

1
2
3 **Dynamical downscaling improves upon**
4 **gridded precipitation products in the Sierra**
5 **Nevada, California**
6

7 **Mimi Hughes¹, Jessica D. Lundquist², and Brian Henn³**

8 ¹University of Colorado, Cooperative Institute for Research in Environmental Science, and
9 NOAA/ESRL/PSD.

10 ²University of Washington, Civil and Environmental Engineering.

11 ³Center for Western Weather and Water Extremes, Scripps Institution of Oceanography,
12 University of California, San Diego, La Jolla, California, USA

13
14 Corresponding author: Mimi Hughes (mimi.hughes@noaa.gov); phone: (303)497-4865; fax:
15 (303)497-6101

16
17 Keywords: precipitation, WRF, convection-permitting, downscaling

18
19 In preparation for *Climate Dynamics* (special issue on high-resolution climate modeling)
20

21

22 **Abstract**

23 Uncertainties in gridded and regional climate estimates of precipitation are large at high
24 elevations, where observations are sparse and spatial variability is substantial. We explore these
25 uncertainties for water year 2008 across California's Sierra Nevada in 10 datasets: six regional
26 climate downscalings generated using the Weather, Research, and Forecast (WRF) model at
27 convection-permitting resolution with differing lateral boundary conditions and microphysical
28 parameterizations, and four gauge-based, interpolation-gridded precipitation datasets.
29 Precipitation from these 10 datasets is evaluated against 95 snow pillows and a precipitation
30 dataset inferred from stream gauges using a Bayesian inference method. During water year 2008,
31 the gridded datasets tend to underestimate frozen precipitation on the windward slope of the
32 Sierra Nevada, particularly in the vicinity of Yosemite National Park. The WRF simulations with
33 single-moment microphysics tend to overestimate precipitation throughout much of the region,
34 whereas the WRF simulations with double-moment microphysics tend to better agree with both
35 the snow pillows and inferred precipitation estimates, although they somewhat overestimate the
36 windward/leeside precipitation contrast in the northern Sierra Nevada. WRF simulations, in
37 particular those with single-moment microphysics, better distinguish spatial patterns of wet-
38 versus-dry pillows and watersheds over the water year than the gridded estimates. Our results
39 suggest treating gauge-based datasets as 'truth' may give a misleading representation of model
40 accuracy, since these gauge-based datasets often have issues of their own.
41

42

43 **1 Introduction**

44 Quantitative precipitation estimates in mountainous areas are essential for hydrologic
45 modeling and water management. Despite being critical, accurate precipitation estimates are
46 notoriously difficult to produce in areas of complex terrain due to limitations of observing
47 systems and large spatial variability in these regions. Ground-based radars can provide reliable
48 estimates of precipitation over regions with homogeneous topography, but suffer from beam
49 blocking and other issues in complex terrain (Nelson et al. 2016; Willie et al 2016). Large spatial
50 variability, due to meteorological response to terrain features, makes it difficult to translate point
51 measurements into gridded estimates, a problem confounded by the difficulty of attaining a
52 dense network of point measurements in the complex terrain. Thus hydrologists often turn to
53 gridded precipitation datasets, created either through statistical methods applied to in situ data
54 (hereafter gridded estimates) or through dynamical downscaling of reanalysis datasets with
55 atmospheric models.

56 Several daily statistically gridded precipitation estimates exist for the continental United
57 States at resolutions as fine as 1 km (e.g., see Table 1 in Lundquist et al. 2015, hereafter L15).
58 These datasets interpolate gauge data to a grid using numerical methods; the majority of these
59 scale their daily values such that their long-term monthly means match the Precipitation-
60 elevation Regressions on Independent Slopes Model (PRISM; Daly et al. 1994) long-term
61 monthly means. Because they rely directly on in situ data, these datasets provide an estimate of
62 precipitation with uncertainties commensurate with the uncertainties in the in situ data
63 themselves where in situ data are dense, and with greater uncertainty where in situ data are
64 sparse, such as at high elevations.

65 With the advent of more computing power, ever-improving atmospheric models, and
66 better-constrained atmospheric reanalyses, dynamical downscaling of reanalysis datasets to
67 convection-permitting resolution has become a viable, if more computationally expensive,
68 alternative to statistical downscaling for high-resolution precipitation estimation. Dynamical
69 downscaling uses a state-of-the-art numerical weather model to estimate the atmospheric state at
70 high resolution given prescribed large-scale conditions (e.g., see review by Xue et al. 2014). To
71 generate high resolution historical data, the prescribed large-scale conditions are generally from
72 reanalysis datasets. Because dynamical downscaling uses a numerical weather model to calculate
73 precipitation based on discretized equations of state, it can represent physical processes not
74 captured by the linear precipitation-elevation regressions used by PRISM.

75 The skill of dynamical downscaling to accurately represent reality is a function of the
76 model physics and accuracy of the large scale conditions, and has the largest potential to improve
77 over coarser scale reanalyses in areas of complex terrain and coastlines (Xue et al. 2014; Feser et
78 al. 2010). Prior studies have shown that microphysical parameterizations have large impacts on
79 precipitation type/phase in the cloud and on the ground (Minder and Kingsmill 2013 and Jankov
80 et al. 2007, 2009, 2011), and that precipitation sensitivity to microphysical scheme is often larger
81 than to other physics parameterizations (Liu et al. 2011). Although not as well studied as
82 sensitivity to physics parameterizations, dynamical downscalings are also sensitive to their
83 lateral boundary conditions (Yang et al. 2012), so uncertainties in reanalysis datasets create
84 uncertainties in the downscaled data.

85 Evaluating high-resolution gridded precipitation datasets is challenging: Statistically
86 gridded datasets often ingest most available in situ precipitation gauge data, making determining
87 their skill versus that of dynamical downscalings tricky. In fact, many studies use PRISM-based
88 gridded datasets to evaluate their dynamical downscalings (e.g. Caldwell et al. 2009), despite the
89 uncertainties in the gridded datasets. These uncertainties can be large: For example, when
90 compared with independent precipitation observations, PRISM-based gridded datasets have been
91 shown to be off by up to a factor of two in mountainous terrain (Gutmann et al. 2012; Jeton et al
92 2006). These uncertainties also have important implications for water resources because they can
93 significantly impact water year total precipitation amounts: Henn et al. (2016b) showed using
94 different gridded datasets in complex terrain could cause up to 40% differences in water year
95 totals.

96 Two data sources that have been largely untapped for precipitation dataset evaluation are
97 stream gauges and snow pillows. Streamflow datasets are not used in any of the gridded datasets,
98 to our knowledge, and thus serve as a completely independent dataset representing basin-mean
99 precipitation; however, using them to validate precipitation estimates requires accounting for
100 hydrologic processes that carry the precipitation to the streams. Henn et al. (2016a; hereafter
101 H16) infer basin-mean precipitation from 56 streamflow gauges in the Sierra Nevada with
102 Bayesian modeling, providing a streamflow-inferred precipitation dataset independent from
103 gridded estimates of precipitation; this dataset includes its own range of uncertainty as part of the
104 Bayesian modeling technique. Snow pillows offer the second widely unused in situ dataset for
105 gridded precipitation validation: L15 quality controlled 20 years of snow pillow data across the
106 Sierra Nevada from the California Department of Water Resources (CA DWR); some of these
107 pillows and nearby snow courses have been used to adjust PRISM climatologies (L15), but the
108 daily SWE amounts are not used directly in any of the gridded datasets.

109 L15 evaluated two gridded datasets, Hamlet (Hamlet et al. 2010; Hamlet and Lettenmaier
110 2005) and Livneh (Livneh et al. 2013), against the CA DWR snow pillows and found median
111 errors for the entire period of $\pm 10\%$. L15 further identified that during some years the two
112 gridded datasets severely underpredicted precipitation during individual storms by as much as
113 50%, leading to water year total errors of $\sim 20\%$; water year 2008 (WY2008) showed the most
114 severe underestimation in their 20-year study period, with two major storms having substantial
115 snow underestimation in each dataset. L15 showed these errors occurred because of an increase
116 in orographic precipitation gradient when the winds were more westerly/northwesterly than
117 typical during precipitating days, and hypothesized that this shift was associated with more time
118 spent in the storm's cold sector.

119 In this manuscript, we extend upon the results of L15 for WY2008, by comparing the
120 snow pillow observations to four gridded datasets -- Hamlet, Livneh, Daymet (Thornton et al.
121 1997; Thornton et al. 2014), and Newman (Newman et al. 2015). We assess whether the
122 additional two gridded datasets, Daymet and Newman, also underpredict snow and total
123 precipitation at high elevations in WY2008. We hypothesize that these two datasets will suffer
124 similar biases in WY2008 as those in Hamlet and Livneh, since they use similar, terrain-based
125 interpolation techniques, although we expect some variation across the gridded datasets because
126 of differing methodological choices. In addition, we identify whether six dynamical
127 downscalings suffer from the same deficiencies in WY2008, and hypothesize that the dynamical
128 downscalings will not suffer from similar biases on the windward slope of the Sierra Nevada

129 because of their ability to represent orographic precipitation processes. We also investigate the
130 sensitivity of the dynamical downscalings to their microphysical parameterizations and lateral
131 boundary conditions to test whether the large scale uncertainties in reanalyses produce
132 differences in precipitation comparable to those from microphysics parameterizations. Testing
133 over the course of a water year (rather than performing case studies of individual storms) allows
134 us to test whether differences in precipitation due to microphysical and large scale uncertainties
135 accumulate over the course of the water year, and to see whether they are systematic or random
136 in different parts of the Sierra Nevada. Finally, we compare the 10 precipitation estimates (4
137 gridded and 6 dynamical downscalings) to inferred basin-mean precipitation from H16. The
138 manuscript is laid out as follows: Section 2 describes the datasets and methods used in the
139 manuscript; Section 3 explores the differences between the 10 precipitation estimates and the
140 snow pillow and inferred precipitation amounts; and Section 4 provides a summary and
141 discussion of the results.

142 **2 Datasets and methods**

143 2.1 WRF simulations

144 Six dynamical downscalings of WY2008 were generated using the Weather, Research,
145 and Forecast (WRF) model, version 3.6 (Skamarock et al. 2008). All simulations were initialized
146 in July 2007 and run continuously through Oct. 1, 2008, with the first three months discarded as
147 model spin-up. All six simulations were identically configured aside from the lateral boundary
148 conditions and microphysics schemes used (Table 1). Half the simulations used lateral boundary
149 conditions (LBCs) from the ERA Interim reanalysis (ERA-I; Dee et al. 2011), and half used the
150 North American Regional Reanalysis (NARR; Mesinger et al. 2006). Each LBC was paired with
151 one of three microphysics schemes: WRF Single-Moment 6-Class Scheme (WSM6; Hong and
152 Lim 2006), the Morrison et al. (2009) double-moment scheme (Morr), or the Thompson et al.
153 (2008) scheme (Thom), resulting in six simulations, which are hereafter identified as E.Morr,
154 E.Thom., E.WSM6, N.Morr, N.Thom, and N.WSM6, where E. and N. refer to ERA-I and
155 NARR, respectively. Approximate computational time per month of simulation is shown in
156 Table 1.

157 The simulations used an 18 km outer domain covering much of the intermountain west
158 and stretching west across the northeastern Pacific Ocean, with a 6 km inner domain that
159 extended across all of California (Fig. 1). Both domains used the Rapid Refresh Transfer Model
160 for GCM applications (RRTMG; Iacono et al. 2008) for shortwave and longwave radiation and
161 the Yonsei University planetary boundary layer scheme (Hong et al. 2006) with revised
162 Mesoscale Model version 5 surface layer physics (Jimenez et al. 2012). The 18 km domain used
163 the Kain-Fritsch convective parameterization (Kain 2004), while in the 6 km domain only
164 resolved convection could occur. Spectral nudging was used in the 18 km domain to prevent
165 simulation drift; nudging was applied with strength 0.0003 s^{-1} on winds and temperature above
166 the 40th model level. The simulations used 82 vertical levels.

167 A yearlong test simulation which used the above E.Morr configuration with the Noah
168 land surface model (Tewari et al. 2004) revealed extremely cold surface temperature biases that
169 developed in springtime (not shown). These biases were attributed to the representation of snow
170 within the Noah land surface model (e.g., Barlage et al. 2015; Pavelsky et al. 2011), and thus the

171 simulations used in this manuscript use the more sophisticated Noah-MP land surface model
172 (Niu et al. 2011), which eliminates the springtime biases (not shown).

173 WRF outputs total precipitation, snow, ice, and graupel. Thus for the comparison with
174 snow pillow data, the sum of snow, ice, and graupel is used as frozen precipitation, whereas for
175 the comparison with the Bayesian estimated precipitation, total precipitation is used.

176 2.2 Statistically gridded precipitation estimates

177 Four datasets that interpolate precipitation and temperature from gauge observations and
178 estimate them on a grid that extends across the Continental United States are used in this
179 manuscript.

180 As discussed in more detail in Lundquist et al. (2015), two of the gridded
181 precipitation/temperature datasets, Livneh (Livneh et al. 2013) and Hamlet (Hamlet et al. 2010;
182 Hamlet and Lettenmaier 2005) (and many other gridded datasets not shown in this manuscript),
183 use the Parameter–Elevation Regressions on Independent Slopes Model (PRISM; Daly et al.
184 1994, 2008) climatology to interpolate precipitation over topography. Both Livneh and Hamlet
185 are available on a $1/16^\circ$ grid. Both datasets used gauges from the National Climatic Data Center
186 (NCDC) Cooperative Observer (COOP) network, although they differ slightly in their criteria for
187 station inclusion. Hamlet uses PRISM to rescale temperature over topography, whereas Livneh
188 uses a constant lapse rate of $6.5\text{ }^\circ\text{C km}^{-1}$ for topographic temperature adjustment. Hamlet has no
189 data available in the northeastern quadrant of our focus region after 2006 (see greyed region in
190 Fig. 3e), but most of our comparisons focus west and south of this area.

191 The third gridded precipitation dataset used is Daymet (Thornton et al. 1997; Thornton et al.
192 2014), which combines a Gaussian weighting filter centered at the observation locations with
193 linear regression to account for elevation changes to solve for both daily gridded precipitation
194 and temperature minimum and maximum. Daymet is available on a 1km grid; prior to
195 interpolation to the WRF grid described below we smooth Daymet with a 5 km-wide centered
196 average. Daymet includes both COOP precipitation stations and stations in the U.S. Natural
197 Resources Conservation Service (NRCS) Snowpack Telemetry (SNOTEL) network.

198 The fourth gridded precipitation dataset, Newman (Newman et al. 2015), uses similar gridding
199 methodologies as the other three datasets, but differs in its inclusion of uncertainty estimates by
200 generating an ensemble of estimates following the methods of Clark and Slater (2006). Newman
201 uses distance dependent weightings from nearby stations with regression methods to generate the
202 gridded precipitation estimates, where the regression residuals are used to generate uncertainty
203 estimates. Topographic slope information was included in the regressions to account in a simple
204 way for windward and leeward slope precipitation differences. The individual Newman
205 precipitation and temperature ensemble members are available on a $1/8^\circ$ grid. Newman includes
206 more gauge data than the other datasets, including COOP stations and SNOTEL as well as
207 gauges from the Community Collaborative Rain, Hail, and Snow (CoCoRaHS) network; and the
208 various automated airport weather stations.

209 Since all four datasets are available on different grids, we interpolate them using nearest-
210 neighbor interpolation to the 6 km WRF grid prior to our analysis. In addition, for the

211 comparison to snow pillow data, daily frozen precipitation was calculated by summing
212 precipitation only on days with minimum temperature (T_{min}) less than 0°C (following L15).
213 Because of its ensemble nature, the Newman dataset required additional processing: Calculation
214 of frozen precipitation in Newman used T_{min} (computed from the dataset-native temperature
215 mean and range) for each ensemble member individually, constructing an ensemble of frozen
216 precipitation. Throughout the manuscript, when only one value is shown for Newman
217 precipitation or frozen precipitation, we are showing the ensemble median. We also show on a
218 few figures the 25th and 75th percentiles of Newman precipitation or frozen precipitation, in
219 addition to the median, to characterize the uncertainty captured by the ensemble.

220 2.3 Snow pillows

221 The CA DWR manages a network of 125 snow pillows, 103 in the Sierra Nevada (Fig. 1, data
222 available from California Data Exchange Center 2014); 95 of these pillows report enough quality
223 data in 2008 for comparison to our frozen precipitation datasets. These are generally located in
224 flat clearings and measure the weight of snow accumulating over an area of about 7 m^2 to
225 determine snow water equivalent (SWE). Because pillows can experience several hours delay
226 in responding to changes in SWE (Beaumont 1965; Johnson and Marks 2004), they are not as
227 reliable at sub-daily resolution, and thus data were analyzed at daily increments. All positive
228 daily changes in measured snow water equivalent, ΔSWE , were taken to be a measure of daily
229 snowfall. An increase in SWE was attributed to snow falling on the pillow, or to liquid water
230 falling on snow already on the pillow and freezing into the snowpack, thereby increasing its
231 density. In freezing locations where a snow pillow was co-located with a precipitation gauge,
232 the timing and amount of ΔSWE closely tracked the total accumulated precipitation. Exceptions
233 occurred where the precipitation gauge suffered severe undercatch (in those cases ΔSWE
234 exceeds measured precipitation) or during warm rain events (when rainwater passes through the
235 snowpack and drains away from the pillow, and measured precipitation exceeds ΔSWE).
236 Snowmelt and/or sublimation also may decrease SWE. Wind redistribution of snow can either
237 augment or decrease SWE, but this effect is slight because most California snow pillows are in
238 sheltered locations (Farnes 1967). In summary, snow pillows are a reliable measure of high-
239 elevation snowfall, and they do not suffer from the undercatch that standard precipitation gauges
240 suffer in such environments (Yang et al. 2005). However, because Sierra snowpacks are
241 typically warm and isothermal, most rain falling on a snow pillow is not retained and therefore,
242 not measured (Lundquist et al. 2008). All snow pillow data were quality controlled as described
243 in L15.

244 2.4 Bayesian precipitation estimates

245 In order to provide another independent estimate against which to validate the modeled
246 precipitation, we use daily streamflow observations and a method for inferring basin-mean
247 precipitation given streamflow. Streamflow observations provide an indirect representation of
248 precipitation patterns, as each basin integrates spatially-distributed precipitation inputs into the
249 streamflow response.

250 Of the 56 stream gauges identified by H16, which measure streamflow from basins that are
251 largely free of upstream diversions and flow regulation, we use a subset of 31 with data in

252 WY2008. We then apply a Bayesian methodology (Henn et al. 2015) to infer the probability
253 distribution of the basin-mean precipitation total for WY2008, given the observed streamflow in
254 each basin. The methodology uses lumped hydrologic models forced by daily precipitation time
255 series, which are scaled using multiplier parameters. These parameters, along with the other
256 hydrologic model parameters, are then inferred in Bayesian model calibration to streamflow
257 observations. Thus, the inferred precipitation from streamflow, P_{inferred} , is the WY2008
258 precipitation total that yields the best match to observed streamflow in each basin. P_{inferred} is
259 given as an ensemble resulting from using six different hydrologic model structures, in order to
260 represent the uncertainty associated with this approach. While the uncertainty of P_{inferred} is
261 substantial, we note that streamflow represents a spatially-integrated response to precipitation,
262 unlike precipitation gauge-based datasets that are derived from point measurements. In areas of
263 high spatial variability of precipitation and sparse gauge networks, streamflow-derived P_{inferred}
264 may capture aspects of this variability missed by gauge-based datasets. For more information on
265 the methodology used to infer precipitation from streamflow, see Henn et al. (2015, 2016a).

266 **3 Differences in frozen and total precipitation across datasets**

267 3.1 Snow pillow comparisons

268 3.1.1 Differences in annual frozen precipitation

269 In this section, we examine how frozen precipitation varies across the different datasets, and how
270 each dataset's frozen precipitation compares with that of the snow pillows and the multi-product
271 mean. We begin by comparing the gridded datasets to the multi-product mean (Figs. 2 and 3); the
272 multi-product mean is the average of the six WRF and four gridded datasets (the Newman
273 median frozen precipitation is used). WY2008 had ~13% lower-than-average snow totals
274 compared with the 20-year average across the pillows of L15, and ~30% lower-than-average
275 precipitation totals for the Sierra Nevada in the Sierra Nevada 8-station index (Ralph et al. 2016);
276 a larger-than-average percentage of WY2008's precipitation fell during westerly wind situations
277 (L15). There are stark differences in where the gridded datasets and WRF tend to put the largest
278 frozen precipitation amounts: WRF places more precipitation on the windward slope of the
279 Sierra Nevada, just east of the 1000 m terrain contour, with less precipitation than the multi-
280 product mean east of the crest and in the foothills. This tendency is most pronounced in the
281 WSM6 simulations, which have frozen precipitation amounts 400 to 500 mm larger on the
282 western slope than the multi-product mean, and least pronounced in EThom, followed by
283 NThom, and EMorr, respectively. All four gridded datasets have a nearly inverse pattern, with
284 less precipitation along the windward slope (up to 300-400 mm less in Livneh and Hamlet) and
285 more precipitation along the foothills and east of the crest (largest in Daymet and Newman).

286 Individual comparison of frozen precipitation at the nearest gridpoint to each snow pillow
287 observation (Fig. 4 and Fig. 5) reveals general biases of the individual datasets. When scattered
288 against the snow pillow water year totals (Fig. 4a), the reduced precipitation along the windward
289 slope in the gridded datasets shows up as a general tendency for these datasets to fall below the
290 1:1 line, especially at pillow locations with greater than 800 mm of snow in WY2008. The bias
291 of the WRF simulations depends strongly on the microphysics scheme: The two WRF
292 simulations with WSM6 microphysics have a large wet bias across much of the region (10-20%,

293 on average, and up to 50-60% at a few locations). The other four WRF simulations with the
294 double-moment microphysics schemes generally fall between the gridded datasets and the
295 WSM6 simulations, with frozen precipitation amounts that are frequently more in line with that
296 observed at the snow pillows, although large biases still exist at some locations. These overall
297 tendencies are reflected by summary statistics of the differences of the datasets' total WY2008
298 frozen precipitation with the snow pillow observations (Table 2). The WRF simulations with
299 double-moment microphysics schemes overall have mean and median differences that are closer
300 to zero than the other datasets, with the smallest median differences in EThom, NThom, and
301 EMorr, respectively. The two WRF simulations with WSM6 microphysics have mean and
302 median differences that are greater than 100 mm or ~13% of gauge-median total precipitation,
303 reflecting their wet bias. The four gridded datasets all have mean and median differences that are
304 negative, with Daymet and the Newman datasets having differences about half as large as the
305 differences of Livneh and Hamlet. The standard deviations of the differences are rather large
306 (greater than 200 mm, or ~26% of gauge-median total precipitation) for all datasets, indicating
307 the large variation of comparisons with individual snow pillows.

308 Linear fits of each dataset with the snow pillow data (Fig. 4b) make clear a few additional
309 details. First, the linear fits for all of the datasets tend to have a slope less than 1, indicating that
310 they do not have a large enough difference between the 'dry' pillows and 'wet' pillows (i.e.,
311 those pillows with a small and large annual total frozen precipitation, respectively). This
312 characteristic is generally worse in the gridded datasets – particularly Daymet – than in the WRF
313 simulations. The two WSM6 WRF simulations, despite their large wet bias, seem to suffer the
314 least of all datasets from this effect. Some of the inability to represent wet versus dry pillows
315 could be due to spatial variability occurring on scales smaller than the 6 km grid, but the variance
316 of this across datasets suggests some datasets are missing the general areas of heaviest and
317 lightest precipitation. The linear fits also reveal a systematic difference between the WRF
318 simulations with different lateral boundary conditions not easily visible in the scatterplot: The
319 NARR-forced runs tend to be slightly wetter than their ERA-I-forced counterparts, and have a
320 steeper (thus more in agreement with the snow pillows) slope, with more frozen precipitation at
321 the stations with more observed snowfall, although this effect is clearly second-order when
322 compared with the effects of microphysics.

323 To tease out the impact terrain forcing has on the distribution of frozen precipitation and errors in
324 the downscaled and gridded datasets, Fig. 5 shows the snow pillow precipitation totals along
325 with percent differences in each datasets' frozen precipitation plotted as a function of zonal
326 terrain gradient. The zonal terrain gradient in longitude/latitude space is shown for reference in
327 Fig. 5a, and has been calculated from the WRF terrain, smoothed by a 7 gridpoint filter to focus
328 on larger-scale terrain features, and multiplied by -1 to facilitate the plotting: Eastward-directed
329 gradients, i.e., those on the windward or west-facing slope, are thus negative values in Fig. 5a
330 and appear on the left half of each panel Fig. 5b-l (in agreement with the gradients in
331 longitude/latitude space), with westward-directed gradients on the right half of each panel. In the
332 northern Sierra Nevada, the snow pillows generally have considerably more precipitation on the
333 windward slope than in the lee, with annual frozen precipitation totals greater than 800 mm
334 where the gradient is sloping up to the east and less than 800 mm where it is sloping up to the
335 west. All 6 WRF downscalings tend to overdo this windward/leeside contrast, with slight
336 (double-moment runs) to moderate (WSM6 runs) positive biases at the windward locations

337 ranging in magnitude from near 0 to 60%, and dry leeside biases of up to 60% in all 6
338 simulations. The windward/leeside contrast in the southern Sierra is generally smaller in the
339 snow pillow totals, and the WRF simulations similarly do not show as consistent of a pattern in
340 their biases in this region. Unlike the WRF simulations, the biases in the gridded datasets do not
341 seem to have any strong relationship with terrain gradient.

342 3.1.2 Differences in event and daily frozen precipitation

343 Our comparison thus far has focused on differences in total WY2008 frozen precipitation, but
344 now we turn our attention to biases in daily precipitation, to get a sense for how these biases
345 evolve over the water year. To do this, we examined cumulative traces of daily snow pillow and
346 dataset snow pillow precipitation at each pillow (not shown). The behavior of the datasets'
347 frozen precipitation with respect to the snow pillows varied substantially from pillow to pillow,
348 with snow at some pillows being very over- or under-estimated by all datasets, well-represented
349 in WRF but underestimated by gridded datasets, or well-represented by most datasets. These
350 large differences at different pillows were not apparent when all pillows were lumped together
351 for error statistics calculation. Thus, we subjectively divided the snow pillows into four groups,
352 based on the general characteristics of the amount of snow pillow frozen precipitation with
353 respect to that of the other datasets. Our four groups are:

354 Group A (21% of pillows): snow pillow and WRF > at least 2 gridded datasets

355 Group B (25% of pillows): snow pillow > 8 or more datasets

356 Group C (28% of pillows): snow pillow near center of all datasets

357 Group D (25% of pillows): snow pillow < 8 or more datasets

358 A list of which snow pillows are in each group is provided in Table 3, and a map of their
359 distribution is included on Fig. 1b. See H16a for a complete list and map of watersheds. Group A
360 pillows largely run down the crest of the central Sierra, with a secondary cluster south of the
361 Merced watershed. Group B is mostly in the lee of the northern Sierra and at a few of the lowest-
362 elevation, windward side locations. Group C has a cluster of locations across the Tuolumne and
363 Cherry-Eleanor watersheds, with additional locations scattered across the entire region. Group D
364 is largely confined to the southern Sierra.

365 Example cumulative frozen precipitation traces from each of the four groups are shown in Fig. 6.
366 These traces are representative of their individual groups, although the 'best' and 'worst' datasets
367 at each pillow vary substantially. All four pillows shown, and 92% of all pillows (not shown),
368 receive more than 50% of their total WY2008 precipitation during three, 3-11 day periods: These
369 three events were identified in L15 as coincident 'missed storms' – i.e., underestimated snow
370 amounts – in Hamlet and Livneh. That such a large fraction of WY2008 snow fell during 3 short
371 periods is consistent with previous work showing that a substantial portion of annual California
372 precipitation tends to fall during a few synoptic-scale events, often containing atmospheric rivers
373 (L15; Dettinger et al. 2011): Dettinger et al. 2011 show that 50% of each year's precipitation
374 accumulates over less than 15 days in the Northern Sierra and less than 10 days in the Southern
375 Sierra, on average. Although each dataset has small errors during most of the precipitation

376 events, the errors are larger for the largest three events, and in many cases large
377 under/overestimation of precipitation in one event leads to the bulk of the WY disagreement. Fig.
378 6 also reveals that the timing of the precipitation is very similar in all datasets. This agreement is
379 not surprising for the gridded datasets, which rely on gauge data. However, the WRF
380 simulations, which are run in a regional climate framework (i.e., initialized in July 2007 and then
381 integrated continuously from then through Oct. 1, 2008) could potentially diverge in their
382 evolution of synoptic features after these features enter the edges of the outer domain. Since the
383 timing agrees so well across the simulations and with observations, the domain configuration and
384 nudging used on the outermost domain are providing strong constraints and keeping the
385 simulations in agreement with the reanalysis lateral boundary conditions. Finally, Fig. 6 reveals
386 information about the uncertainty estimate from the Newman dataset: The interquartile range
387 (IQR) of the Newman precipitation at 3 of the 4 pillows shown exceeds the total range of the 10
388 frozen precipitation estimates and observed snowfall, and this is true for 86% of all pillows.

389 The question remains as to whether each datasets' biases are consistent for each storm leading to
390 the total precipitation, or rather, if individual storm biases tend to vary in sign and cancel out
391 over the WY. We address this question using histograms of errors for each group (Fig. 7) and
392 boxplots of the errors for the three largest storms periods for each group (Fig. 8). For all datasets
393 in all groups, the histogram peaks lie between -10 and 0 mm, indicating a general tendency for
394 all datasets to underestimate small precipitation events. However, the majority of error in WY
395 total frozen precipitation is driven by larger events, which impact the skewness and width of the
396 histograms, and show up more clearly in boxplots from the three largest storms (Fig. 8).

397 Biases in group B in all 10 datasets seem to be fairly systematic, with the histograms of daily
398 errors and the errors for the three largest storms lying mostly below zero, indicating that all the
399 datasets have a tendency to underestimate frozen precipitation at these locations. In contrast,
400 biases in group D are quite skewed, with a heavy positive tail (Fig. 7), and the overprediction at
401 these pillows is mainly due to overprediction of frozen precipitation during the 3 big storm
402 events of the WY (Fig. 8); the errors in this group are a bit worse in the WRF simulations than
403 the gridded datasets, and are largest for the first storm period, especially in the NARR-forced
404 simulations. The histograms and boxplots for group A and, to a lesser degree group C, illustrate
405 the tendency for all the gridded datasets to underestimate precipitation in the central Sierra
406 Nevada, and the daily error statistics suggest this is largely a systematic problem for these two
407 groups of stations, with the majority of the histogram probability lying very close to or below
408 zero for all four datasets in these groups. The difference between group A and group C in the
409 gridded datasets is largely a result of Storm 1: in group A the boxes and most of the whiskers for
410 all three storms are consistently below zero, and the net result is large underestimation of WY
411 total frozen precipitation by the gridded datasets. In group C the gridded datasets' Storm 1 errors
412 are more consistently positive, thus compensating for the general underestimation and resulting
413 in smaller biases in WY total frozen precipitation. In both groups A and C, the biases in all 6
414 WRF simulations are more centered around zero than those of the gridded datasets, although in
415 group C there are more large positive outliers. The positive WY total biases in the two WSM6
416 simulations appear as a slight shift in the histograms (Fig. 7) and storm total barplots (Fig. 8).
417 The storm total barplots also reveal large differences in biases from storm to storm in the NARR-
418 forced simulations, with Storm 1 showing very large positive errors in Groups A, C, and D for all

419 microphysics parameterizations; the ERA-Interim-forced WRF biases are more consistent from
420 storm to storm.

421 These statistics for the daily and storm-total errors suggest that error outliers contribute strongly
422 to the overall biases for each ensemble member, and in some cases these biases can vary
423 significantly from one storm to another. This lack of a more systematic bias possibly suggests
424 that case studies of individual storm events of microphysics scheme performance may sometimes
425 lead to incorrect conclusions about their overall tendencies, since the biases are reflected in the
426 higher-order statistics rather than being systematic. California's tendency to receive most of each
427 water year's precipitation in a few concentrated storm periods means that these errors in
428 individual storms can strongly influence the water year biases. In addition, the different behavior
429 of the daily errors in the different groups, which cluster in localized regions (Fig. 1), mean that
430 bias tendencies can be highly variable spatially; thus conclusions about overall bias need to
431 either take this spatial variability into account or be drawn for rather large areas.

432 3.2 Comparison to Bayesian estimated precipitation

433 We now turn our attention to a comparison of the datasets' WY2008 precipitation with P_{inferred}
434 (Section 2.4). Although P_{inferred} is limited temporally to WY total precipitation and spatially to
435 basin-mean precipitation amounts, it provides an independent verification when combined with
436 the snow pillow dataset used in the previous section. We begin this comparison with maps of the
437 differences between P_{inferred} and basin-mean precipitation for each dataset (Fig. 9). Because
438 P_{inferred} includes an uncertainty estimate, we use this estimate in our comparison, and rather than
439 showing absolute or percentage-wise difference maps, we categorically compare with P_{inferred} .
440 The precipitation estimates for a large number of basins for all datasets fall within the IQR of
441 P_{inferred} . However, for the differences falling outside this uncertainty range, we see some patterns
442 quite similar to those we saw with the snow pillow comparisons. For instance, in the central
443 Sierra, the Yosemite-area watersheds of the Tuolumne River and/or Cherry and Eleanor Creeks
444 are generally drier than P_{inferred} in the four gridded datasets. The pillows in this region are mostly
445 Group C (i.e., snow pillow near center of all datasets); however, for most of the pillows in these
446 regions, at least 2 of the gridded datasets greatly underestimated frozen precipitation (not
447 shown), and Hamlet and Daymet both underestimate frozen precipitation at more pillows than
448 Livneh and Newman, consistent with P_{inferred} . In addition, the two WSM6 WRF runs are wetter
449 than P_{inferred} in several basins across the region, similar to the snow pillow comparisons. The
450 North Fork of the American River, in the northern Sierra, is consistently underestimated by the
451 gridded datasets and also tends to be underestimated by the WRF simulations, although in WRF
452 the underestimation is within the range of uncertainty for all but one ensemble member. Many of
453 the WRF simulations and gridded datasets overestimate precipitation in several of the smaller
454 watersheds throughout the region, in particular, the San Joaquin basins of Pitman, Bear, and
455 Bishop Creeks, and the Mokelumne basins; the pattern of which watersheds are overestimated is
456 more consistent across the WRF simulations than across the gridded datasets. Finally, in the
457 southern Sierra, the Kern River watershed precipitation is overestimated in all WRF simulations
458 but those using Thompson microphysics, as well as in Daymet, consistent with the differences
459 against snow pillows seen earlier in Fig. 5.

460 We see similar correspondence to the snow pillow comparisons when the datasets' precipitation
461 is scattered against P_{inferred} (Fig. 10), although as we saw with the mapped differences, many of
462 the datasets' estimates fall within the IQR of P_{inferred} . Displayed in this way it becomes clear that
463 all datasets tend to overestimate basin-mean precipitation in watersheds with the smallest
464 amounts of WY2008 precipitation inferred from streamflow. The WSM6 simulations
465 systematically overestimate precipitation in many watersheds, whereas the Morrison and
466 Thompson simulations more consistently fall within the IQR of P_{inferred} . The NARR-driven
467 simulations are consistently wetter than the ERA Interim-driven simulations, although again, this
468 effect is secondary compared to the impact of microphysics on total precipitation. The gridded
469 datasets consistently underestimate precipitation in the wettest watersheds; this underestimation,
470 when combined with the overestimation of drier watersheds, is reflected in a flatter slope of
471 linear fits against P_{inferred} (Fig. 10b) similar to that seen in the snow pillow comparisons. All the
472 WRF simulations do a consistently better job in distinguishing wet and dry basins; curiously,
473 despite their consistent wet bias and consistent with the snow pillow results, the slope of the
474 WSM6 simulations best matches that seen in P_{inferred} .

475 **4 Summary and Discussion**

476 In this manuscript, we explore the uncertainties during WY2008 in the Sierra Nevada of
477 California's high-elevation precipitation in 10 datasets: six WRF regional climate downscalings
478 with differing lateral boundary conditions and microphysical parameterizations, and four gauge-
479 based, interpolation-gridded precipitation datasets: Livneh, Hamlet, Daymet, and Newman. We
480 first compare frozen precipitation from these 10 datasets with positive daily changes in snow
481 water equivalent from a network of 95 snow pillows across the Sierra Nevada, then follow this
482 with a comparison of total precipitation with a precipitation dataset inferred from stream gauges
483 using a Bayesian inference method. Most of the manuscript focuses on comparisons of WY total
484 precipitation, but we also compare daily error statistics with the snow pillow data.

485 During WY2008, the gridded datasets, especially Livneh and Hamlet, tend to underestimate
486 frozen precipitation on the windward slope of the Sierra Nevada, particularly in the vicinity of
487 Yosemite National Park (Fig. 1). The WRF simulations consistently place more precipitation on
488 the windward slope than the gridded datasets, although the amount of precipitation along the
489 windward slope depends strongly on microphysical parameterization: the WRF simulations with
490 single-moment microphysics tend to overestimate precipitation along the windward slope,
491 whereas those with double-moment microphysics tend to better agree with the snow pillows at a
492 large proportion of the snow pillow locations. WRF simulations with NARR as lateral boundary
493 conditions are slightly wetter than those with ERA Interim boundary conditions, but this effect is
494 second order compared to microphysical differences.

495 All six WRF simulations somewhat overestimate the windward/leeside precipitation contrast in
496 the northern Sierra Nevada. This problem is more pronounced in the single-moment simulations,
497 which produce significantly more graupel than the double-moment schemes (not shown, e.g., see
498 Jankov et al. 2009), suggesting it could be partially due to too-efficient fallout of hydrometeors
499 on the windward side of the Sierra Nevada. However, since the issue appears in all six WRF
500 simulations irrespective of microphysics, another factor is probably also contributing, and we
501 speculate that this may be due to insufficient embedded convection in these simulations during
502 post-cold-frontal periods: Since embedded convection during post-cold-frontal storms tends to

503 result in an increase in leeside snow (Geerts et al. 2011), if these simulations have too little
504 embedded convection (due perhaps to their horizontal grid spacing) it could lead to not enough
505 precipitation being lofted into the lee. This hypothesis is also consistent with the location of this
506 bias in the northern Sierra Nevada, since the trajectory of most wintertime extratropical cyclones
507 hitting California means the northern Sierra Nevada spends more time in the cold sector than the
508 southern Sierra Nevada. The position of the most prominent ‘leeside’ underestimation east of
509 Lake Tahoe also raises the possibility that an additional issue is insufficient resolution of the
510 very narrow Carson range to the east of Lake Tahoe along with potential lake effects.

511 Daily errors were investigated by sorting the pillows into four different groups based on how
512 much snow pillow precipitation fell with respect to the other datasets; three of these four groups
513 clustered in highly localized regions. Daily and storm-total error statistics for the gridded
514 datasets were fairly consistently dry-biased, with outliers determining the overall site biases;
515 WRF’s biases were more centered around zero, and similarly, outliers contributed to the overall
516 tendencies for the different schemes. All 10 datasets underestimated small precipitation events.
517 WRF’s zero-centered daily biases and varying storm-total biases suggest case studies of
518 microphysical bias need to be interpreted carefully, since individual events may not sample the
519 distribution adequately to capture the error distribution. In addition, geographical tendencies of
520 biases varied widely with topographic characteristics. Thus conclusions about region-mean bias
521 need to be drawn for rather large areas, although region-mean bias is likely not adequate to
522 understand the dynamics leading to the biases.

523 Finally, the WRF simulations, in particular those with single-moment microphysics, better
524 distinguish wet-versus-dry pillows and watersheds than the gridded estimates. Even though the
525 WRF simulations with single-moment microphysics have a large wet bias, this wet bias was
526 fairly systematic across the region. The double-moment WRF simulations were less biased
527 overall, but the differences between wet and dry pillows/watersheds were not as large as
528 observed in these simulations. The gridded datasets have the least contrast between wet and dry
529 pillows and watersheds, with positive biases at dry pillows and watersheds and large negative
530 biases at the wettest pillows and watersheds; this result was also seen in H16. The differences
531 between the gridded datasets and WRF simulations in this respect are likely caused in part by
532 limitations of the statistical gridding methodologies used by the gridded datasets: These datasets
533 use linear regression techniques based on climatology to extrapolate gauge precipitation amounts
534 to regions without measurements, and thus embed these climatological relationships in their
535 estimates. When precipitation patterns differ from climatology, these gridded datasets would
536 tend toward climatology, and that introduces biases; in general, any statistical interpolation
537 technique will likely produce a smoother solution than the underlying data. All 10 datasets
538 somewhat underestimated the wet and dry contrast, and this consistent underestimation may be
539 related to small-scale terrain features unrepresented by the 6km grid (Minder et al. 2008). More
540 work should be done to understand whether this wet-versus-dry bias is systematic across
541 multiple water years. Further, until this issue is better understood, uncertainty in precipitation
542 should be explicitly considered when conducting research that uses precipitation as an input
543 (e.g., hydrological or ecological science). The reason for better wet-versus-dry ratios in the
544 single-moment than double-moment WRF simulations is unclear, and will require an in-depth
545 investigation of storm dynamics and cloud microphysical properties, which is beyond the scope
546 of the present manuscript.

547 The present study found systematic differences in the error statistics of WRF simulations based
548 on microphysical scheme sophistication, and also between these WRF dynamical downscalings
549 and gridded estimates of precipitation, with double-moment microphysics WRF simulations
550 outperforming both single-moment WRF simulations and the gridded datasets in most respects in
551 the Sierra Nevada. We highlight three significant results of this manuscript: 1) Many model
552 evaluations use PRISM-based gridded datasets as truth; had we taken this approach our
553 conclusion would be that all the WRF datasets have a wet bias, which the comparisons to snow
554 pillows and P_{inferred} show is not the case. 2) Our investigation of both water year total and single
555 storm precipitation biases revealed that the water year total biases were in some cases quite
556 dependent on biases from one major water year storm: Case studies of model configuration
557 performed for individual storm events could lead to incorrect conclusions about the model's
558 overall tendencies, since precipitation biases are reflected in the higher-order statistics rather
559 than being systematic. 3) Our focus both Sierra-wide and at smaller scales (e.g., watershed scale)
560 reveals that very different biases can exist at highly localized scales. These three results provide
561 guidance for future research, suggesting care be taken regarding spatial and temporal scales and
562 with the "observations" used for model evaluation: Without this care, studies may reach incorrect
563 conclusions about model performance and where to focus future model development. Our results
564 for the Sierra Nevada should be largely transferrable to other mid-latitude mountainous regions
565 that receive most of their precipitation from orographically enhanced synoptic scale events, with
566 that caveat that the performance of the gauge-based datasets is likely sensitive to gauge density,
567 and that California's tendency to receive most precipitation in a few large events per year
568 increases the risk that individual case studies may not represent overall biases.

569 Our results are limited by the single water year of available WRF output. Further, we emphasize
570 that the water year presented was particularly problematic for gridded datasets (L15); further
571 work is needed to examine whether the patterns in biases that we show here are consistent during
572 water years with different conditions, particularly years with more extreme precipitation, and to
573 investigate the dynamic and thermodynamic causes for interannual changes in orographic
574 precipitation gradient. Finally, although using WRF to downscale reanalysis data is shown in this
575 manuscript to improve over gridded datasets during certain water years, it is a computationally
576 expensive option, and may not always be feasible for all applications. Two possible and
577 promising alternative approaches are hybrid techniques that combine statistical and dynamical
578 downscaling approaches (e.g., Sun et al. 2015) or simpler and less computationally demanding
579 dynamical models, such as the Intermediate Complexity Atmospheric Research Model (Gutmann
580 et al. 2016), although more work is needed in the development and testing of such tools. These
581 approaches could also potentially be used to improve gauge-based gridded datasets,
582 incorporating dynamical information to improve upon the weaknesses of purely statistical and
583 elevation-based gridding techniques.

584 **Acknowledgments and Data**

585 This work was supported by the National Science Foundation (NSF) grant EAR – 1344595. This
586 work utilized the Janus supercomputer, which is supported by the National Science Foundation
587 (award number CNS-0821794) and the University of Colorado Boulder. The Janus
588 supercomputer is a joint effort of the University of Colorado Boulder, the University of Colorado
589 Denver and the National Center for Atmospheric Research. The authors acknowledge the NOAA

590 Research and Development High Performance Computing Program for providing computing and
591 storage resources that have contributed to the research results reported within this paper. URL:
592 <http://rdhpcs.noaa.gov>. We would like to thank three anonymous reviewers as well as Rob Cifelli
593 and Andrea Thorstensen for providing comments which improved the manuscript.

594 All data used in this study are publicly available. ERA Interim data were downloaded from
595 ECMWF using tools available on their website ([http://apps.ecmwf.int/datasets/data/interim-full-
596 daily/levtype=sfc/](http://apps.ecmwf.int/datasets/data/interim-full-daily/levtype=sfc/)). The North American Regional Reanalysis data were retrieved from the
597 Research Data Archive at the National Center for Atmospheric Research, Computational and
598 Information Systems Laboratory. <http://rda.ucar.edu/datasets/ds608.0/>, accessed 15 Aug 2014.
599 Daymet (Thornton et al. 2014) was retrieved from the Distributed Active Archive Center of Oak
600 Ridge National Laboratory. Hamlet data are housed by the University of Washington Climate
601 Impacts Group (<http://ces.washington.edu/cig/data/wus.shtml>). Livneh data are available from
602 an ftp site (www.hydro.washington.edu/Lettenmaier/Data/livneh/livneh.et.al.2013.page.html).
603 Newman data are available on NCAR's Earth System Grid
604 (https://www.earthsystemgrid.org/dataset/gridded_precip_and_temp.html). The California
605 Department of Water Resources snow pillow data are available from the California Data
606 Exchange Center (CDEC; <http://cdec.water.ca.gov>). The WRF simulations are available through
607 personal communication with the corresponding author (mimi.hughes@noaa.gov).

608 **References**

- 609 Accessed 10 May 2014.: California Data Exchange Center—Snow.
- 610 Barlage, M., M. Tewari, F. Chen, G. Miguez-Macho, Z. L. Yang, and G. Y. Niu, 2015: The
611 effect of groundwater interaction in North American regional climate simulations with
612 WRF/Noah-MP. *Climatic Change*, 129, 485-498.
- 613 Beaumont, R., 1965: Mt. Hood pressure pillow snow gage. *Journal of Applied Meteorology*, 4,
614 626-631.
- 615 Caldwell, P., H. N. S. Chin, D. C. Bader, and G. Bala, 2009: Evaluation of a WRF dynamical
616 downscaling simulation over California. *Climatic Change*, 95, 499-521.
- 617 Clark, M. P., and A. G. Slater, 2006: Probabilistic quantitative precipitation estimation in
618 complex terrain. *Journal of Hydrometeorology*, 7, 3-22.
- 619 Daly, C., R. P. Neilson, and D. L. Phillips, 1994: A STATISTICAL TOPOGRAPHIC MODEL
620 FOR MAPPING CLIMATOLOGICAL PRECIPITATION OVER MOUNTAINOUS
621 TERRAIN. *Journal of Applied Meteorology*, 33, 140-158.
- 622 Daly, C., and Coauthors, 2008: Physiographically sensitive mapping of climatological
623 temperature and precipitation across the conterminous United States. *International Journal of
624 Climatology*, 28, 2031-2064.
- 625 Dee, D. P., and Coauthors, 2011: The ERA-Interim reanalysis: configuration and performance of
626 the data assimilation system. *Quarterly Journal of the Royal Meteorological Society*, 137, 553-
627 597.
- 628 Dettinger, M. D., F. M. Ralph, T. Das, P. J. Neiman, and D. R. Cayan, 2011: Atmospheric
629 Rivers, Floods and the Water Resources of California. *Water*, 3, 445-478.
- 630 Farnes, P., 1967: Criteria for determining mountain snow pillow sites. *Proc. 35th Western Snow
631 Conf*, 59-62.
- 632 Feser, F., B. Rockel, H. von Storch, J. Winterfeldt, and M. Zahn, 2011: Regional Climate Models
633 Add Value To Global Model Data A Review and Selected Examples. *Bulletin of the American
634 Meteorological Society*, 92, 1181-1192.
- 635 Geerts, B., Q. Miao, and Y. Yang, 2011: Boundary Layer Turbulence and Orographic
636 Precipitation Growth in Cold Clouds: Evidence from Profiling Airborne Radar Data. *Journal of
637 the Atmospheric Sciences*, 68, 2344-2365.
- 638 Gutmann, E. D., and Coauthors, 2012: A Comparison of Statistical and Dynamical Downscaling
639 of Winter Precipitation over Complex Terrain. *Journal of Climate*, 25, 262-281.

- 640 Gutmann, E. D., I. Barstad, M. P. Clark, J. R. Arnold, and R. M. Rasmussen (2016), The
641 Intermediate Complexity Atmospheric Research Model, *J. Hydrometeorol*, 17, 957–973, doi:
642 10.1175/JHM-D-15-0155.1.
- 643 Hamlet, A., and Coauthors, 2010: Final project report for the Columbia basin climate change
644 scenarios project. Report, University of Washington, Seattle, WA. Available at: [http://www.
645 hydro.washington.edu/2860/report/](http://www.hydro.washington.edu/2860/report/). [Accessed November 20, 2015].
- 646 Hamlet, A. F., and D. P. Lettenmaier, 2005: Production of temporally consistent gridded
647 precipitation and temperature fields for the continental United States. *Journal of
648 Hydrometeorology*, 6, 330-336.
- 649 Henn, B., M. P. Clark, D. Kavetski, and J. D. Lundquist, 2015: Estimating mountain basin-mean
650 precipitation from streamflow using Bayesian inference. *Water Resources Research*, 51, 8012-
651 8033. doi:10.1002/2014WR016736
- 652 Henn, B., M. P. Clark, D. Kavetski, A. J. Newman, M. Hughes, B. McGurk, and J. D. Lundquist,
653 2016a: Spatiotemporal patterns of precipitation inferred from streamflow observations across the
654 Sierra Nevada mountain range. *Journal of Hydrology* (in press).
655 doi:dx.doi.org/10.1016/j.jhydrol.2016.08.009.
- 656 Henn, B., Newman, A. J., Livneh, B., Daly, C., & Lundquist, J. D., 2016b. An assessment of
657 differences in gridded precipitation datasets in complex terrain. *Journal of Hydrology* (in press).
658 doi:<http://dx.doi.org/10.1016/j.jhydrol.2017.03.008>.
- 659 Hong, S.-Y., and J.-O. J. Lim, 2006: The WRF single-moment 6-class microphysics scheme
660 (WSM6). *J. Korean Meteor. Soc*, 42, 129-151.
- 661 Hong, S.-Y., Y. Noh, and J. Dudhia, 2006: A new vertical diffusion package with an explicit
662 treatment of entrainment processes. *Monthly Weather Review*, 134, 2318-2341.
- 663 Iacono, M. J., J. S. Delamere, E. J. Mlawer, M. W. Shephard, S. A. Clough, and W. D. Collins,
664 2008: Radiative forcing by long-lived greenhouse gases: Calculations with the AER radiative
665 transfer models. *Journal of Geophysical Research-Atmospheres*, 113.
- 666 Jankov, I., P. J. Schultz, C. J. Anderson, and S. E. Koch, 2007: The impact of different physical
667 parameterizations and their interactions on cold season QPF in the American River basin. *Journal
668 of Hydrometeorology*, 8, 1141-1151.
- 669 Jankov, I., J. W. Bao, P. J. Neiman, P. J. Schultz, H. L. Yuan, and A. B. White, 2009: Evaluation
670 and Comparison of Microphysical Algorithms in ARW-WRF Model Simulations of Atmospheric
671 River Events Affecting the California Coast. *Journal of Hydrometeorology*, 10, 847-870.
- 672 Jankov, I., and Coauthors, 2011: An Evaluation of Five ARW-WRF Microphysics Schemes
673 Using Synthetic GOES Imagery for an Atmospheric River Event Affecting the California Coast.
674 *Journal of Hydrometeorology*, 12, 618-633.

- 675 Jeton, A. E., S. A. Watkins, and J. Huntington, 2006: Evaluation of Precipitation Estimates from
676 PRISM for the 1961-90 and 1971-2000 Data Sets, Nevada. US Geological Survey.
- 677 Jimenez, P. A., J. Dudhia, J. F. Gonzalez-Rouco, J. Navarro, J. P. Montavez, and E. Garcia-
678 Bustamante, 2012: A Revised Scheme for the WRF Surface Layer Formulation. *Monthly*
679 *Weather Review*, 140, 898-918.
- 680 Johnson, J. B., and D. Marks, 2004: The detection and correction of snow water equivalent
681 pressure sensor errors. *Hydrological Processes*, 18, 3513-3525.
- 682 Kain, J. S., 2004: The Kain-Fritsch convective parameterization: An update. *Journal of Applied*
683 *Meteorology*, 43, 170-181.
- 684 Liu, C. H., K. Ikeda, G. Thompson, R. Rasmussen, and J. Dudhia, 2011: High-Resolution
685 Simulations of Wintertime Precipitation in the Colorado Headwaters Region: Sensitivity to
686 Physics Parameterizations. *Monthly Weather Review*, 139, 3533-3553.
- 687 Livneh, B., and Coauthors, 2013: A Long-Term Hydrologically Based Dataset of Land Surface
688 Fluxes and States for the Conterminous United States: Update and Extensions. *Journal of*
689 *Climate*, 26, 9384-9392.
- 690 Lundquist, J. D., P. J. Neiman, B. Martner, A. B. White, D. J. Gottas, and F. M. Ralph, 2008:
691 Rain versus snow in the Sierra Nevada, California: Comparing Doppler profiling radar and
692 surface observations of melting level. *Journal of Hydrometeorology*, 9, 194-211.
- 693 Lundquist, J. D., M. Hughes, B. Henn, E. D. Gutmann, B. Livneh, J. Dozier, and P. Neiman,
694 2015: High-Elevation Precipitation Patterns: Using Snow Measurements to Assess Daily
695 Gridded Datasets across the Sierra Nevada, California. *Journal of Hydrometeorology*, 16, 1773-
696 1792.
- 697 Mesinger, F., and Coauthors, 2006: North American regional reanalysis. *Bulletin of the*
698 *American Meteorological Society*, 87, 343-+.
- 699 Minder, J. R., and D. E. Kingsmill, 2013: Mesoscale Variations of the Atmospheric Snow Line
700 over the Northern Sierra Nevada: Multiyear Statistics, Case Study, and Mechanisms. *Journal of*
701 *the Atmospheric Sciences*, 70, 916-938.
- 702 Minder, J. R., D. R. Durran, G. H. Roe, and A. M. Anders, 2008: The climatology of small-scale
703 orographic precipitation over the Olympic Mountains: Patterns and processes. *Quarterly Journal*
704 *of the Royal Meteorological Society*, 134, 817-839.
- 705 Morrison, H., G. Thompson, and V. Tatarskii, 2009: Impact of Cloud Microphysics on the
706 Development of Trailing Stratiform Precipitation in a Simulated Squall Line: Comparison of
707 One- and Two-Moment Schemes. *Monthly Weather Review*, 137, 991-1007.

- 708 Nelson, B. R., O. P. Prat, D. J. Seo, and E. Habib, 2016: Assessment and Implications of NCEP
709 Stage IV Quantitative Precipitation Estimates for Product Intercomparisons. *Weather and*
710 *Forecasting*, 31, 371-394.
- 711 Newman, A. J., and Coauthors, 2015: Gridded Ensemble Precipitation and Temperature
712 Estimates for the Contiguous United States. *Journal of Hydrometeorology*, 16, 2481-2500.
- 713 Niu, G. Y., and Coauthors, 2011: The community Noah land surface model with
714 multiparameterization options (Noah-MP): 1. Model description and evaluation with local-scale
715 measurements. *Journal of Geophysical Research-Atmospheres*, 116.
- 716 Pavelsky, T. M., S. Kapnick, and A. Hall, 2011: Accumulation and melt dynamics of snowpack
717 from a multiresolution regional climate model in the central Sierra Nevada, California. *Journal of*
718 *Geophysical Research-Atmospheres*, 116.
- 719 Skamarock, W., J.B. Klemp, J. Dudhia (more), 2008: A Description of the Advanced Research
720 WRF Version 3. . NCAR Technical Note NCAR/TN-475+STR.
- 721 Sun F, D Walton, and A Hall, 2015: A hybrid dynamical-statistical downscaling technique, part
722 II: End-of-century warming projections predict a new climate state in the Los Angeles region.
723 *Journal of Climate*, 28(12): 4618-4636. DOI: 10.1175/JCLI-D-14-00197.1.
- 724 Tewari, M., and Coauthors, 2004: Implementation and verification of the unified NOAH land
725 surface model in the WRF model. 20th conference on weather analysis and forecasting/16th
726 conference on numerical weather prediction, pp. 11–15.
- 727 Thompson, G., P. R. Field, R. M. Rasmussen, and W. D. Hall, 2008: Explicit Forecasts of Winter
728 Precipitation Using an Improved Bulk Microphysics Scheme. Part II: Implementation of a New
729 Snow Parameterization. *Monthly Weather Review*, 136, 5095-5115.
- 730 Thornton, P. E., S. W. Running, and M. A. White, 1997: Generating surfaces of daily
731 meteorological variables over large regions of complex terrain. *Journal of Hydrology*, 190, 214-
732 251.
- 733 Thornton, P.E., M.M. Thornton, B.W. Mayer, N. Wilhelm, Y. Wei, R. Devarakonda, and R.B.
734 Cook. 2014. Daymet: Daily Surface Weather Data on a 1-km Grid for North America, Version 2.
735 ORNL DAAC, Oak Ridge, Tennessee, USA. Accessed February 01, 2016. Time period: 2007-
736 10-01 to 2008-09-30. <http://dx.doi.org/10.3334/ORNLDAAC/1219>
- 737 Willie, D., and Coauthors, 2016: Evaluation of Multisensor Quantitative Precipitation Estimation
738 in Russian River Basin. *Journal of Hydrologic Engineering*, E5016002.
- 739 Xue, Y. K., Z. Janjic, J. Dudhia, R. Vasic, and F. De Sales, 2014: A review on regional
740 dynamical downscaling in intraseasonal to seasonal simulation/prediction and major factors that
741 affect downscaling ability. *Atmospheric Research*, 147, 68-85.

742 Yang, D. Q., D. Kane, Z. P. Zhang, D. Legates, and B. Goodison, 2005: Bias corrections of long-
743 term (1973-2004) daily precipitation data over the northern regions. *Geophysical Research*
744 *Letters*, 32.

745 Yang, H. W., and B. Wang, 2012: Reduction of systematic biases in regional climate
746 downscaling through ensemble forcing. *Climate Dynamics*, 38, 655-665. Barlage, Michael, et al.
747 "The effect of groundwater interaction in North American regional climate simulations with
748 WRF/Noah-MP." *Climatic Change* 129.3-4 (2015): 485-498.
749

750

751 **Figure 1.** (a) WRF terrain (m) and extent of 18 km (edge of color fill) and 6 km (red outline)
 752 domains. Green line shows focus area of manuscript. (b) 2-min USGS terrain (black contours at
 753 0, 1000, and 3000m), watershed extent for stream gauges (purple outlines and green color fill),
 754 and locations of snow pillows (colored markers) for green-outlined region of (a). Color and
 755 shape of snow pillow markers indicate their group in section 3: Group A, blue circles; Group B,
 756 green triangles; Group C, red diamonds; Group D, orange squares. Letters A, B, C, and D
 757 indicate locations of example pillows. Numbers 1-5 indicate watersheds highlighted in the text,
 758 with names given in upper right.

759 **Figure 2.** WY2008 total (mm): a) snow pillow observed snowfall, b) multi-product mean frozen
 760 precipitation, c) difference of multi-product mean frozen precipitation and snow pillow observed
 761 snowfall, d) gridded dataset mean frozen precipitation, e) WRF mean frozen precipitation, f)
 762 gridded dataset difference from multi-product mean, and g) WRF difference from multi-product
 763 mean. WRF amounts (e, g) are the sum of snow and graupel; gridded datasets (d, f) are the sum
 764 of precipitation on all days when $T_{min} < 0^{\circ}\text{C}$ 2-minute terrain is plotted every 1000m starting at
 765 0m.

766 **Figure 3:** Precipitation difference of individual dataset WY2008 total frozen precipitation from
 767 multi-product mean (Fig. 2a; mm). WRF amounts (a, b, c, f, g, h) are the sum of snow and
 768 graupel; gridded datasets (d, e, I, j) are the sum of precipitation on all days when $T_{min} < 0^{\circ}\text{C}$. 2-
 769 minute terrain is plotted at 0, 1000, and 3000 m.

770 **Figure 4.** (a) Scatterplot of frozen precipitation versus snow pillow data at nearest gridpoint
 771 (mm). b) Linear regressions for scatterplots of (a).

772 **Figure 5.** a) Meridional gradient of smoothed terrain (color fill, m 6km-1), terrain from WRF
 773 simulation (black contours, every 1000m), and locations of snow pillows (black dots). b) Snow
 774 pillow water year total snow (mm) versus smoothed zonal terrain gradient (x-axis, as in (a)) and
 775 latitude (y-axis). (c-l) As in (b), but percent difference between frozen precipitation and snow
 776 pillow snow (%).

777 **Figure 6.** Cumulative traces of daily snow pillow snow and frozen precipitation for examples
 778 from each of the 4 groups of snow pillows (A-D) outlined in the text. Purple and cyan arrows
 779 show start dates of ‘missed storms’ from Lundquist et al. (2015) in Livneh and Hamlet datasets,
 780 respectively. Locations of example pillows are shown with letters in Fig. 1b. Group A (21% of
 781 pillows): snow pillow and WRF > at least 2 gridded datasets; group B (25% of pillows): snow
 782 pillow > 8 or more datasets; group C (28% of pillows): snow pillow near center of all datasets;
 783 group D (25% of pillows): snow pillow < 8 or more datasets.

784 **Figure 7.** Histogram of errors (gridded-snow pillow) of smoothed daily ‘frozen’ precipitation for
 785 each of the 4 groups of snow pillows (A-D) outlined in the text, on days with smoothed observed
 786 snow > 5 mm. Group A (21% of pillows): snow pillow and WRF > at least 2 gridded datasets;
 787 group B (25% of pillows): snow pillow > 8 or more datasets; group C (28% of pillows): snow
 788 pillow near center of all datasets; group D (25% of pillows): snow pillow < 8 or more datasets.

789 **Figure 8.** Boxplots of errors (gridded-snow pillow) of ‘storm’ total ‘frozen’ precipitation (mm)
790 for each of the 4 groups of snow pillows (A-D) outlined in the text, for three major storm periods
791 highlighted with arrows in Fig. 5: Storm1: Jan 3-8, 2008; Storm 2: Jan. 26 – Feb. 5, 2008; and
792 Storm 3: Feb. 19-26, 2008. Outliers shown with red + signs are +/- 2 standard deviations. Group
793 A (21% of pillows): snow pillow and WRF > at least 2 gridded datasets; group B (25% of
794 pillows): snow pillow > 8 or more datasets; group C (28% of pillows): snow pillow near center
795 of all datasets; group D (25% of pillows): snow pillow < 8 or more datasets.

796 **Figure 9.** (a) Median basin-mean WY2008 Pinferred (mm). (b-k) Categorical difference of
797 gridded dataset basin-mean precipitation (P) and Pinferred. Categories are: 1: $P < \min \text{Pinferred}$,
798 2: $\min \text{Pinferred} < P < 25\% \text{ Pinferred}$, 3: $25\% \text{ Pinferred} < P < 50\% \text{ Pinferred}$, 4: $50\% \text{ Pinferred} < P < 75\% \text{ Pinferred}$, 5: $75\% \text{ Pinferred} < P < \max \text{Pinferred}$, 6: $\max \text{Pinferred} < P$. Note that categories 3 and 4 are within the interquartile range of uncertainty.

801 **Figure 10.** (a) Basin-mean precipitation (mm) from gridded datasets (see legend), as a function
802 of Pinferred (black crosses). Large black crosses show median and gray shading bounded by
803 smaller black crosses show interquartile range (IQR) of Pinferred. (b) Linear fit for each dataset.
804 Black solid line shows Bayesian median and gray shading bounded by black dashed lines show
805 IQR.

806

807

808

809

810

811

812

813

814 **Table 1:** Details of the WRF simulations. References for microphysics schemes and lateral
 815 boundary conditions can be found in Section 2.1.

Simulation name	Summary description	Resolution of lateral boundary conditions (LBCs)	Compute hours per 30 days of simulation (run on 120 CPUs)
E.Morr	ERA Interim LBCs and Morrison microphysics	~80 km	~4000
E.Thom	ERA Interim LBCs and Thompson microphysics	~80 km	~4000
E.WSM6	ERA Interim LBCs and WSM6 microphysics	~80 km	~3700
N.Morr	NARR LBCs and Morrison microphysics	~38 km	~4000
N.Thom	NARR LBCs and Thompson microphysics	~38 km	~4000
N.WSM6	NARR LBCs and WSM6 microphysics	~38 km	~3700

816

817 **Table 2:** Mean, median, and standard deviation of differences between total WY2008 frozen
 818 precipitation for each dataset and the snow pillow observations.

819

Dataset	Mean Difference (mm)	Median Difference (mm)	Standard Deviation of Difference (mm)
E.Morr	-2.4	36.5	222.0
E.Thom	-62.8	-16.6	210.4
E.WSM6	103.0	132.3	255.5
N.Morr	29.5	68.6	236.1
N.Thom	-22.9	19.0	221.0
N.WSM6	146.1	174.2	274.3
Livneh	-142.8	-111.1	259.2
Hamlet	-164.3	-141.1	222.1
Daymet	-60.7	-47.7	227.7
Newman	-81.1	-59.0	215.6

820

821

822

823 **Table 3.** Snow pillows in each of the 4 groups of Figs. 5 and 6. Group A: snow pillow and WRF
 824 > at least 2 gridded datasets; Group B: snow pillow > 8 or more datasets; Group C: snow pillow
 825 near center of all datasets; Group D: snow pillow < 8 or more datasets

	Name	Elev.	Latitude	Longitude
G r o u p A P i l l o w s	Agnew Pass	2880	37.728	-119.143
	Bloods Creek	2195	38.45	-120.033
	Caples Lake (DWR)	2438	38.71	-120.042
	Gianelli Meadow	2560	38.205	-119.892
	Green Mountain	2408	37.555	-119.238
	Hagans Meadow	2438	38.853	-119.94
	Chilkoot Meadow	2179	37.41	-119.49
	Dana Meadows	2987	37.897	-119.257
	Ebbetts Pass	2652	38.561	-119.808
	Highland Meadow	2652	38.49	-119.805
	Mud Lake	2408	38.615	-120.14
	Poison Flat	2408	38.501	-119.631
	Poison Ridge	2103	37.403	-119.52
	Schneiders	2667	38.747	-120.068
	Stanislaus Meadow	2362	38.5	-119.937
	Squaw Valley Gold Coast	2499	39.194	-120.276
	Van Vleck	2042	38.945	-120.305
	Ward Creek 3	2057	39.137	-120.22
	Echo Peak 5	2377	38.849	-120.079
Graveyard Meadow	2103	37.465	-119.29	
G r o u p B P i l l o w s	Alpha (SMUD)	2316	38.805	-120.215
	Bucks Lake	1753	39.85	-121.242
	Blue Canyon	1609	39.276	-120.708
	Casa Vieja Meadows	2530	36.2	-118.268
	Cottonwood Lakes	3094	36.483	-118.177
	Dismal Swamp	2149	41.993	-120.165
	Four Trees	1570	39.813	-121.321
	Forni Ridge	2316	38.805	-120.213
	Gem Pass	3277	37.78	-119.17
	Heavenly Valley	2682	38.929	-119.917
	Independence Lake (SCS)	2576	39.435	-120.322
	Kettle Rock	2225	40.14	-120.715
	Lobdell Lake	2804	38.44	-119.377
	Monitor Pass	2545	38.67	-119.615
	Marlette Lake	2438	39.173	-119.905
	Mount Rose Ski Area	2713	39.326	-119.902
	Pascoes	2789	35.967	-118.35
	Quaking Aspen	2195	36.117	-118.54
	Rattlesnake	1859	40.125	-121.043
	Slate Creek	1737	41.045	-122.478
Snow Mountain	1814	40.778	-121.782	
Upper Burnt Corral	2957	37.183	-118.937	
Big Meadows (SCS)	2652	39.458	-119.946	
Grizzly Ridge	2103	39.917	-120.645	
	Name	Elev.	Latitude	Longitude

G r o u p C P i l l o w s	Adin Mountain	1890	41.237	-120.792
	Blackcap Basin	3139	37.067	-118.77
	Chagoopa Plateau	3139	36.497	-118.442
	Central Sierra Snow Laboratory	2103	39.325	-120.367
	Deadman Creek	2819	38.332	-119.653
	Gin Flat	2149	37.767	-119.773
	Huntington Lake (USBR)	2134	37.228	-119.221
	Horse Meadow	2560	38.158	-119.662
	Independence Creek	1981	39.494	-120.293
	Kaiser Point	2804	37.3	-119.1
	Mammoth Pass (USBR)	2835	37.61	-119.033
	Pilot Peak (Dwr)	2073	39.786	-120.875
	Robbs Saddle	1798	38.912	-120.378
	Lower Relief Valley	2469	38.243	-119.758
	Slide Canyon	2804	38.092	-119.43
	Sonora Pass Bridge	2667	38.318	-119.601
	Tunnel Guard Station	2713	36.367	-118.288
	Volcanic Knob	3063	37.388	-118.903
	Cedar Pass	2164	41.583	-120.303
	Blue Lakes	2438	38.613	-119.931
	Black Springs	1981	38.375	-120.192
	Humbug	1981	40.115	-121.368
	Huysink	2012	39.282	-120.527
	Lower Kibbie Ridge	2042	38.032	-119.877
Medicine Lake	2042	41.592	-121.61	
Paradise Meadow	2332	38.047	-119.67	
Ostrander Lake	2499	37.637	-119.55	
G r o u p D P i l l o w s	Beach Meadows	2332	36.127	-118.293
	Crabtree Meadow	3261	36.563	-118.345
	Giant Forest (USACE)	2027	36.562	-118.765
	Independence Camp	2134	39.454	-120.299
	Leavitt Meadows	2195	38.305	-119.552
	Robbs Powerhouse	1570	38.903	-120.375
	Rock Creek Lakes	3048	37.455	-118.743
	Rubicon Peak 2	2286	39.001	-120.14
	South Lake	2926	37.176	-118.562
	Sawmill	3109	37.162	-118.562
	Tahoe City Cross	2057	39.172	-120.154
	Truckee 2	1951	39.3	-120.194
	Virginia Lakes Ridge	2835	38.077	-119.234
	West Woodchuck Meadow	2774	37.03	-118.918
	Big Pine Creek	2987	37.128	-118.475
	Bishop Pass	3414	37.1	-118.557
	Mitchell Meadow	3018	36.737	-118.712
	Silver Lake	2164	38.678	-120.118
	Spratt Creek	1875	38.667	-119.818
	Tamarack Summit	2301	37.165	-119.2
Tuolumne Meadows	2621	37.873	-119.35	
Upper Tyndall Creek	3475	36.65	-118.397	
Wet Meadows	2728	36.348	-118.572	
Gold Lake	2057	39.675	-120.615	

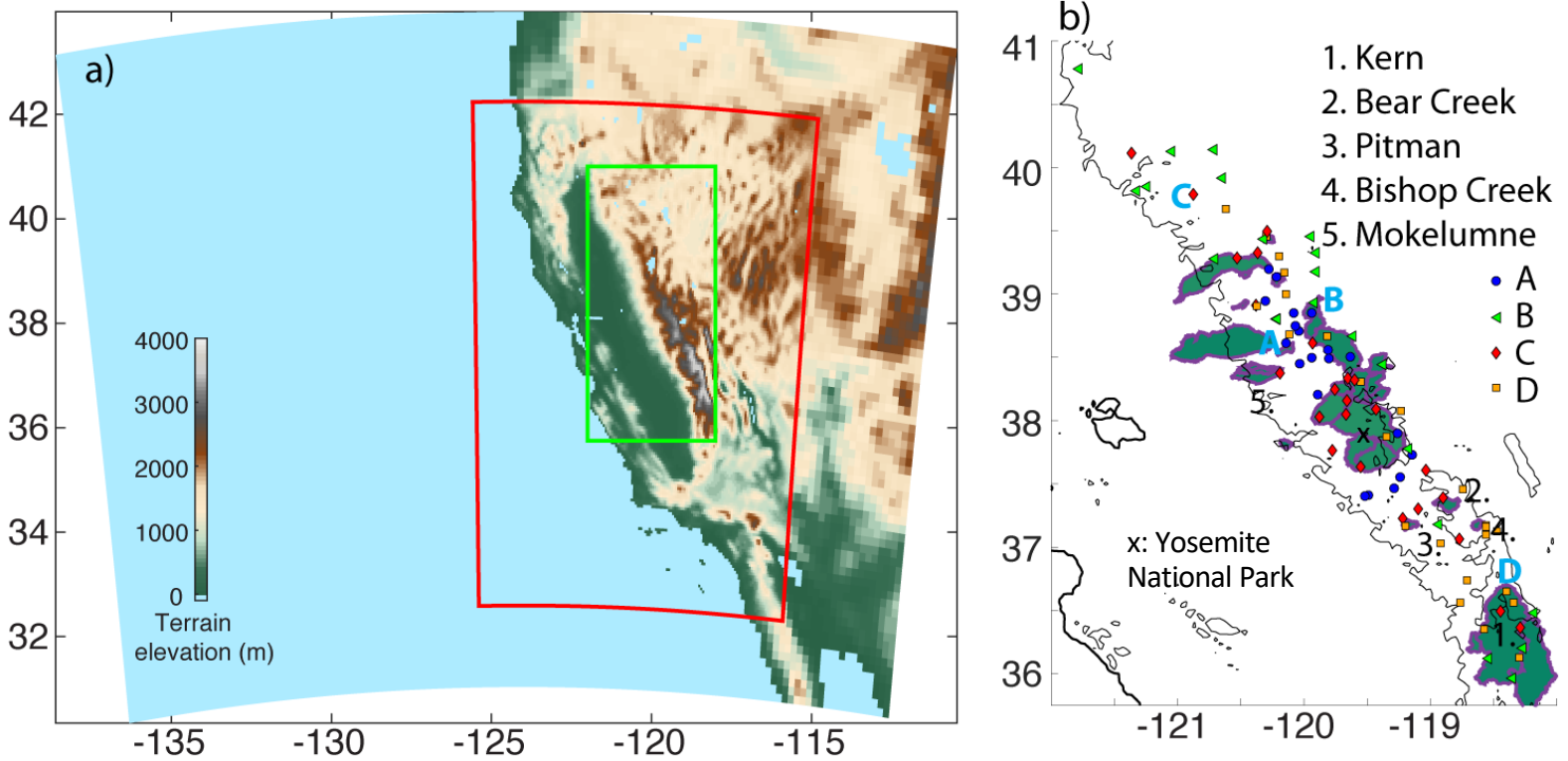


Fig. 1: (a) WRF terrain (m) and extent of 18 km (edge of color fill) and 6 km (red outline) domains. Green line shows focus area of manuscript. (b) 2-min USGS terrain (black contours at 0, 1000, and 3000m), watershed extent for stream gauges (purple outlines and green color fill), and locations of snow pillows (colored markers) for green-outlined region of (a). Color and shape of snow pillow markers indicate their group in section 3: Group A, blue circles; Group B, green triangles; Group C, red diamonds; Group D, orange squares. Letters A, B, C, and D indicate locations of example pillows. Numbers 1-5 indicate watersheds highlighted in the text, with names given in upper right.

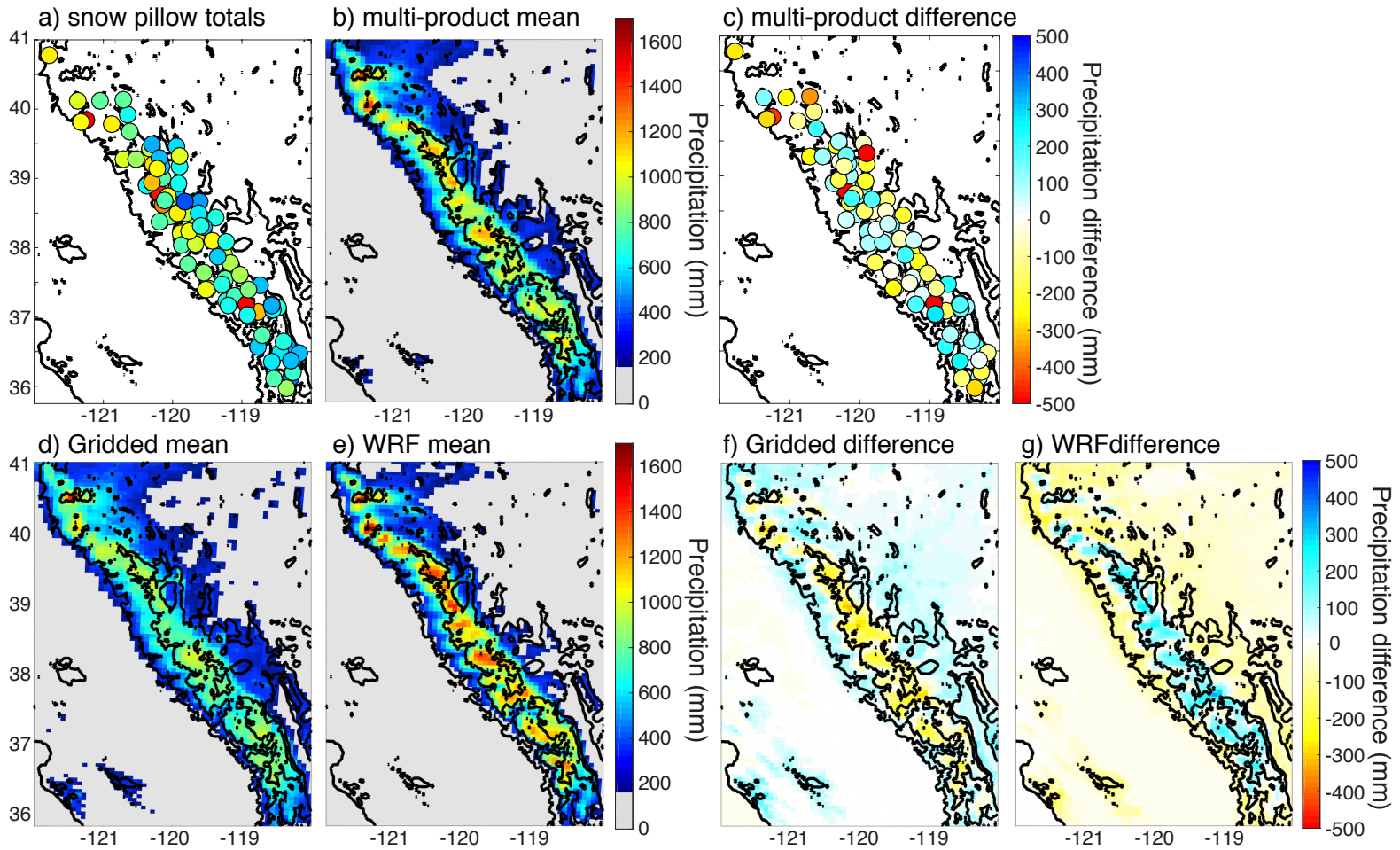


Fig. 2: WY208 total (mm): a) snow pillow observed snowfall, b) multi-product mean frozen precipitation, c) difference of multi-product mean frozen precipitation and snow pillow observed snowfall, d) gridded dataset mean frozen precipitation, e) WRF mean frozen precipitation, f) gridded dataset difference from multi-product mean, and g) WRF difference from multi-product mean. WRF amounts (e, g) are the sum of snow and graupel; gridded datasets (d, f) are the sum of precipitation on all days when $T_{min} < 0\text{ }^{\circ}\text{C}$ 2-minute terrain is plotted every 1000m starting at 0m.

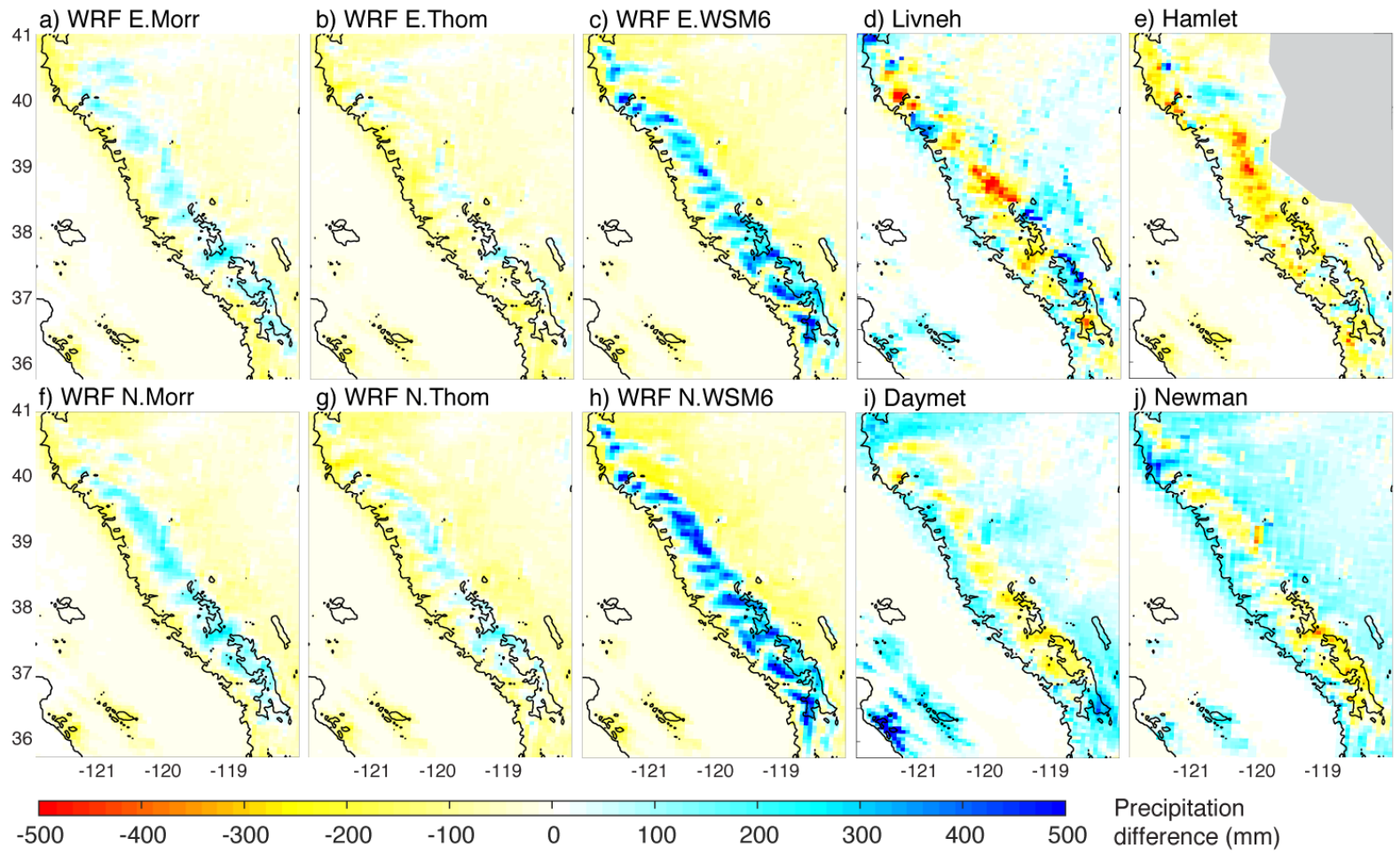


Fig. 3: Precipitation difference of individual dataset WY2008 total frozen precipitation from multi-product mean (Fig. 2a; mm). WRF amounts (a, b, c, f, g, h) are the sum of snow and graupel; gridded datasets (d, e, i, j) are the sum of precipitation on all days when $T_{min} < 0$ °C. 2-minute terrain is plotted at 0, 1000, and 3000 m.

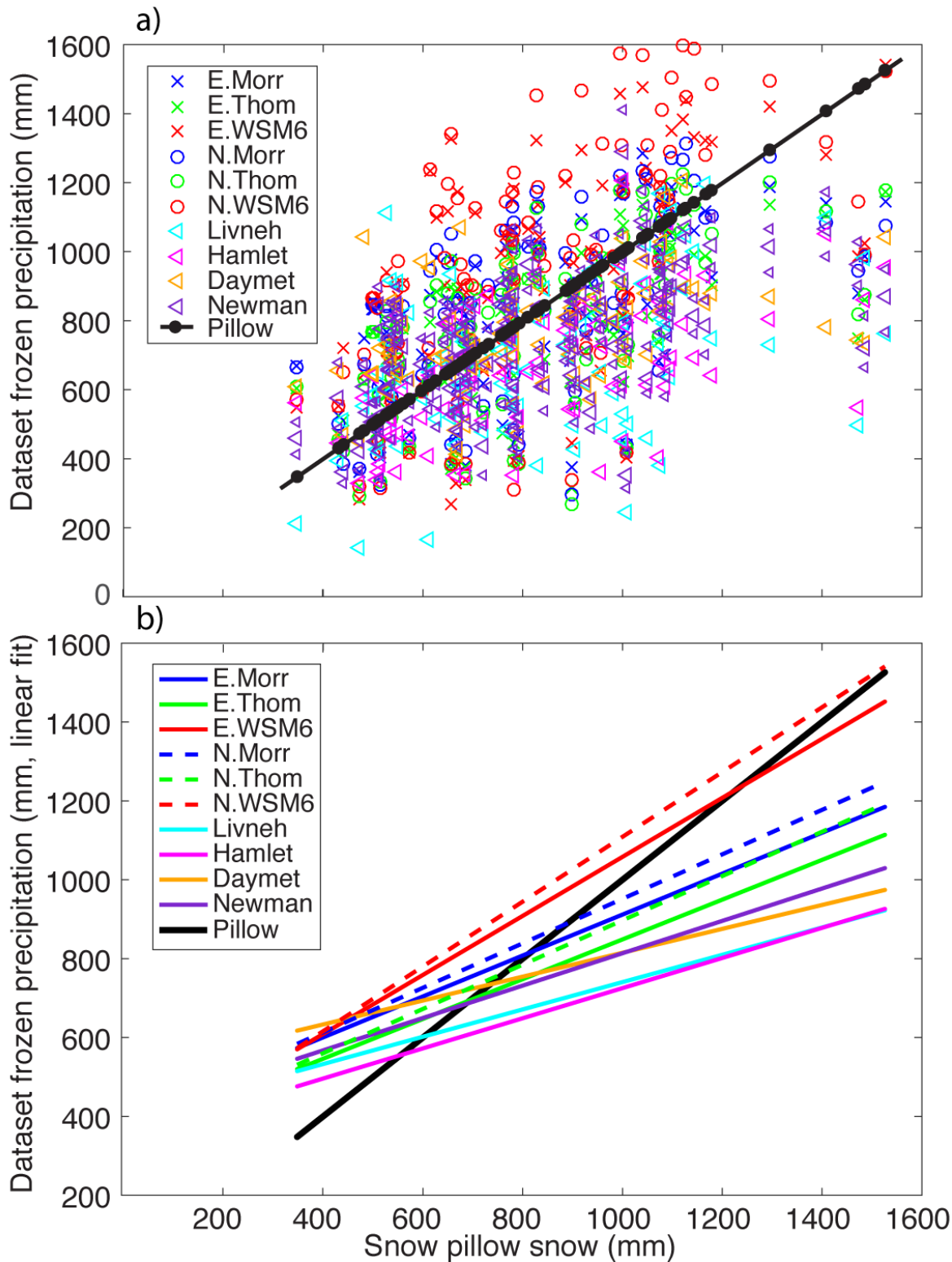


Fig. 4: a) Scatterplot of frozen precipitation versus snow pillow data at nearest gridpoint (mm). b) Linear regressions for scatterplots of (a).

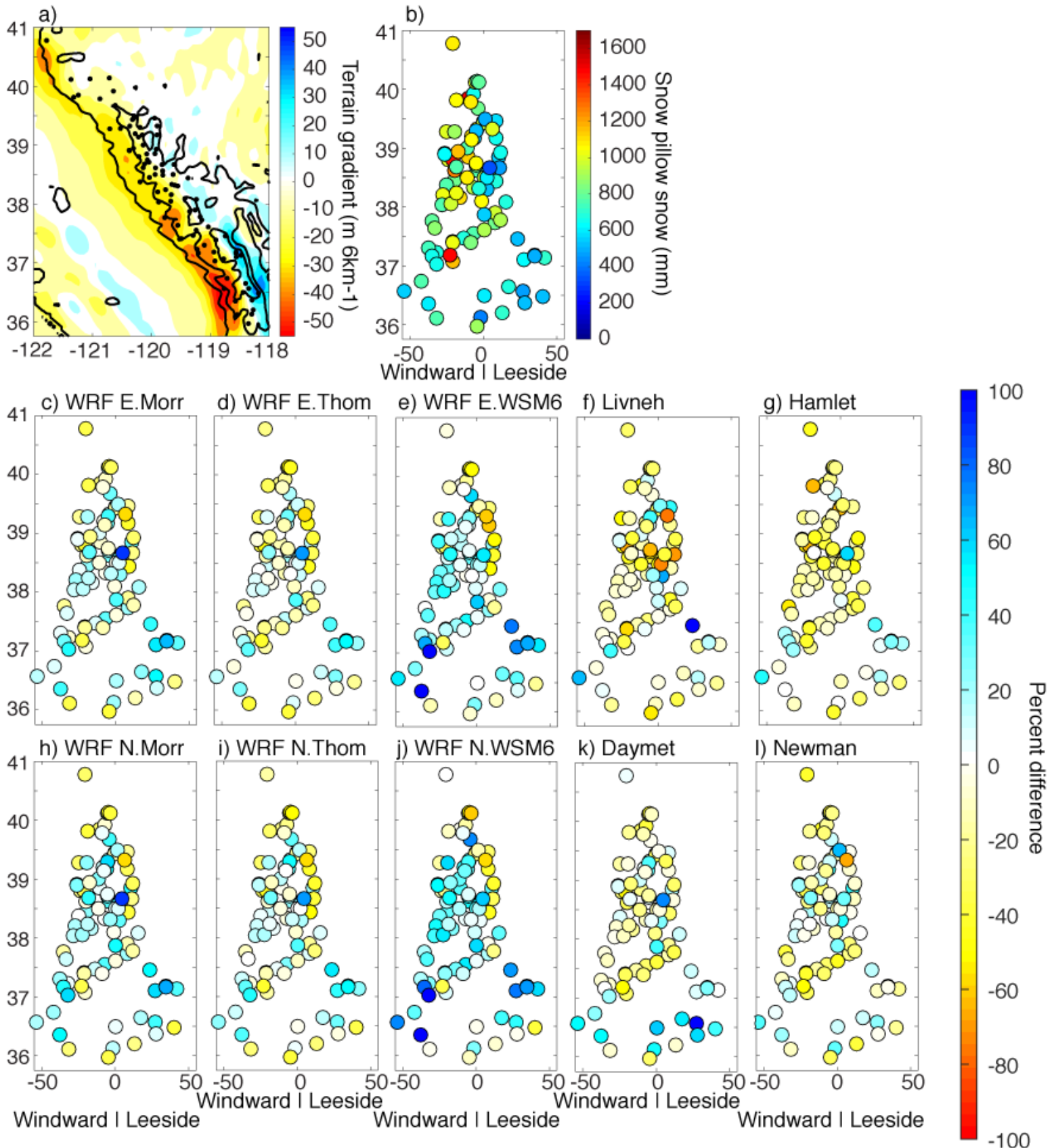


Fig. 5: a) Meridional gradient of smoothed terrain (color fill, $\text{m } 6\text{km}^{-1}$), terrain from WRF simulation (black contours, every 1000m), and locations of snow pillows (black dots). b) Snow pillow water year total snow (mm) versus smoothed zonal terrain gradient (x-axis, as in (a)) and latitude (y-axis). (c-l) As in (b), but percent difference between frozen precipitation and snow pillow snow (%).

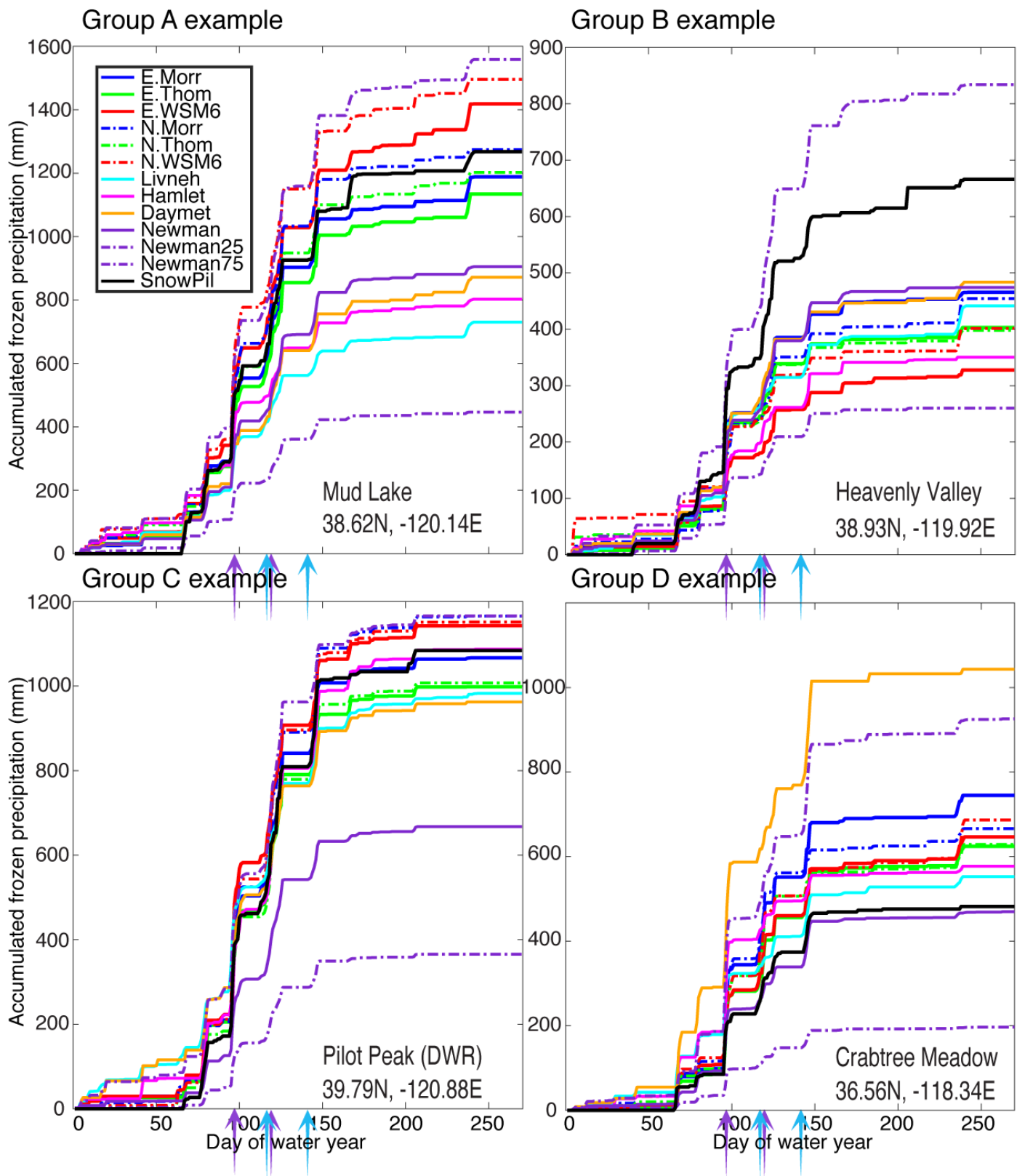


Fig. 6: Cumulative traces of daily snow pillow snow and frozen precipitation for examples from each of the 4 groups of snow pillows (A-D) outlined in the text. Purple and cyan arrows show start dates of ‘missed storms’ from Lundquist et al. (2015) in Livneh and Hamlet datasets, respectively. Locations of example pillows are shown with letters in Fig. 1b. Group A (21% of pillows): snow pillow and WRF > at least 2 gridded datasets; group B (25% of pillows): snow pillow > 8 or more datasets; group C (28% of pillows): snow pillow near center of all datasets; group D (25% of pillows): snow pillow < 8 or more datasets.

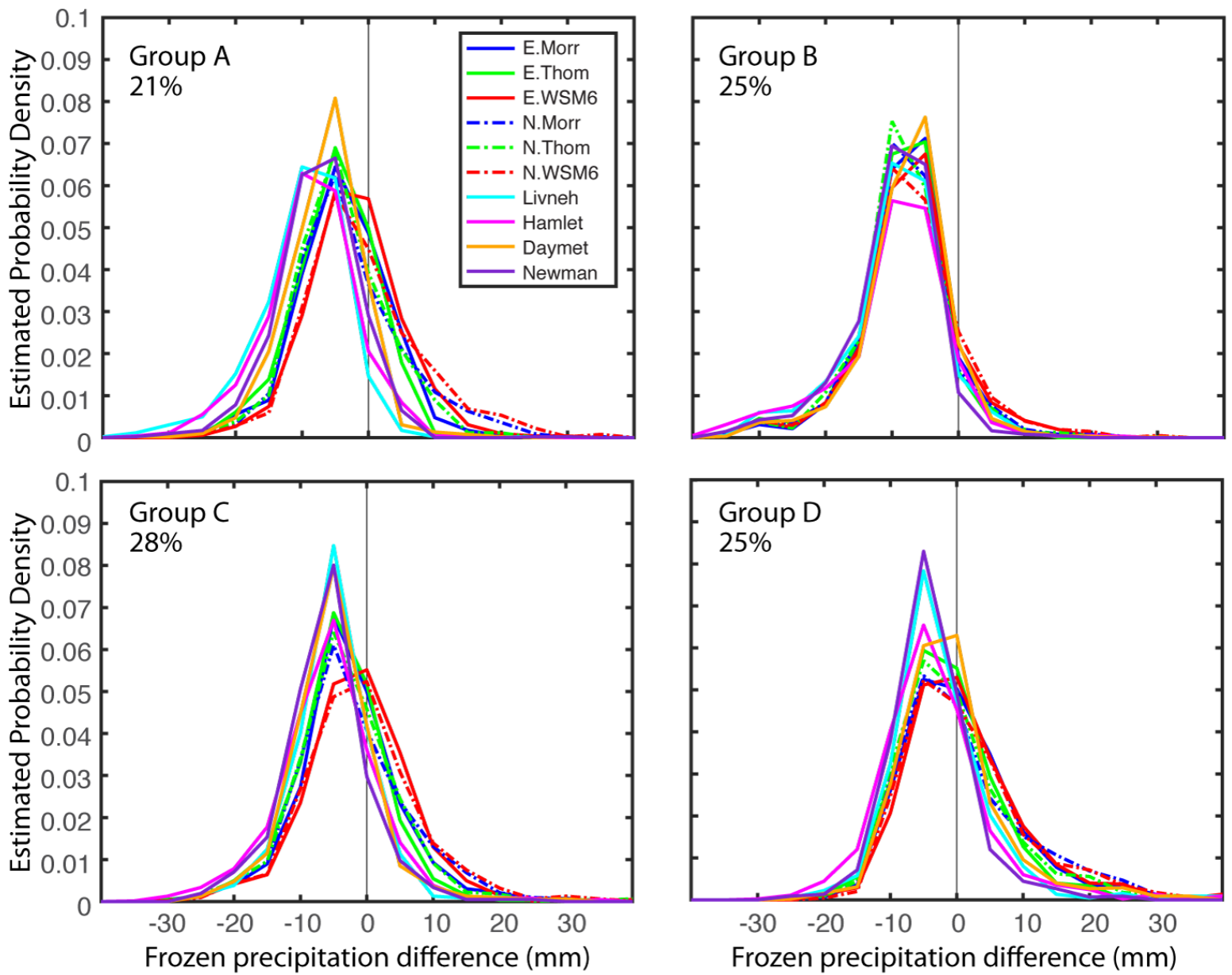


Fig. 7: Histogram of errors (gridded-snow pillow) of smoothed daily ‘frozen’ precipitation for each of the 4 groups of snow pillows (A-D) outlined in the text, on days with smoothed observed snow > 5 mm. Group A (21% of pillows): snow pillow and WRF > at least 2 gridded datasets; group B (25% of pillows): snow pillow > 8 or more datasets; group C (28% of pillows): snow pillow near center of all datasets; group D (25% of pillows): snow pillow < 8 or more datasets.

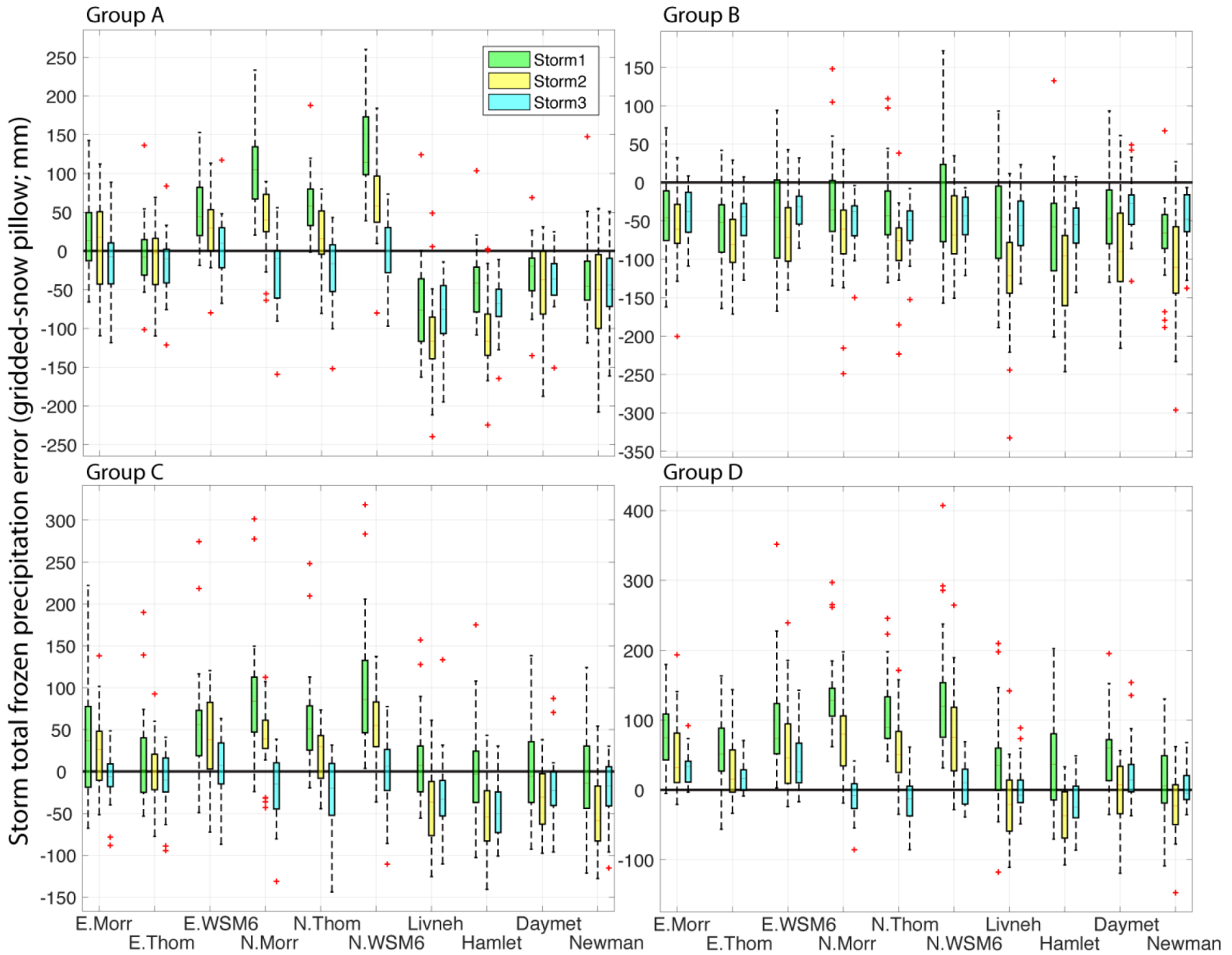


Fig. 8: Boxplots of errors (gridded-snow pillow) of ‘storm’ total ‘frozen’ precipitation (mm) for each of the 4 groups of snow pillows (A-D) outlined in the text, for three major storm periods highlighted with arrows in Fig. 5: Storm1: Jan 3-8, 2008; Storm 2: Jan. 26 – Feb. 5, 2008; and Storm 3: Feb. 19-26, 2008. Outliers shown with red + signs are +/- 2 standard deviations. Group A (21% of pillows): snow pillow and WRF > at least 2 gridded datasets; group B (25% of pillows): snow pillow > 8 or more datasets; group C (28% of pillows): snow pillow near center of all datasets; group D (25% of pillows): snow pillow < 8 or more datasets.

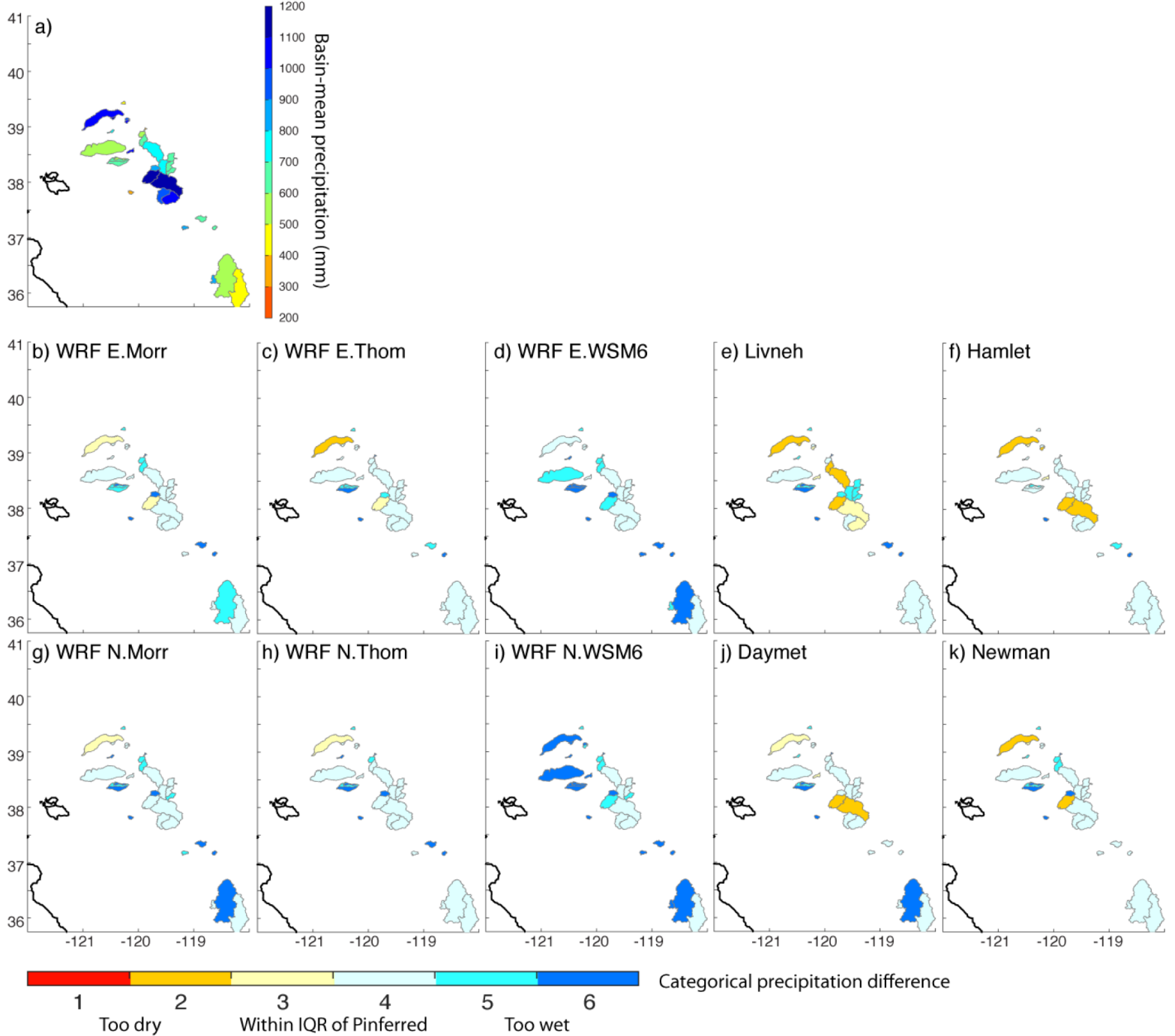


Fig. 9: (a) Median basin-mean WY2008 $P_{inferred}$ (mm). (b-k) Categorical difference of gridded dataset basin-mean precipitation (P) and $P_{inferred}$. Categories are: 1: $P < \min P_{inferred}$, 2: $\min P_{inferred} < P < 25^{th} \% P_{inferred}$, 3: $25^{th} \% P_{inferred} < P < 50^{th} \% P_{inferred}$, 4: $50^{th} \% P_{inferred} < P < 75^{th} \% P_{inferred}$, 5: $75^{th} \% P_{inferred} < P < \max P_{inferred}$, 6: $\max P_{inferred} < P$. Note that categories 3 and 4 are within the interquartile range of uncertainty.

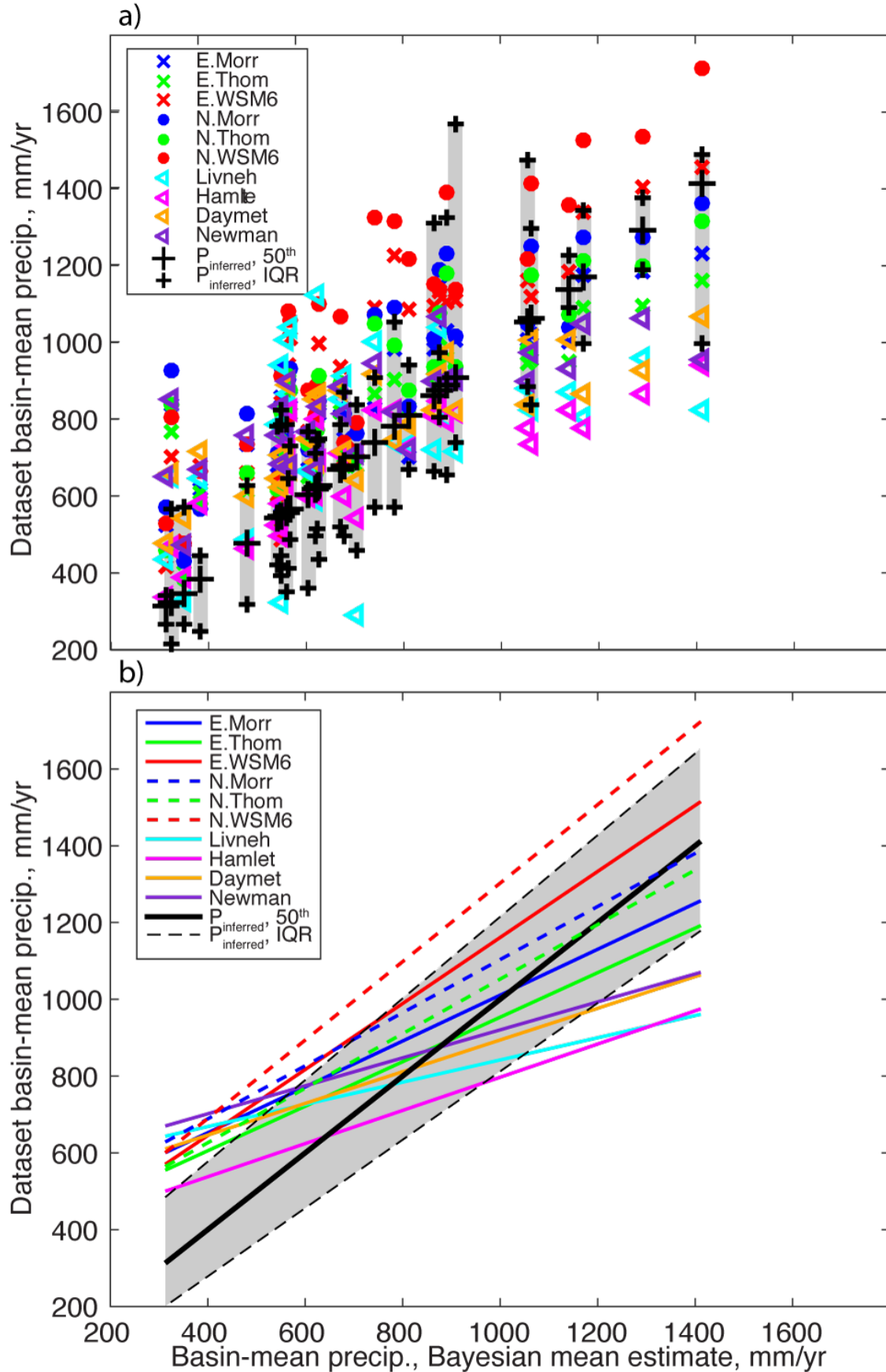


Fig. 10: (a) Basin-mean precipitation (mm) from gridded datasets (see legend), as a function of P_{inferred} (black crosses). Large black crosses show median and gray shading bounded by smaller black crosses show interquartile range (IQR) of P_{inferred} . (b) Linear fit for each dataset. Black solid line shows Bayesian median and gray shading bounded by black dashed lines show IQR.