MNSNALB	Minimum snow albedo [-]	0.55	0.45-0.65	0.596	Aguado (1985); Dirmhirn & Eaton (1975)

MXSNALB	Maximum snow albedo [-]	0.84	0.75-0.95	0.853	Aguado (1985); Essery & Etchevers (2004)
OMEGAS1	Two stream parameter omega for snow [-]	0.8	0.64-0.96	0.802	--
OMEGAS2	Two stream parameter omega for snow [-]	0.4	0.32-0.48	0.402	--
RSURF_SNOW	Surface resistance for snow [s/m]	50.0	40.0-60.0	49.851	--
SNDECAYEXP	Exponent in snow decay albedo relationship [h ⁻¹]	0.01	0.001-0.10	0.0338	Essery & Etchevers (2004)
SSI	Liquid water holding capacity for snowpack [m ³ /m ³]	0.03	0.01-0.08	0.0398	Amorocho & Espildora (1966); Anderson (1973)
SWEMX	New snow mass to fully cover old snow [mm]	1.0	0.5-5.0	2.280	Xia et al. (2012)
T_LLIMIT	Lower temperature limit in rain-snow partitioning [C]	0.5	0.0-2.0	0.707	--
T_MLIMIT	Middle temperature limit in rain-snow partitioning [C]	2.0	0.5-3.0	1.724	--
T_ULIMIT	Upper temperature limit in rain-snow partitioning [C]	2.5	1.0-5.0	3.393	--

ZOSNO	Snow surface roughness length [m]	0.002	0.0001-0.01	0.00298	Marks & Dozier (1992); Reba et al. (2014)
SNOWF_SCALEF	Snowfall scale factor [-]	N/A	0.1-10.0	1.159	--

1203
1204
1205
1206
1207

Table 2. Details of six evaluation points, including location, elevation, and percent tree canopy cover.

Evaluation Point	Latitude	Longitude	Elevation (m)	Tree Canopy Cover (%) [*]
Senator Beck	37.91° N	107.73° W	3721	14
Niwot	40.03° N	105.58° W	3185	79
Fool Creek	39.87° N	105.87° W	3400	89
Cameron Pass	40.52° N	105.89° W	3129	83
Rock Creek	38.98° N	107.03° W	3395	19
Skyway/Grand Mesa	39.05° N	108.06° W	3245	71

1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229

^{*} Tree canopy cover calculated from Landsat 7 ETM+ data at 30 m spatial resolution (Hansen et al., 2013).

Table 3. Snow depth bias and RMSE for calibrated and uncalibrated Noah-MP simulations compared to UA and SNODAS for six SnowEx field site locations and the full western

1230 Colorado domain. Bold indicates better performance, and for the overall domain
 1231 comparisons, an asterisk (*) indicates a statistically significantly difference between the
 1232 two model performances

Evaluation Point	UA				SNODAS			
	Noah-MP-Cal Snow Depth Bias (m)	Noah-MP-Def Snow Depth Bias (m)	Noah-MP-Cal Snow Depth RMSE (m)	Noah-MP-Def Snow Depth RMSE (m)	Noah-MP-Cal Snow Depth Bias (m)	Noah-MP-Def Snow Depth Bias (m)	Noah-MP-Cal Snow Depth RMSE (m)	Noah-MP-Def Snow Depth RMSE (m)
Senator Beck	0.247	-0.183	0.418	0.315	0.163	-0.267	0.349	0.428
Niwot	- 0.0663	-0.245	0.211	0.387	- 0.0398	-0.218	0.279	0.397
Fool Creek	-0.131	- 0.0769	0.229	0.207	-0.327	-0.119	0.484	0.227
Cameron Pass	-0.104	-0.142	0.211	0.254	- 0.0267	-0.0647	0.170	0.197
Rock Creek	-0.180	-0.251	0.308	0.400	-0.187	-0.259	0.346	0.428
Skyway/ Grand Mesa	- 0.0229	-0.176	0.173	0.320	-0.100	-0.253	0.246	0.425
Overall Domain	- 0.00229*	-0.0362	0.131*	0.146	- 0.00980*	-0.0437	0.186	0.179*

1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241
 1242
 1243
 1244
 1245

1246 Table 4. Comparison of seasonal bias and RMSE of snow depth for evaluation points from Noah-MP-Def and Noah-MP-Cal for
 1247 the accumulation (December-February), peak snow (March-April), and ablation (May-July) seasons. Bold indicates better
 1248 performance.

Point	Simulation	Accumulation Season Bias (cm)	Accumulation Season RMSE (cm)	Peak Season Bias (cm)	Peak Season RMSE (cm)	Ablation Season Bias (cm)	Ablation Season RMSE (cm)
Cameron Pass	Noah-MP-Def	-15.51	19.43	-30.35	33.86	-19.63	37.66
	Noah-MP-Cal	-11.15	16.17	-24.77	29.12	-14.14	30.49
Fool Creek	Noah-MP-Def	7.05	15.10	15.60	29.76	11.73	29.71
	Noah-MP-Cal	-17.89	23.37	-31.59	39.20	-11.55	22.69
Niwot Ridge	Noah-MP-Def	-34.03	40.62	-51.37	57.61	-26.02	45.72
	Noah-MP-Cal	-8.08	21.59	-11.68	30.11	-10.71	26.11
Rock Creek	Noah-MP-Def	-40.47	47.98	-64.26	68.27	-15.64	31.85
	Noah-MP-Cal	-22.79	30.72	-46.41	51.51	-17.59	32.96
Senator Beck	Noah-MP-Def	-28.24	35.50	-49.51	54.09	-12.79	27.68
	Noah-MP-Cal	27.11	36.84	32.71	49.86	41.22	61.32
Skyway (Grand Mesa)	Noah-MP-Def	-15.82	20.26	-48.90	55.45	-19.85	40.17
	Noah-MP-Cal	4.42	12.22	-11.97	25.51	-6.20	24.41

1249
 1250
 1251
 1252

1253 Table 5. Error metrics for snow depth and SWE comparing Noah-MP-Cal and Noah-MP-Def
 1254 with snow pit observations from Grand Mesa and Senator Beck SnowEx study sites from
 1255 the February 2017 field campaign. Bold indicates better performance between the two
 1256 Noah-MP configurations.
 1257

Study Site	Simulation	Snow Depth				SWE			
		Mean Bias (cm)	MAE (cm)	Mean % diff (abs. value)	RMSE (cm)	Mean Bias (mm)	MAE (mm)	Mean % diff (abs. value)	RMSE (mm)
Grand Mesa	Noah-MP-Cal	-12.1	28.1	20.0%	34.9	-23.0	106.4	23.6%	132.9
	Noah-MP-Def	-48.2	48.8	32.2%	54.4	-160.6	162.9	32.6%	185.4
	SnowModel	-25.3	34.6	23.8%	41.1	-35.5	88.5	20.2%	112.2
Senator Beck	Noah-MP-Cal	92.6	92.6	84.4%	102.5	386.9	388.4	111.5 %	413.0
	Noah-MP-Def	-13.1	43.1	34.0%	49.7	-32.1	142.5	37.6%	167.6

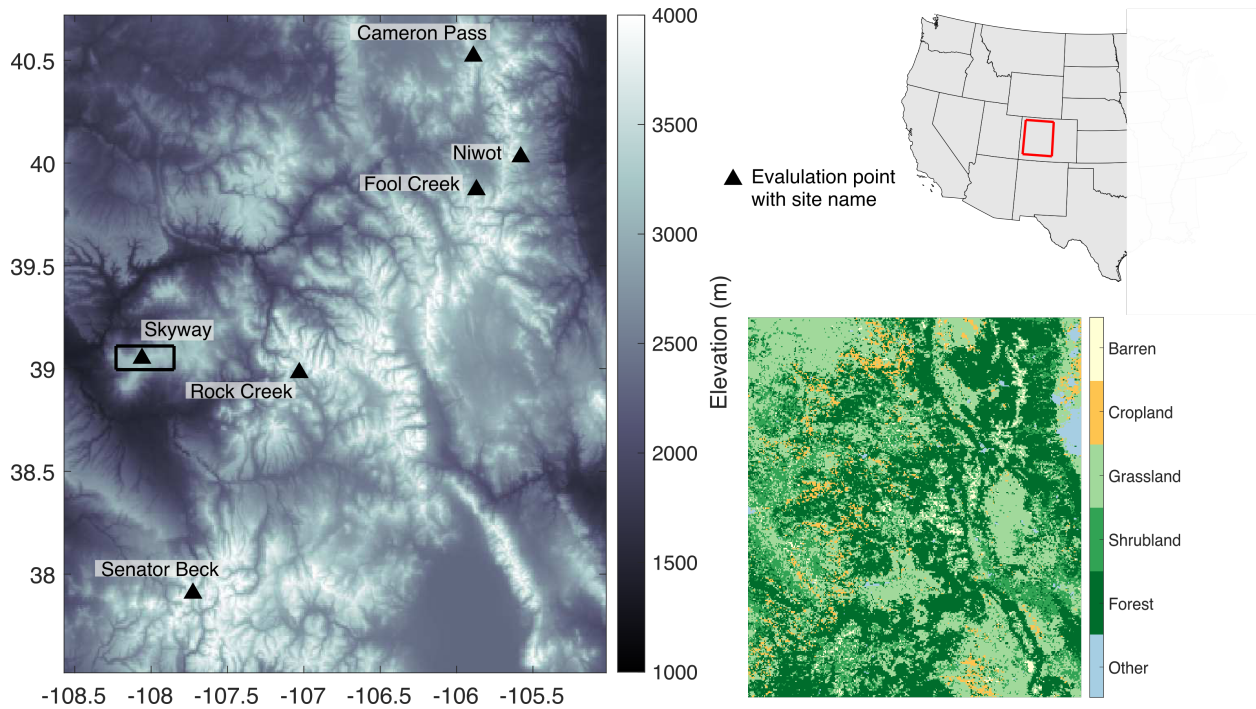
1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282

1283 Table 6. Nash-Sutcliffe Efficiency values for calibrated and uncalibrated Noah-MP
 1284 simulations for six subbasins, as described by their USGS streamgage ID. Included in the
 1285 number of Noah-MP grid cells within each subbasin. Bold indicates better performance.
 1286 Asterisk indicates where monthly streamflow difference between Noah-MP-Cal and Noah-
 1287 MP-Def is statistically significant at the 95% confidence level.

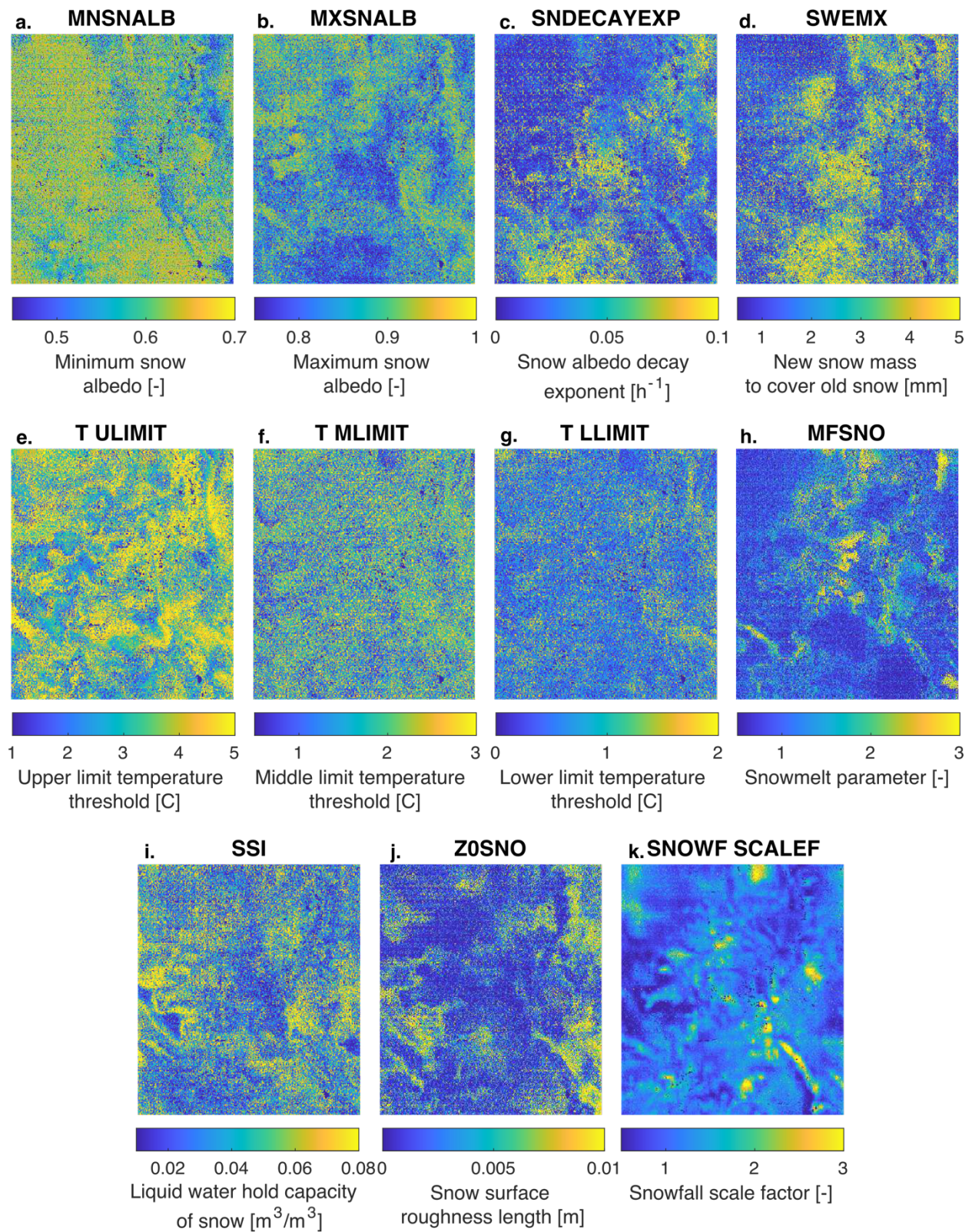
Basin	Basin ID	Noah-MP-Cal NSE	Noah-MP-Def NSE	Basin Area (km ²)
Colorado River at Glenwood Springs, CO	9072500	0.56*	-0.065	11784
Taylor River below Taylor Park Reservoir, CO	9109000	-1.03	-1.96	656
Gunnison River below Blue Mesa Dam, CO	9124700	-3.62	-0.71	8933
Gunnison River below Crystal Reservoir, CO	9127800	-2.95	-0.59	10264
Crystal River above Avalanche Creek, CO	9081600 (CAMELS)	0.72	0.43	436
Taylor River at Taylor Park, CO	9107000 (CAMELS)	-0.13	-1.14	331

1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295
 1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308

1309
1310 **11. Figures**
1311
1312



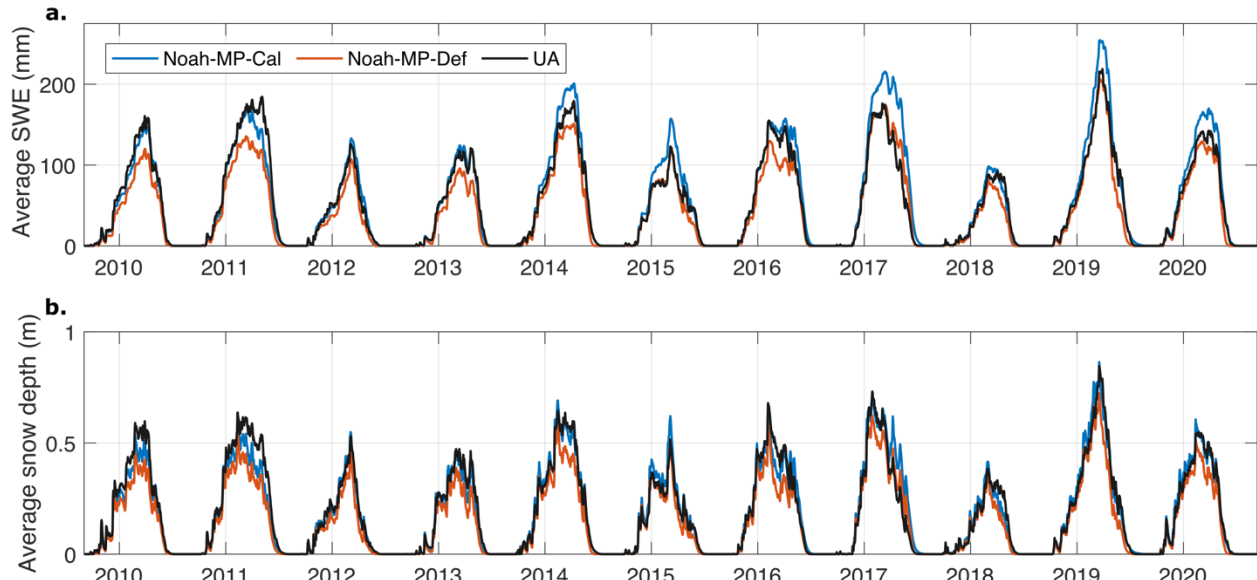
1313
1314 Figure 1. Elevations of western Colorado Noah-MP domain. Black box indicates the Grand
1315 Mesa intensive observation period field site from the NASA SnowEx 2017 field campaign.
1316 Triangles mark the six evaluation points and are labeled with the evaluation site name. The
1317 inset map shows the western Colorado domain with respect to the western United States.
1318 The bottom right plot shows the land classes for the model domain.
1319



1320
1321
1322
1323

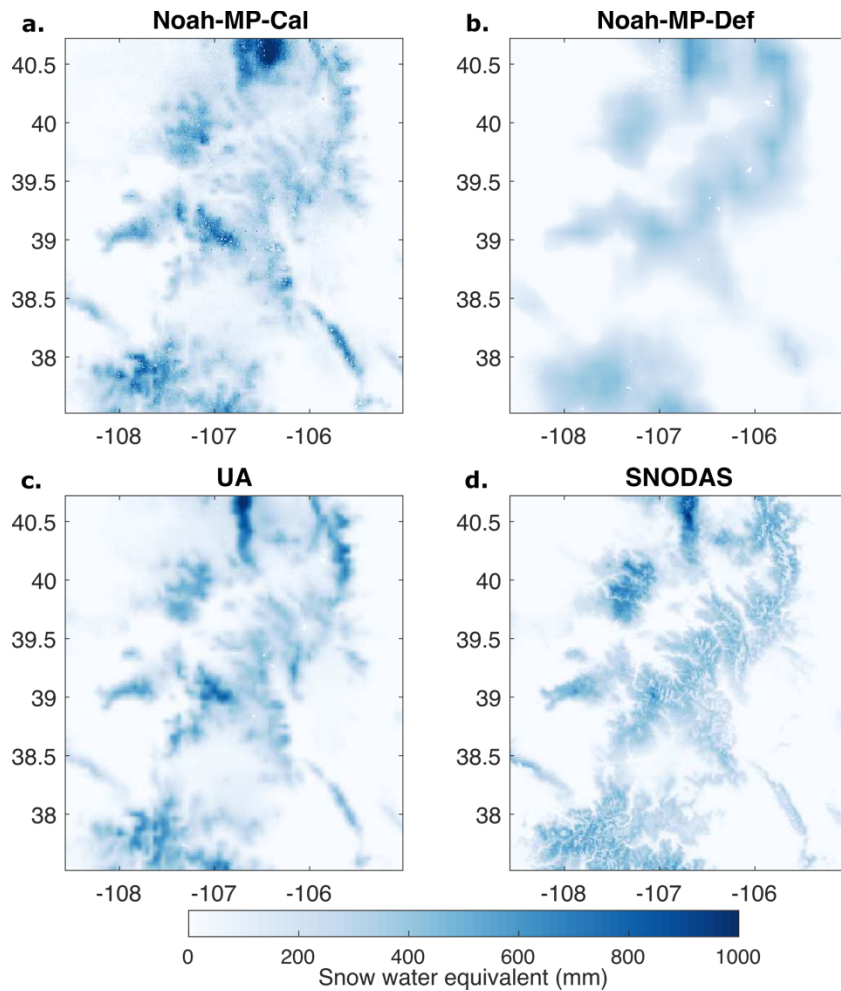
Figure 2. Calibrated parameters after the genetic algorithm procedure. Shown here are the 11 parameters that are most sensitive to calibration.

1324
1325
1326
1327

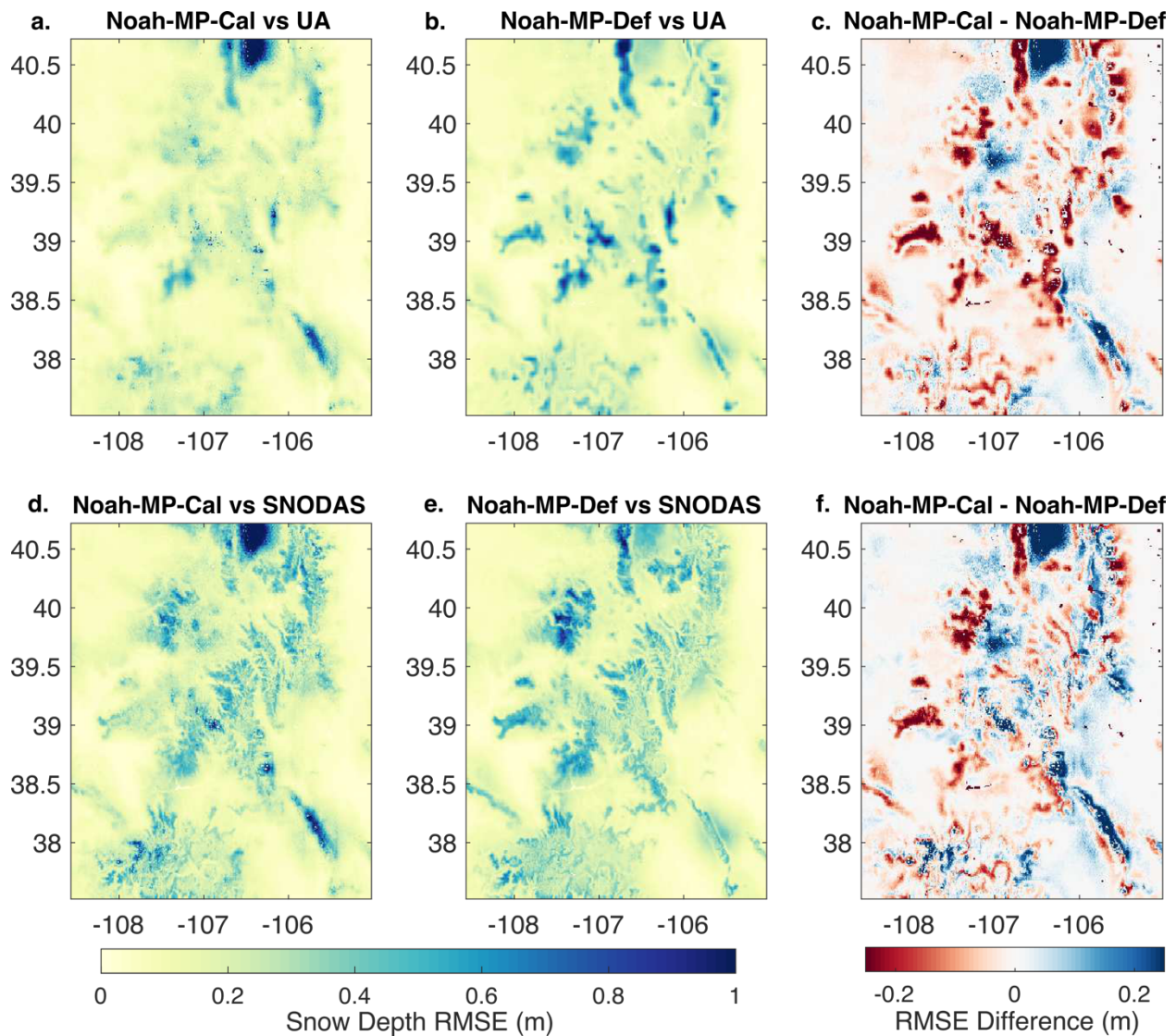


1328
1329
1330
1331
1332
1333
1334
1335

Figure 3. Time series of average SWE, in mm, and average snow depth, in m, over the full domain for calibrated (blue), uncalibrated (orange), and UA (black) estimates. Note that at the time of this analysis, UA ends in 2017 but the Noah-MP simulations continue through 2020.

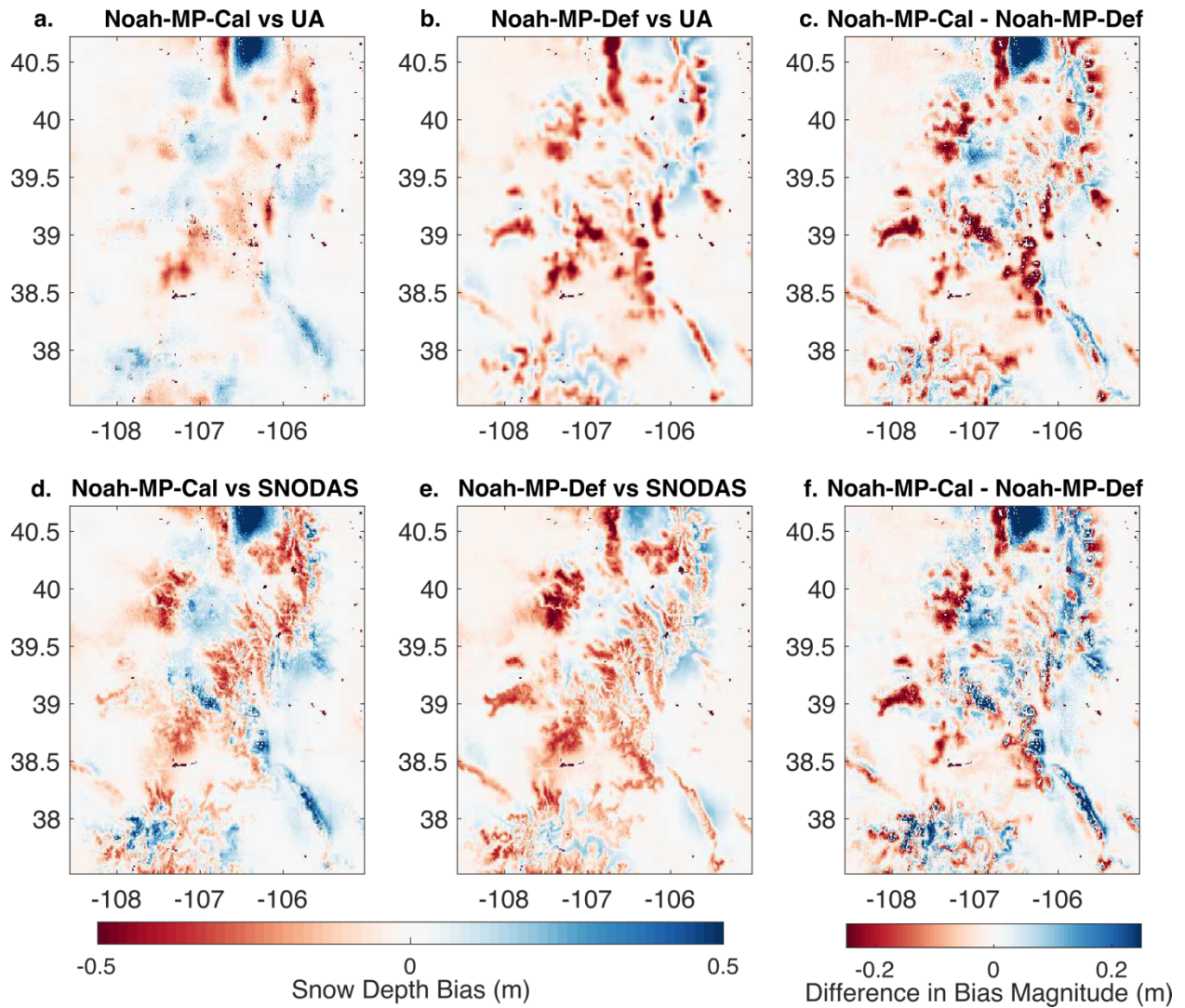


1336
 1337 Figure 4. Average April 1 SWE, in mm, for the calibrated (a) and uncalibrated (b)
 1338 simulations. (c) April 1 SWE difference, where blue indicates grid cells where the calibrated
 1339 simulation has larger SWE and red indicates where the uncalibrated simulation has larger
 1340 SWE.
 1341
 1342

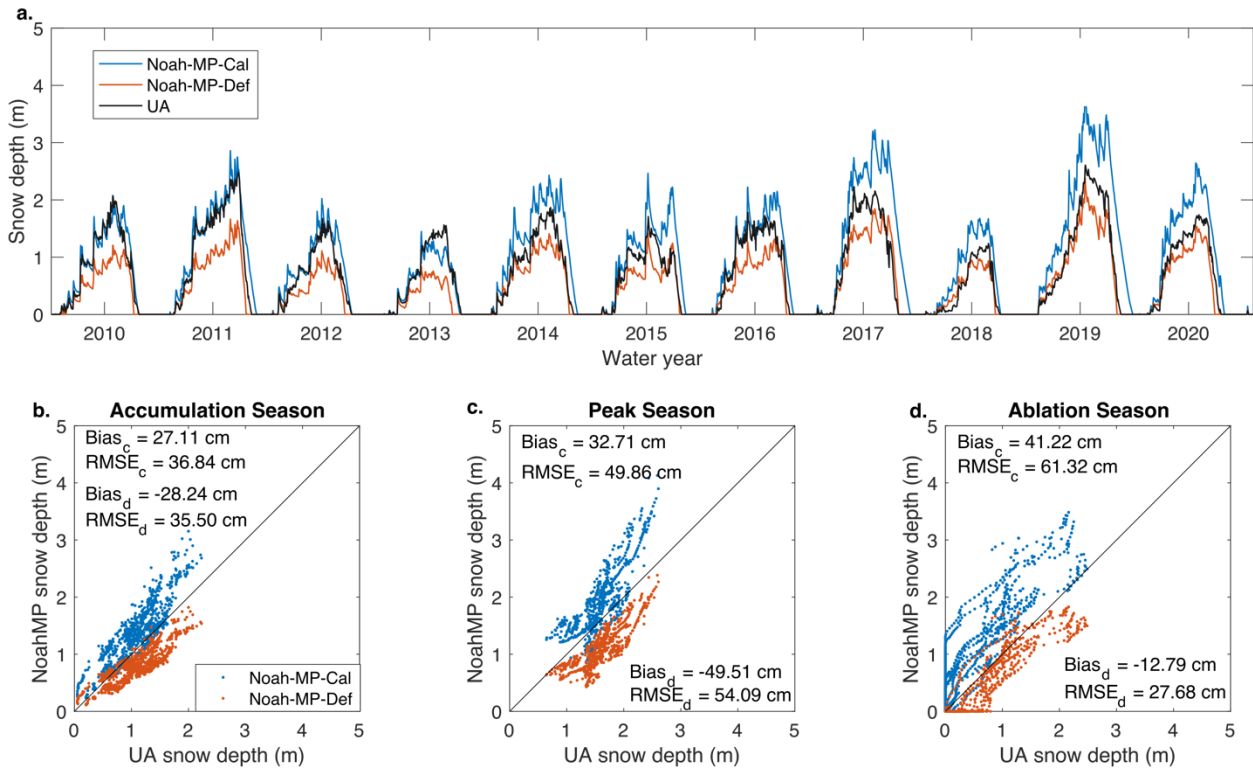


1343
 1344
 1345
 1346
 1347
 1348
 1349

Figure 5. (top row) Snow depth RMSE for Noah-MP-Cal and Noah-MP-Def compared to UA for the full analysis period. The right column shows the difference in RMSE values between Noah-MP-Cal and Noah-MP-Def. (bottom row) Same as top row except compared against SNODAS.

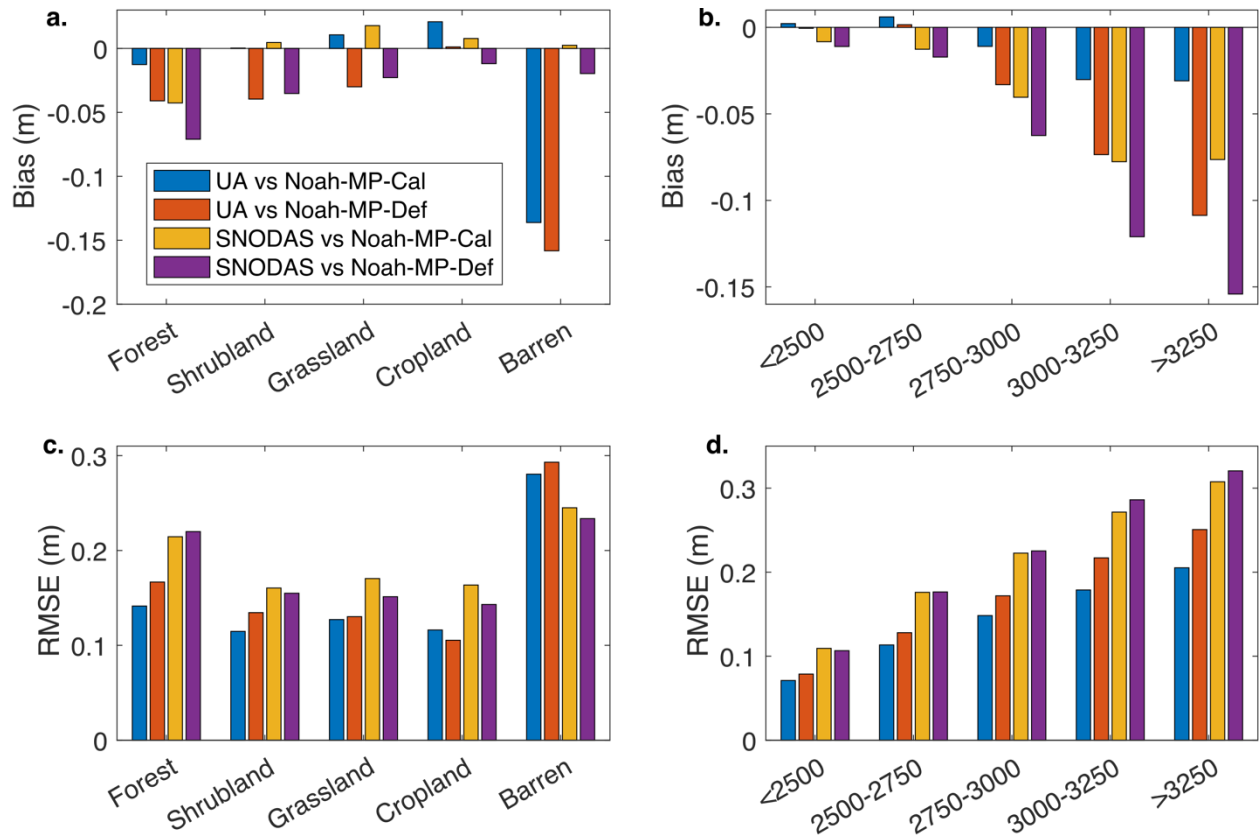


1350
 1351 Figure 6. As for Figure 5 but for snow depth bias.
 1352
 1353

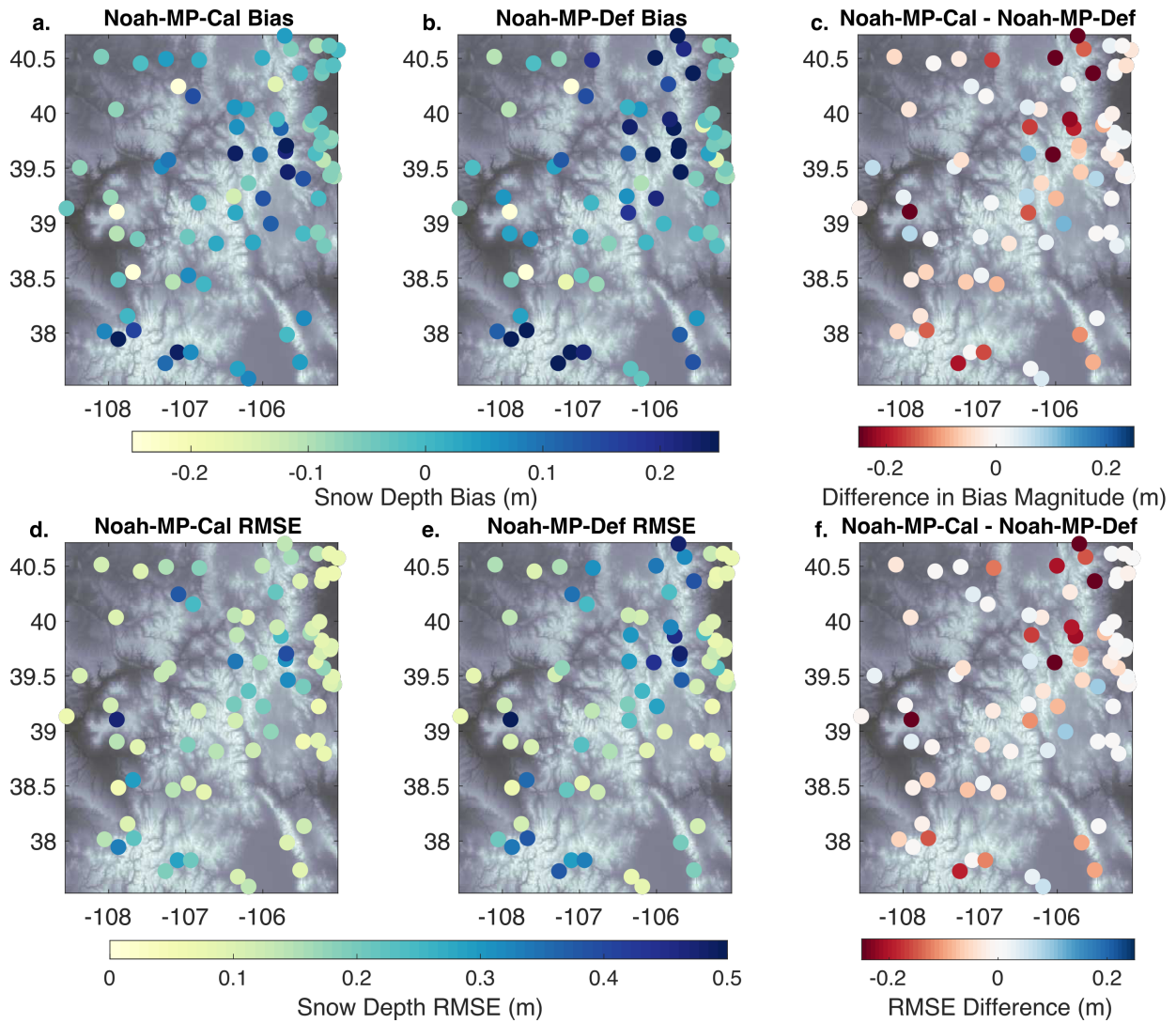


1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362

Figure 7. Evaluation of calibrated and uncalibrated Noah-MP over a single point in the Senator Beck basin. (a) Time series of daily snow depth over the grid cell that contains the Senator Beck study site. (bottom row) Scatter plot of Noah MP simulated snow depth versus UA snow depth for both calibrated (blue) and uncalibrated (orange) simulations, separated in accumulation (b; December-February), peak (c; March-April), and ablation (d; May-July) seasons.

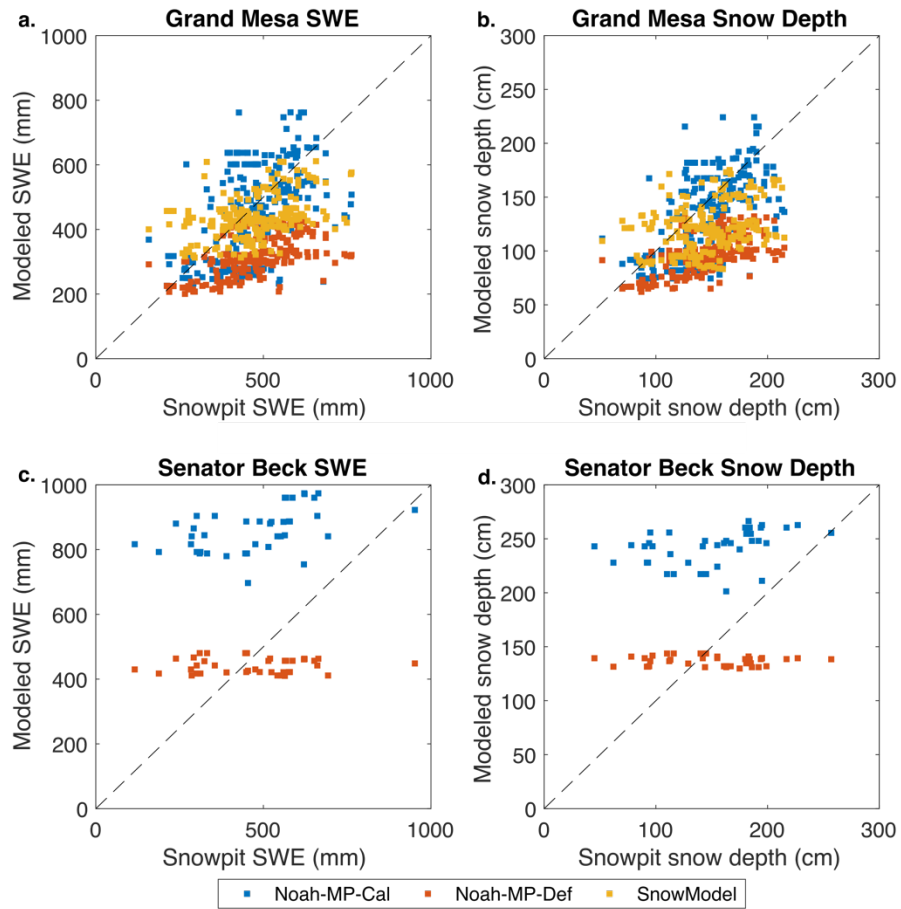


1363
 1364 Figure 8. Snow depth bias (a) and RMSE (c) from Noah-MP-Cal and Noah-MP-Def compared
 1365 to UA and SNODAS for five aggregated land cover categories. Snow depth bias (b) and
 1366 RMSE (d) for the forest land cover category separated into elevation bands. Bias and RMSE
 1367 values are temporal averages from the full analysis period.
 1368
 1369
 1370
 1371
 1372
 1373



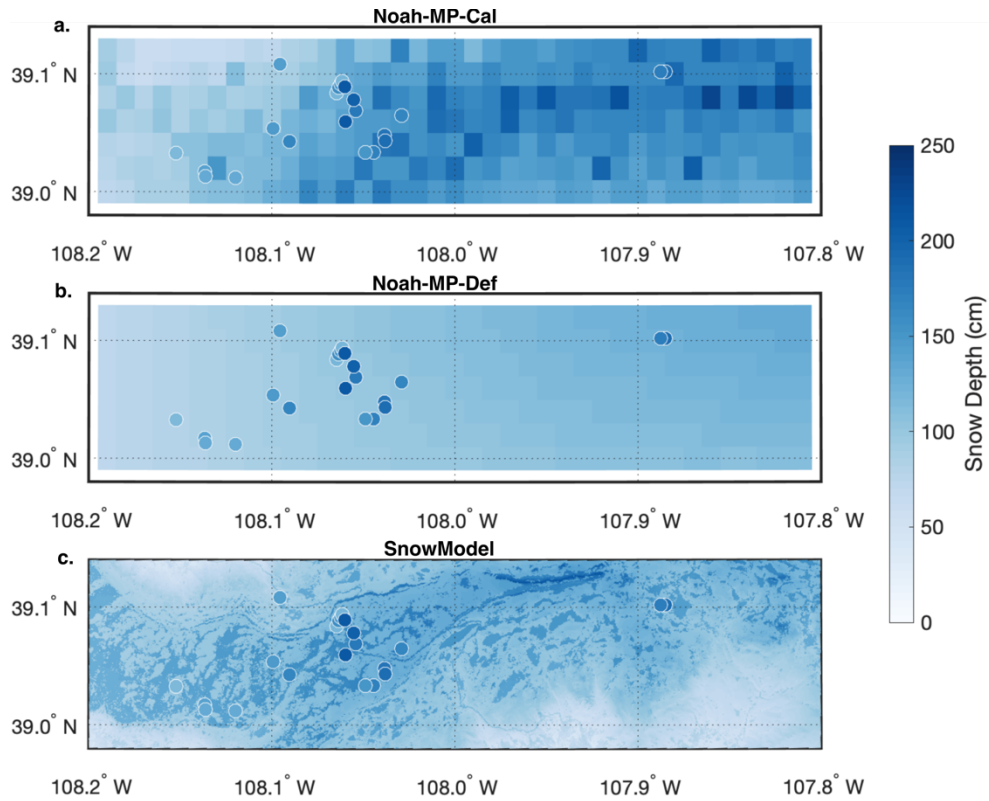
1374
 1375
 1376
 1377
 1378
 1379
 1380

Figure 9. Snow depth bias (top row) and RMSE (bottom row) from Noah-MP-Cal and Noah-MP-Def compared to GHCN station observations. The right column shows the difference between Noah-MP-Cal and Noah-MP-Def for bias and RMSE.



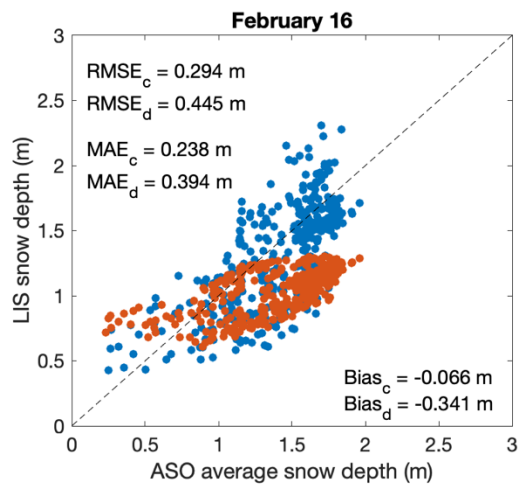
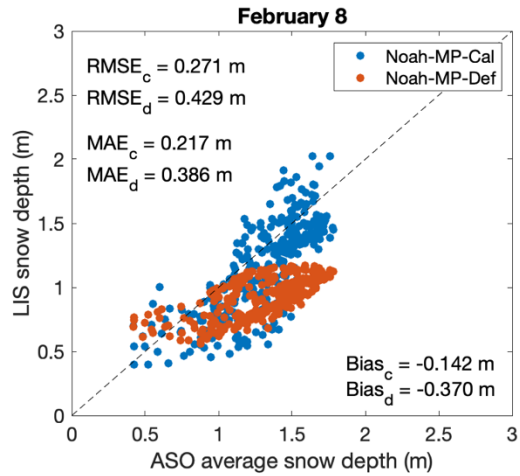
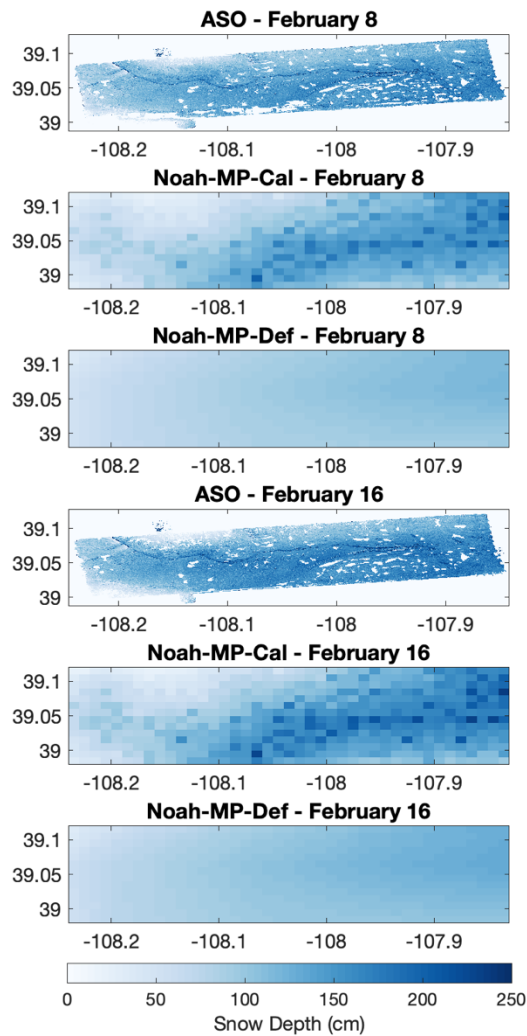
1381
 1382 Figure 10. Comparison of SWE and snow depth between Noah-MP and SnowModel
 1383 simulations and observations from snow pits during the SnowEx 2017 field campaign. In all
 1384 plots, blue squares are calibrated Noah-MP, orange squares are default Noah-MP, and
 1385 yellow squares are SnowModel (at native 30 m resolution). Plots (a) and (b) compare snow
 1386 pit measurements for Grand Mesa and plots (c) and (d) compare for Senator Beck.

1387
 1388
 1389
 1390
 1391



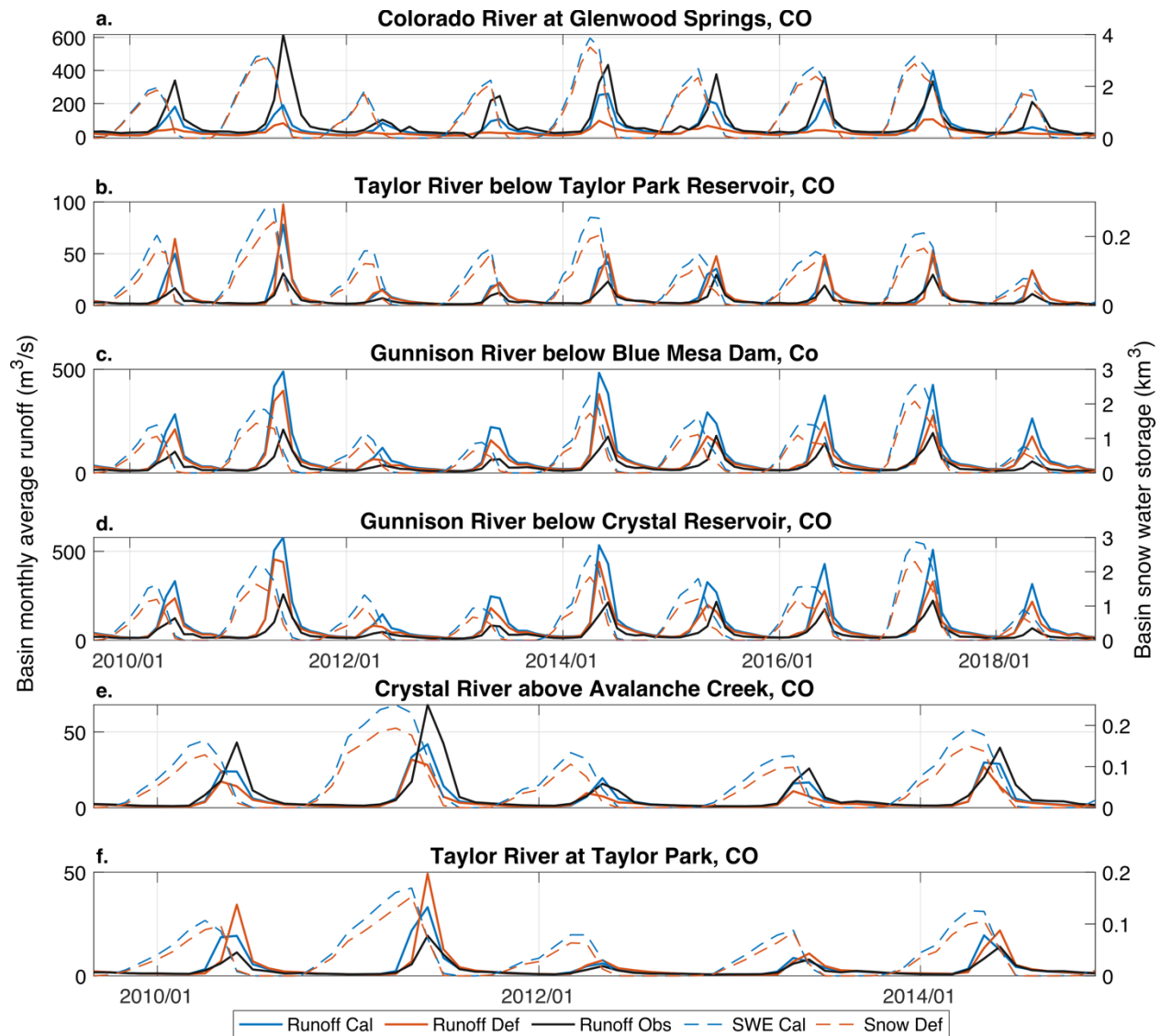
1392
 1393 Figure 11. Comparison of Noah-MP-Cal (a), Noah-MP-Def (b), and SnowModel (c) snow
 1394 depth estimates with snow pit observations for February 22, 2017, over the SnowEx Grand
 1395 Mesa field campaign site. SnowModel is shown at 30 m spatial resolution. Snow pit depths
 1396 and model depths are on the same color scale.

1397
 1398
 1399
 1400



1401
 1402 Figure 12. Comparison of snow depth from ASO flights with Noah-MP-Cal and Noah-MP-Def
 1403 over Grand Mesa for February 8 and 16, 2017. Spatial maps are all at native resolution: 3 m
 1404 for ASO and 0.01° for Noah-MP simulations. Scatter plots compare Noah-MP simulations to
 1405 ASO observations, where ASO has been upscaled to 0.01° resolution.

1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414



1415
 1416
 1417
 1418
 1419
 1420
 1421

Figure 13. (a-d) Comparison of runoff, in m^3 /month, between Noah-MP simulations and estimates of naturalized flow for four subbasins in the Upper Colorado River Basin. (e-f) Comparison of runoff, in m^3 /month, between Noah-MP simulations and observed streamflow from USGS streamgages for small unmanaged subbasins, selected from the CAMELS database. Streamgauge locations are shown on Figure S4. Dashed lines in all plots show basin snow water storage, in km^3 .