1

2

3

4        Article type      : Technical Paper

5

6

# Evaluating Global Climate Models for Hydrological Studies of the

# Upper Colorado River Basin

9

*David W. Pierce, Daniel R. Cayan, Jordan Goodrich, Tapash Das, Armin Munévar*

Division of Climate, Atmospheric Sciences, and Physical Oceanography (Pierce, Cayan), Scripps Institution of Oceanography, La Jolla, California, USA; School of Science (Goodrich), University of Waikato, Hillcrest, Hamilton, New Zealand; Water Resources and Resilience (Das, Munévar), Jacobs, San Diego, California, USA (Correspondence to Pierce: dpierce@ucsd.edu).

**Research Impact Statement**: The latest CMIP6 generation of climate models still have biases in the Upper Colorado River Basin but show clear improvements over previous generations after a simple bias correction is performed.

**ABSTRACT:** Three generations of Global Climate Models (GCMs), CMIP3, CMIP5, and CMIP6, are evaluated for performance simulating seasonal mean and annual-to-decadal variability of temperature and precipitation in the Upper Colorado River Basin. Low-frequency precipitation variability associated with drought is a particular focus and found to be a significant model shortcoming. The evaluation includes remote teleconnected atmospheric responses to the Pacific Ocean, including the El Niño/Southern Oscillation (ENSO) and Pacific Decadal Oscillation (PDO). GCMs have improved their simulation of the Upper Basin over model

25 generations, but primarily in atmospheric circulation metrics. Persistent winter precipitation

26 biases have changed little, including in multiyear precipitation variability. Users generally bias-

27 correct GCM data before use; evaluation using a simple spatially and temporally averaged bias

28 correction shows that the CMIP6 models outperform earlier generations after the bias correction,

29 although more complex precipitation biases remain even after the simple bias correction. These

30 model rankings will be useful when selecting GCMs for a variety of hydrological and ecological

31 climate studies in the Upper Basin.

32 (KEYWORDS: Global Climate Models; Colorado River; Upper Colorado River Basin;

33 model evaluation; winter precipitation bias, regional climate model evaluation)

# INTRODUCTION

35 The Upper Colorado River Basin (UCRB hereafter) drains a multi-state area in the

36 Southwestern United States, stretching from southwest Wyoming through western Colorado and

37 eastern Utah into portions of northern Arizona and New Mexico. The UCRB is a vital source of

38 water in this largely arid region, supplying water to nearly 40 million inhabitants, irrigation to

39 5.5 million acres of farmland, and water flow to numerous wildlife refuges and national parks

40 (USBR, 2012). The Colorado River also is an important source of hydropower, capable of

41 supplying 4,200 megawatts of electricity generation to the region (USBR, 2012). Simulating

42 multi-year precipitation variability, drought, and future climate changes in the UCRB is therefore

43 of substantial societal and economic importance.

44 A significant amount of research has examined how the UCRB's annual discharge,

45 typically measured at Lees Ferry, might respond to warming, either historical or projected future

46 changes (e.g., McCabe et al. 2017; Udall and Overpeck 2017; Xiao et al. 2018; Hoerling et al.

47 2019). The propensity for drought and long-term reliability of the water supply are other

48 important concerns. Different methods have been used to examine these questions. For example,

49 some approaches use estimates of the Upper Basin's flow sensitivity to regional temperature and

50 precipitation variations, then apply projected climate changes to estimate the Basin's response

51 (e.g., Barnett and Pierce, 2008; Rajagopalan et al. 2009; Vano et al. 2012; for a review see Vano

52 et al. 2014). This approach can use runoff sensitivity estimates from observational studies or land

53 surface models to consider possible future water shortfalls (e.g., Bennett et al. 2018), and may

54 use future temperature and precipitation trends indicted by global climate model (GCMs)

55    projections (e.g., Dettinger et al. 2015). However, many regional impact studies use data from

56    one or more GCM projections as the basis for analysis, with the GCMs often being statistically

57    or dynamically downscaled to the Upper Basin to better capture important details of the regional

58    topography (e.g., Barnett et al. 2004; Christensen et al. 2004; Christensen and Lettenmaier, 2007;

59    Cayan et al. 2010; Dawadi and Ahmad, 2012; Seager et al. 2012; Ficklin et al. 2013; Tillman et

60    al. 2017). In such cases it is best to select GCMs that do a credible job of simulating the

61    historical climate and its variability in the Upper Basin, as poorly performing GCMs could

62    misrepresent processes that are important to how Upper Basin drought and discharge could

63    change in the future.

64         The purpose of this work is to evaluate the performance of GCMs in reproducing the

65    historical mean climate and variability of temperature and precipitation in the Upper Basin, with

66    an emphasis on studies of hydrology and water management in the region. Climate measures in

67    the immediate region of the Upper Basin are examined as well as remote teleconnected signals

68    associated with Upper Basin climate fluctuations, such as those originating from the tropical

69    Pacific Ocean through the El Nino/Southern Oscillation (ENSO). Seasonal, annual, and multi-

70    year timescales are considered, as they all have important implications for regional ecosystems

71    and the existing water management infrastructure.

72         There has been substantial previous work on the evaluation of GCMs using a wide

73    variety of metrics for both global and regional applications, although few that have focused

74    specifically on the Upper Basin (cf. Tamaddun et al. 2019). For example, Gleckler et al. 2008

75    evaluated global measures of performance using GCMs from the Coupled Model

76    Intercomparison Project version 3 (CMIP3) archive. A similar global analysis for the subsequent

77    generation of GCMs, CMIP5, appears in Flato et al. 2013. These global evaluations consider

78    such aspects as the Earth's radiation fields, surface precipitation and temperature, and winds,

79    pressures, and temperatures at key vertical pressure levels in the atmosphere. Pierce et al. 2009

80    performed a similar evaluation using the CMIP3 models but focusing on the western United

81    States; the procedure used here is based on the approach developed in that work. Rupp et al.

82    (2013) performed a CMIP5 GCM evaluation for the Pacific Northwest, and subsequently for the

83    Southwest United States as reported by California Department of Water Resources (CA DWR

84    2015). Knutti et al. (2017) examined how GCMs can be weighted to take into account model

85  quality scores, such as developed here, when analyzing a multi-model ensemble. Lorenz et al.

86  explore weighted multimodel ensemble predictions of summer maximum temperature over North

87  America, and Brunner et al. (2019) examine temperature and precipitation projections over

88  Europe using a model weighting scheme that incorporates model performance and independence.

89        The evaluation shown here differs from those previous efforts in several key areas. First,

90  it is focused on the UCRB specifically. Second, the analysis includes both the older-generation

91  CMIP3 and CMIP5 models as well as the newer CMIP6 models. Third, we compare the

92  performance of different variables in an absolute sense, an effort that previously has generally

93  been avoided in favor of relative measures of performance. This point will be explained in more

94  detail below.

95        Since one of our key purposes is to evaluate hydroclimatic features, we include

96  evaluations of multi-year precipitation variability important to drought processes. Numerous

97  other studies that have examined such processes in various regions, such as Rupp et al. 2013 and

98  Abatzoglou and Rupp 2017 in the Pacific Northwest, and Moon et al. (2018) and Ukkola et al.

99  (2018) for global evaluations. Global climate teleconnections are also included since they are of

100  first-order importance to climate variability in the region. Such teleconnected responses to the

101  western U.S. have previously been considered by Pierce et al. (2009) and Rupp et al. (2013), for

102  example.

103        One question we examine is whether the performance and quality metrics indicate that

104  previous generation models (such as from the CMIP3 archive) should be discarded from

105  consideration. This question is relevant because the CMIP3 models show, on average, drier

106  future conditions in the UCRB than the more recent models (Ficklin et al. 2015). Either

107  arbitrarily excluding or unjustifiably including the CMIP3 models could bias understanding of

108  future drought in the Upper Basin. Evaluations of model skill improvement over the CMIP

109  generations that examine global measures, rather than the UCRB-specific metrics considered

110  here, can be found in Bock et al. (2020) and Fasullo et al. (2020).

111        The current work focuses on the GCMs' representation of temperature and precipitation

112  in the UCRB and teleconnected responses to the tropical and North Pacific. Land surface models

113  (LSMs) are a key part of GCMs and have evolved considerably over the model generations.

114  Because of their importance, LSMs and their responses to climate change in the CMIP models

115    have been examined in their own right (e.g., Boone et al. 2009; Dirmeyer et al. 2013, Li et al.

116    2018; Li et al. 2021). LSM fields such as runoff and streamflow are not examined here but were

117    addressed as part of this project and will be reported at a later date.

118    Results from this analysis can be used to inform model selection for a variety of climate

119    impact studies in the Upper Basin. Although our focus is on hydrology and drought, ecosystems

120    are also strongly affected by local temperature and precipitation so GCM selection is important

121    for ecological application as well. Likewise, human health and regional energy demand will be

122    impacted by future temperature changes, so GCM-based studies in those fields could employ the

123    evaluation developed here.

# DATA AND METHODS

125    *Variable Selection*

126    We obtained GCMs data from three generations of GCMs from the Climate Model

127    Intercomparison Project (CMIP), referred to as CMIP3 (Meehl et al., 2007), CMIP5 (Taylor et

128    al., 2012), and CMIP6 (Eyring et al. 2016). GCMs produce a wide variety of variables describing

129    the state of the atmosphere, ocean, land surface, and cryosphere, although it is not feasible to

130    save all variables in the CMIP archives and variable coverage is smaller in the earlier CMIP

131    generations. The current work focuses on the GCMs' performance in temperature, precipitation,

132    and teleconnections associated with temperature and precipitation. The relatively coarse

133    resolution of CMIP GCMs yield a poor simulation of land surface processes such as snowpack

134    and soil moisture in a topographically diverse region such as the UCRB. Dynamically or

135    statistically downscaled data, not examined here, is generally better suited to examining such

136    surface fields in a geographically limited, rugged region.

137    *Global Climate Models*

138    We evaluate data from 82 GCMs: 16 CMIP3, 35 CMIP5, and 31 CMIP6 models, as

139    shown in Table 1. The last column shows the approximate spatial resolution of the model's

140    atmospheric data files as they appear in the CMIP archive. The North American CMIP3 and

141    CMIP5 data were obtained from the U.S. Bureau of Reclamation (Reclamation hereafter) archive

142    of climate model output available from Lawrence Livermore National Laboratory's Green Data

143    Oasis archive (https://gdo-dcp.ucllnl.org/downscaled_cmip_projections/dcpInterface.html).

144  Additional CMIP3 and CMIP5 data, and all the CMIP6 data, were downloaded from the Earth

145  System Grid (e.g., https://esgf-node.llnl.gov/search/cmip6/) in mid-to-late 2020. We only include

146  models that provide daily fields of minimum and maximum temperature (Tmin and Tmax) and

147  precipitation, required for hydrological modeling work not described here. Additionally, we only

148  include models that have data for both a historical and future climate change simulation. Several

149  CMIP6 models have historical data available but no future shared socioeconomic pathway (SSP;

150  Raihi et al. 2017) simulation, and so were excluded. Other models lack daily data over the

151  historical or future period, and likewise were not analyzed (in particular, at the time of writing

152  the CESM2 family of models do not provide daily Tmin/Tmax over the historical period).

153      The CMIP3, CMIP5, and CMIP6 generations used different historical periods, ending in

154  1999, 2005, and 2014, respectively. As a compromise between excluding recent data and using a

155  different analysis period for all 3 generations, we used a historical period of 1950-1999 for

156  CMIP3 and 1950-2005 for the CMIP5 and CMIP6 models. Monthly mean daily average

157  temperature (Tavg) was formed as the mean of monthly averaged Tmin and Tmax and is the

158  temperature quantity analyzed here.

159      The ensemble members used are shown in Table 1. Sea level pressure from only the first

160  realization was available for the CMIP3 models. Only historical realizations are shown in the

161  table since the model/observations comparison only uses data over the historical period. Each

162  ensemble member was evaluated on all metrics, and then the final metric for each model was

163  taken as the mean of the values for all the ensemble members. This approach prevents models

164  with many ensemble members from having undue influence on the results. Additionally, the

165  spread across the ensemble members was used to quantify uncertainty.

166      [TABLE 1 GOES HERE]

167      Data in Reclamation's archive had been re-gridded to a common 2-by-2-degree latitude-

168  longitude grid for the CMIP3 models and a common 1-by-1-degree grid for the CMIP5 models.

169  To examine the models on the same grid and explore the effect of spatial resolution on our

170  results, we interpolated the CMIP6 and CMIP3 models to the same 1-by-1-degree grid via

171  bilinear interpolation and aggregated the CMIP5 and CMIP6 data to the 2-by-2-degree grid. We

172  found that whether the 1x1 or 2-by-2-degree grid is used makes only a minor difference in the

173    final ranked model quality results, so most of the results here will be shown using the 1x1-degree

174    gridded data.

175    *Observations*

176         Daily temperature and precipitation over North America were obtained from Livneh et al.

177    2015 (Livneh hereafter), a gridded product based on airport and cooperative weather stations.

178    The data cover central Mexico through southern Canada at a 16th-degree latitude-longitude

179    resolution over the period 1950-2013, which was trimmed to 1950-2005 to match the CMIP5 and

180    CMIP6 model historical periods. Values were aggregated to the same common 1x1 and 2x2

181    degree grids as the models. Daily minimum and maximum temperature were averaged to

182    produce daily average temperature, then averaged to monthly values to match the GCM data.

183    Massmann (2020) shows that Livneh does well in representing temperature and precipitation

184    across the CONUS for the purposes of hydrological modeling. Pierce et al. (2021) find that

185    Livneh precipitation extremes on a daily timescale are distorted by the data processing

186    methodology, but this does not affect the monthly-averaged analysis performed here.

187         For global observations of monthly sea level pressure (SLP) and temperature we used the

188    ERA5 reanalysis (Hersbach et al. 2018) over the period of 1950-2005. Although historical

189    station-based estimates of monthly temperature exist they are not spatially complete, so the

190    reanalysis data was used in preference. In comparisons with the older NCEP reanalysis product

191    (Kalney et al., 1996), some minor differences in model ranking were found in the metrics

192    sensitive to global SLP when using the ERA5 vs. NCEP reanalysis. This shows that

193    observational uncertainty can affect model ranking, but this aspect of uncertainty is not explored

194    in the current work (cf. Lorenz et al. 2018).

195              CULLING OF GCMs BASED ON GLOBAL METRICS

196         GCMs have a wide range of performance, and it is not consistent which model performs

197    best on which metric (e.g., Gleckler et al., 2008). However, some models perform systematically

198    worse than the other GCMs across a range of global metrics. To alleviate the concern that a

199    highly targeted, Upper Basin-centric analysis might select models that do well in this small

200    region but poorly in simulating the overall Earth's climate, an initial culling was performed using

201    published hemispheric to global scale metrics to eliminate the bottom-performing 25 percent of

202    models. This was done separately for the CMIP3, CMIP5, and CMIP6 GCMs. Most of the

203    figures in the main text use this culled set of models. For key figures (called out below), the

204    supplementary information contains figures made using the full, un-culled set of models for

205    comparison.

206    The CMIP3 culling was based on Gleckler et al. (2008), specifically the model

207    performance in the Northern Hemisphere extra-tropics (their Figure 3d). This resulted in the

208    elimination of the following four models: ipsl_cm4, giss_model_e_r, ncar_pcm, and imncm3_0.

209    The CMIP5 culling was based on Flato et al. (2013), which eliminated the following nine

210    models: GISS-E2-R, IPSL-CM5A-MR, inmcm4, FGOALS-g2, bcc-csm1-1-m, MIROC-ESM-

211    CHEM, MIROC-ESM, IPSL-CM5A-LR, and IPSL-CM5B-LR. The CMIP6 culling was based

212    on three sources: Brunner et al. 2020, who evaluated the GCMs for performance and

213    independence from each other; Tokarska et al. 2020, who used emergent constraints to identify

214    models that have historical warming inconsistent with observations; and the online analysis by

215    the Program for Climate Model Diagnosis and Intercomparison (PCMDI), available at

216    https://cmec.llnl.gov/results/physical.html (accessed Jan 27, 2021). A subjective evaluation of

217    the results of those three studies resulted in the following 8 models being eliminated: CanESM5,

218    HadGEM3-GC31-LL, NESM3, UKESM1-0-LL, FGOALS-g3, INM-CM4-8, NorCPM1, and

219    BCC-ESM1. After the culling, 61 GCMs remained for the subsequent analysis (12 CMIP3, 26

220    CMIP5, and 23 CMIP6). Several Supporting Information figures show key results for the entire,

221    un-culled set of models for interested readers.

222    As mentioned previously, the CMIP3 model projections tend to show drier end-of-

223    century conditions in the Upper Basin than the CMIP5 and CMIP6 models. Does the culling,

224    which is based only on historical data, affect this outcome? This is examined in Figure 1, which

225    shows histograms of GCM-projected temperature (red) and precipitation (green) changes in the

226    Upper Basin, both before (top row) and after (bottom row) the global culling. Models with

227    multiple ensemble members use the ensemble-mean result, so that different models are weighted

228    equally in the figure even if they have different numbers of ensemble members. The distribution

229    of projected temperature increases becomes narrower after the global culling, with both the

230    greatest and least warming models culled. The central peak also becomes notably more

231    pronounced. The precipitation distribution, by contrast, is less affected although some of the

232     extreme wet models are culled. A similar pattern is seen in the seasonal results (Supporting

233     Information, Figure S1), with the culling affecting the standard deviation of the projected

234     temperature change considerably more than the projected precipitation trend. The seasonal

235     results also show the most warming in summer as the surface dries, while the greatest

236     precipitation trend is in winter. Natural variability affects the spread of projected trends,

237     especially in a region as small (in a global sense) as the UCRB. Future work using models with

238     large-ensemble simulations could help distinguish between spread based on natural variability

239     and spread due to differences between models.

240        [FIGURE 1 GOES HERE]

## Model Metrics

242        Metrics were developed to evaluate model historical performance in simulating regional

243     precipitation and temperature characteristics and teleconnections of Pacific surface temperature

244     and large scale atmospheric sea level pressure that relate to Upper Basin precipitation variations.

245     The metrics used for these evaluations are based on spatial fields of z-scores, i.e., a spatial field

246     of differences between the model and observations normalized by a measure of variation in the

247     observed value. The normalization allows for model-observed differences to be sensibly

248     evaluated—are they large or small compared to observed variability?

249        Specifically, for seasonal-mean quantities, such as DJF seasonal mean precipitation, the

250     z-scores comparing model to observations are calculated as follows:

251        $$z_{score}(x,y) = \frac{(\ <model\ (x,y,t)> -\ <obs(x,y,t)>\ )}{\text{stddev}(obs(x,y,t))} \qquad\qquad \text{Eq. 1}$$

252        Where the angle brackets $<>$ indicate averaging the time sequence of seasonal means

253     over time, stddev() indicates the standard deviation over time, and *model(x,y,t)* and *obs(x,y,t)* are

254     the time series of seasonal mean data in the model and observations, respectively. This approach

255     non-dimensionalizes the errors, so that measures with different units (e.g., temperature and

256     precipitation measures) can be sensibly compared.

257        The z-score is calculated at each point in the domain, yielding a 2-dimensional map of z-

258     scores. The final overall skill score ss, or metric value, is then calculated as:

259        $$ss = 1 - \text{RMS}(z_{score})$$

260     where RMS indicates the root mean square spatial average over the domain. We do not

261     include area weighting of the grid cells in this work because of the limited domain size, but it

262     could be included if analyzing a more extensive region. This formulation follows the traditional

263     skill score scaling with 1 being a perfect skill, and more negative values indicating less

264     agreement with the observations, measured by the yardstick of observational variability. A skill

265     score of zero means that the model and observations have an RMS mean difference of 1 standard

266     deviation over the domain.

267     As a detailed example, consider JJA mean precipitation over the region evaluated for a

268     CMIP5 GCM. There are 56 seasons of observations for this quantity (since the historical period

269     is 1950 to 2005). At each point in the domain, the model and observed mean are calculated over

270     the 56 seasons. The sample standard deviation of the observations is then calculated at each point

271     from the 56 seasonal values. The difference between the model and observed mean at each point,

272     divided by the sample standard deviation at that point, yields the z-score at that point. The final

273     skill score is 1 minus the RMS of the z-scores over the domain.

274     Measures of variability require a slightly modified approach to calculate the observed

275     standard deviation. For example, consider the metric of standard deviation of winter (December,

276     January, February [DJF]) precipitation averaged into 10-year blocks, which has a direct

277     relationship to drought. The monthly data are first averaged over consecutive, non-overlapping,

278     10-year blocks and the standard deviation computed at each point. This value is easily obtained

279     from the observations and models, but then the denominator in the z-score needs to be

280     calculated, i.e., what is the typical spread in the estimate of the standard deviation of 10-year

281     blocks of precipitation? This was estimated using a block bootstrap method where 100 random,

282     50-year long sequences of 10-year blocks were constructed (with replacement) from the actual

283     sequence of years, with the 10 years in each block being sequential calendar years. The sample

284     standard deviation of 10-year blocks is computed for each random sequence, and the distribution

285     of standard deviations examined. We use a similar bootstrap method to estimate the spread in

286     observed quantities whenever that quantity was needed and not available by direct computation.

287     Note that some of the metrics relating to teleconnected climate responses (for example, from

288     Pacific Ocean sea surface temperatures to UCRB precipitation) are multivariate. In these cases,

289  for each random trial the same randomly-determined temporal ordering was used for all relevant
290  variables to preserve the temporal coherence of the fields.

291  With only a limited time span of data available to analyze (56 years), the sampling errors
292  of the low frequency (5- and 10-yr averaged) metrics are higher than the seasonal metrics. This
293  reduces the reliability of the low-frequency metrics compared to the seasonal metrics. The
294  current analysis does not attempt to down-weight the low frequency metrics to account for their
295  reduced reliability, although such an approach could be useful. For example, Rupp et al. (2013)
296  discussed the issue of metric reliability and accomplished this by effectively setting the weight of
297  unreliable metrics to zero. Metric reliability and how to combine it with the other measures of
298  metric and model ranking uncertainty shown below are not considered here, although it would be
299  a useful direction for further research.

300  The use of z-scores as the basis for our metrics is ultimately why we can compare
301  dissimilar variables, allowing us to do absolute comparisons of the error in different metrics
302  rather than only relative errors. For example, the question of whether a model does a better job
303  simulating mean winter precipitation or summer temperature variability can be quantitatively
304  answered by noting that (as a hypothetical example) a model's mean winter precipitation field
305  may be, on average, 2 standard deviations away from the observed value, while the summer
306  temperature variability errs by only 1 standard deviation. It is in this sense that we can say that
307  the simulation of the precipitation field is worse than the simulation of temperature variability.

308  The absolute approach to metrics developed here is a departure from most previous
309  efforts at evaluating model quality. More commonly, relative measures of error are used, such
310  that error measures in different variables are normalized to have the same range. However, the
311  relative approach has the disadvantage that a metric with a wide range of values, spanning a
312  range from a very good to a very poor simulation, has as much influence on the model ranking as
313  a metric that is well simulated by all models. Therefore, the relative approach discards useful
314  information that could better discriminate between models.

315  Uncertainty in estimating metrics from the limited available time period of observations
316  can contribute to a large sample standard deviation in the observations, yielding a lower z-score
317  (all else being equal). Therefore, a smaller error in a well-known or stable quantity can

318   potentially yield a poorer (larger) z-score than a larger error in a poorly known or highly variable

319   quantity.

320   *Seasonal Means and Variability in the Upper Basin*

321   Of fundamental importance to the evaluation performed here is the ability of the GCMs

322   to simulate the mean climate and variability in the Upper Basin. Accordingly, 32 metrics are

323   used to evaluate the seasonal (DJF, March, April, May [MAM], JJA, September, October,

324   November [SON]) mean temperature and precipitation in the region, and the standard deviation

325   of those seasonal quantities evaluated in 1-, 5-, and 10-year blocks. These values are compared to

326   observations at each point in the domain, which is shown in Figure 2.

327   [FIGURE 2 GOES HERE]

328   An example of the metric fields for winter (DJF) precipitation is shown in Figure 3. The

329   upper left shows the observed winter precipitation, in millimeters per day, while the lower left

330   shows the standard deviation of the observations, used to form the z-score. The middle column

331   shows a model that does well on this measure, CanESM2, with the precipitation field in the top

332   row and the z-score in the bottom row. By contrast, the right column shows a model that does

333   poorly on this measure, FIO-ESM. In this and the following examples (including in the

334   Supporting Information) we do not always select the "best" and "worst" models to show, which

335   indeed our results indicate are subject to uncertainty, but rather show results from a variety of

336   well- and poorly-performing models rather than repeating the same models. Although even

337   CanESM2 does not capture much of the spatial pattern, at least the values are reasonably close to

338   the observations. FIO-ESM, on the other hand, is like many other GCMs in simulating far more

339   winter precipitation than observed. The occurrence of relatively large errors in mean

340   precipitation across many of the GCMs is likely due to the poor representation of topography. An

341   analogous example for summer temperature variability is shown in the Supporting Information,

342   Figure S2.

343   [FIGURE 3 GOES HERE]

344   *Amplitude and Phase of the Seasonal Cycle*

345   The amplitude and phase of the seasonal cycle are calculated for monthly temperature

346   and precipitation, yielding another four metrics. These are calculated from the best-fit annual

347  sinusoid. However, it should be noted that precipitation in the Upper Basin has a relatively weak

348  seasonal cycle. The seasonal cycle of precipitation is a useful metric even though the observed

349  seasonal cycle is weak since a model that had a pronounced seasonal cycle (unlike the

350  observations) would yield a poor simulation of the region. The observed spread in the amplitude

351  and phase estimates were formed by the bootstrap method described above.

352  *El Nino/Southern Oscillation and the Pacific Decadal Oscillation*

353      Various modes of climate variability are associated with temperature and precipitation

354  variations in the Upper Basin, notably the El Nino/Southern Oscillation (ENSO; e.g., Hidalgo

355  and Dracup, 2003) and the Pacific Decadal Oscillation (PDO; e.g., McCabe and Wolock 2020).

356  The fidelity of the models' depiction of these modes was evaluated in three aspects: their mean

357  expression in sea surface temperature (SST) anomalies in the Pacific Ocean; the variance

358  spectrum of the SST pattern; and the associated response in temperature and precipitation over

359  the western United States. SST is computed from near-surface air temperature by including air

360  temperature in ocean regions only and limiting the lowest allowable temperature to the freezing

361  point of seawater, $-1.8$ °C. Lower values generally indicate the presence of sea ice. We chose

362  this route rather than trying to download ocean model SSTs directly because of the substantial

363  gain in efficiency gained by only downloading one set of temperature files, combined with fact

364  that GCMs have very similar SST and 2-meter temperature fields when evaluated on a common

365  1x1 degree grid.

366      The observed and model ENSO indices are taken as the principal component associated

367  with the leading empirical orthogonal function (EOF) of monthly SST anomalies over the region

368  135 E to 80 W, 10 S to 10 N. The EOF is taken in preference to a construction based on the so-

369  called Nino regions, such as Nino 3.4 (170 W to 120 W, 5 S to 5 N), because models do not

370  necessarily capture the correct spatial pattern of ENSO. Using the EOF means that the model's

371  own representation of ENSO SST patterns is used as the basis of that model's ENSO index. The

372  observed SST pattern, as well as examples from an ensemble member of models that have a

373  relatively good (CESM1-CAM5) and bad (CSIRO-Mk3-6-0) simulations of the ENSO pattern

374  are shown in Figure 4. Both models extend the variability too far to the west, a common problem

375  with GCMs, but CSIRO-Mk3-6-0 does this much more than CESM1-CAM5. The z-scores for

376　CSIRO-Mk3-6-0 in the far western tropical Pacific exceed 5 standard deviations, compared to <

377　3 for CESM1-CAM5.

378　　　　[FIGURE 4 GOES HERE]

379　　　　The variance spectrum of the principal component associated with the leading EOF is

380　used as a metric in addition to the SST anomaly pattern. This is a useful metric of model quality

381　because some GCMs simulate a very regular 2-year ENSO cycle that is unlike the observed

382　irregular cycle with enhanced variability at periods between 2 and 7 years. Rather than a spatial

383　pattern of differences between the model and observed value, as done with the measures of

384　model quality described previously, the difference between the logarithms of the model and

385　observed power at frequencies between 2 and 7 years per cycle is computed. The logarithms are

386　used so that the model having twice the power as the observations gives the same error as the

387　observations having twice the power as the model.

388　　　　The teleconnected precipitation and temperature response associated with the SST pattern

389　is taken over the entire North American domain west of 105 W (25.5 N to 52.5 N, 150 W to the

390　coast). This broader domain was used in preference to only the Upper Basin domain because

391　inspection showed significant structure in the teleconnected fields over this wider region, while

392　the Upper Basin tends to straddle the zero line of the response. These teleconnections are

393　evaluated over the cold season only (ONDJFM), when the teleconnected precipitation and

394　temperature signal to North America is strongest. Teleconnections have ramifications over the

395　warm season as well but are poorly simulated during this season in current models (Jong et al.

396　2021). The teleconnected response pattern over the Upper Basin is determined by linear

397　regression between the leading principal component and the response field of interest over the

398　Upper Basin (precipitation or temperature). This approach assumes that the teleconnected

399　response is linear in the associated SST pattern (e.g., the El Nino response is the opposite of the

400　La Nina response), which is likely untrue, but our attempt at using composites to capture this

401　non-linearity resulted in a low signal to noise ratio due to splitting the data into three pieces. An

402　example for the teleconnected response in precipitation to ENSO variability is shown in the

403　Supporting Information, Figure S3.

404　　　　The PDO is evaluated in the same way as ENSO, with the leading EOF of SST anomalies

405　taken as the PDO index, but in this case the domain is 145 E to 110 W, 20 N to 55 N. The

406     observed SST pattern and examples of well and poorly performing models are shown in the

407     Supporting Information, Figure S4. An example for the teleconnected response in precipitation to

408     PDO variability is shown in the Supporting Information, Figure S5.

409     Each of the climate modes (ENSO, PDO) contributes four metrics (mean SST pattern,

410     spectrum, teleconnected cold season response over the western United States in temperature and

411     precipitation), for a total of eight metrics.

412     *Remote Correlations with Upper Basin Precipitation*

413     There are two ways to evaluate the connection of Upper Basin climate variability to

414     wider hemispheric or global fluctuations. One method, described in the last section, is to

415     examine the effect of known climate modes of variability (such as ENSO and the PDO) on the

416     Upper Basin. The other way is to start with precipitation fluctuations in the Upper Basin and

417     examine how other fields correlate with those fluctuations.

418     The latter approach was implemented by forming the time series of cold (ONDJFM) and

419     warm (AMJJAS) season precipitation in the Upper Basin, then correlating those time series with

420     temperature and sea level pressure fluctuations elsewhere around the globe. Examination of the

421     observed correlation maps suggested that a suitable domain to evaluate the correspondence

422     between model and observations is 100 E to 60 W, 10 S to 60 N. The variability was evaluated

423     by the bootstrap method. This field is an example where relatively large sampling variability

424     with respect to the signal leads to low z-scores with a comparatively weak ability to distinguish

425     between models. Examples of the patterns for well and poorly performing models are shown in

426     the Supporting Information, Figures S6 through S9.

# Results

428     Before describing the results, we emphasize a few key points on interpreting skill scores.

429     1) There is no absolute guide as to which metrics to pick to describe diverse aspects of the

430     climate system. This must be guided by experience with the study domain and the aspects of

431     climate relevant to the problem of interest. 2) Skill score values are 1 minus the RMS average z-

432     score of the model error, averaged over the Upper Basin domain. Loosely, positive skill scores

433     indicate that the model biases are no larger than typical fluctuations due to natural variability,

434     while negative skill scores mean that biases are larger than typical variability.

435    The overall portrait plot of model skill scores for each metric is shown in Figure 5. A
436    summary model skill score across all the metrics, indicated by the number in the parenthesis after
437    the model name, is constructed as the Euclidian distance between the perfect model skill point
438    (1, 1, 1, …, 1) and the model's skill scores in all the metrics. Lower values therefore signify
439    better models. This is referred to as Dss, signifying the distance in skill score space. The models
440    are ordered in Figure 5 such at that the best models are at the bottom of the plot (smallest
441    distance to the perfect skill point, therefore lowest Dss), and the worst models at the top of the
442    plot (largest Dss). For models with more than one ensemble member, the ensemble mean value is
443    shown. Uncertainty in the model rankings estimated by spread across the ensemble members will
444    be shown below (Figure 8 and Figure 11).

445    It was previously noted that metrics with large observational uncertainty (particularly the
446    poorly sampled 10-year average metrics) give a large denominator in Eq. 1, yielding skill scores
447    near zero. In other words, if the observational uncertainty is large, it cannot be definitively
448    concluded that model results are inconsistent with the observations, leading to skill scores that
449    are near zero. By contrast, it can be easier to conclude that model-observational differences are
450    large in well-observed quantities with low observational uncertainty, leading to negative skill
451    scores (i.e., it is known that the models are inconsistent with the observations). This can be seen
452    in Figure 5, where the skill scores in 1-yr variability tend to be more negative than the skill
453    scores in the 5- and 10-yr averages but is an outcome of larger uncertainty in the poorly sampled
454    low-frequency metrics. For example, Abatzoglou and Rupp (2017) found that GCM fidelity was
455    generally lower-on multi-year timescales than seasonal or annual timescales when evaluating
456    CMIP5 GCM simulations of drought in the Pacific Northwest.

457    An analogous figure including all models (no culling) is given in the Supporting
458    Information, Figure S10. The culling eliminates some models that would otherwise score well in
459    the UCRB. For example, HadGEM3-GC31-LL is culled on the basis of poor global performance
460    in the evaluation of Brunner et al. (2020), their Figure 4. This finding indicates the importance of
461    considering global metrics even for regional GCM applications, as models that perform poorly
462    on global metrics may be doing well in the UCRB but for the wrong physical reasons.

463    [FIGURE 5 GOES HERE]

464    One striking aspect of Figure 5 is that some metrics have consistently low skill across

465    many models, seen as dark blue vertical columns, particularly metrics associated with winter

466    precipitation variability. Somewhat ironically, the models that stand out for being much better

467    than normal on the winter precipitation variability metric, EC-Earth3 and EC-Earth3-Veg, do

468    unusually poorly on winter and spring temperature variability. The uneven range of skill score

469    variability is illustrated in Figure 6, which shows the distribution of metric values, sorted by the

470    mean metric value, with better simulated metrics having higher means (closer to the perfect

471    value of 1).

472    [FIGURE 6 GOES HERE]

473    The best simulated metric is the annual phase of Upper Basin temperature (tas phase),

474    which is not surprising since the phase is largely controlled by solar insolation and the tilt of the

475    earth's surface with respect to the sun, quantities that are specified in the models. Of the 8 worst

476    metrics, 6 are associated with precipitation variability. Winter precipitation variability on the 1-

477    yr time scale is by far the worst simulated quantity, with a mean metric value less than -5. This

478    likely is influenced by poor model treatment of topography in the Upper Basin. As noted earlier,

479    some metrics with a relatively large spread in the observations, such as the estimated spectral

480    power in ENSO and the PDO, do relatively well in the sense that the model values cannot be

481    shown to be outside the wide range of uncertainty.

482    *Estimating Uncertainty using Ensemble Members*

483    Models with multiple ensemble members can be used to explore how sampling and

484    model-simulated natural climate variability affect the metric scores and overall model ranking.

485    Supporting Information Figure S11 shows the estimated standard deviation for each metric,

486    ranked from most to least certain. This was estimated in two ways: 1) by calculating the standard

487    deviation of each metric from every model with at least three ensemble members, then taking the

488    mean of the model estimates as the final standard deviation (referred to as the mean of the model

489    values); 2) by forming, for each metric, the anomalies of each model's skill scores with respect

490    to that model's mean skill score, then pooling anomalies from all the models and calculating the

491    standard deviation of the result (referred to as the pooled method). The difference between these

492    approaches is minor.

493       There are some similarities between the uncertainty in each metric and the mean value of

494   each metric (Figure 6). For example, the phase of the annual cycle of temperature in the Upper

495   Basin (tas phase) is both the best-simulated and least uncertain metric, while the precipitation

496   variability metrics tend to be the worst simulated and most uncertain. However, there are

497   interesting differences as well. For example, the precipitation variability metrics tend to be the

498   least well simulated (Figure 6), but the uncertainties (Figure S11) are substantially influenced by

499   the averaging period in the metric, with longer averaging periods yielding fewer independent

500   samples and more uncertainty.

501       The spread of values in the overall metric skill score, Dss, for each model with at least 3

502   ensemble members is shown in the Supporting Information, Figure S12. Individual models

503   exhibit a range of spreads, including a standard deviation of 2.67 for CNRM-CM5 (n=5

504   ensemble members), and 0.28 for cccma_cgcm3_1 (n=5). A Monte-Carlo simulation indicates

505   that this nearly order-of-magnitude discrepancy would happen only about 2.5% of the time by

506   chance under the null hypothesis that all the models have the same standard deviation of Dss

507   values.

508       Later figures that display uncertainty in the model's Dss scores are based on this analysis.

509   The uncertainty in models that have less than 3 ensemble members is estimated as the multi-

510   model mean from models with at least three ensemble members. Given that the evidence

511   suggests different models have different levels of variability, this should be considered a rough

512   estimate.

513   *Redundancy in the Metrics*

514       The skill scores presented up to now have been exhaustive, often measuring similar

515   aspects of model performance (for example, the variability of precipitation averaged into 1-, 5-,

516   and 10-year blocks). This redundancy of information can be addressed by forming the EOFs of

517   the skill score matrix (Pierce et al. 2009; Rupp et al. 2013). Computing the EOFs forms optimal

518   combinations of metrics that best describe the model variability, taking covariability between the

519   metrics into account. The number of EOFs to retain is usually chosen so that only modes above

520   the noise floor are kept (Wilks, 2011). Here we choose six modes, which account for 89.1% of

521   the variance.

522    The leading two EOFs (which describe the weighting of each metric) and associated

523    principal components (PCs, which describe the weighting of each model) of the skill score

524    matrix are shown in Figure 7. The EOF weightings show that the leading mode describes poor

525    model performance in simulating seasonal precipitation variability (large negative peaks for DJF,

526    MAM, and JJA precipitation standard deviation), and a large number of models have this

527    problem, led by GFDL-ESM2M and FIO-ESM. Unsurprisingly, these models do not do well in

528    the metric evaluation (Figure 5). The second mode shows co-varying behavior in the quality of

529    model simulations of summer precipitation and temperature variability, and winter precipitation

530    variability. The associated PC shows that 4 models express this behavior strongly: GFDL-

531    ESM2M, GFDL-ESM2G, gfdl_cm2_1, and gfdl_cm2_0. Since this represents two model

532    generations from the same institution, the PC suggests a common physical parameterization or

533    coding approach gives rise to this behavior (c.f. Knutti et al., 2013). However, the CMIP6

534    models from GFDL (GFDL-ESM4 and GFDL-CM4) do not express this relationship, suggesting

535    that a recent change in the model physics or microphysics has altered this behavior.

536    [FIGURE 7 GOES HERE]

537    The overall model rankings (Dss) after the EOF processing are shown in Figure 8. Given

538    the uncertainty in Dss calculated from the ensemble members (indicated by the horizontal red

539    bars), the model with the best overall ranking, EC-EARTH, is not significantly different from

540    any of the other 5 best models. We evaluate the significance of the difference in means using the

541    method of Lanzante 2005, which properly accounts for the joint or pooled uncertainties when

542    estimating the statistical significance of the difference in means of two uncertain quantities. The

543    first model that EC-EARTH is significantly better than is cnrm_cm3, which is rank 6. The curve

544    bends at higher Dss values, indicating that there is a broad and indistinguishable range of

545    relatively good models, but the poorly performing models are more distinct. A version of Figure

546    8 with no culling is shown in Supporting Information Figure S13.

547    [FIGURE 8 GOES HERE]

548    *The Dominance of Precipitation Errors*

549    The results up to now indicate that precipitation errors, especially in winter, are the most

550    problematic aspect of Upper Basin simulations. Since the z-scores normalize by observed natural

551 variability, this is not an artifact of the large variability in this quantity. Given the coarse

552 resolution of most GCMs and the importance of topography to generating precipitation in the

553 Upper Basin, it is not surprising that simulated precipitation in the region often has significant

554 biases. Many applications use bias corrected precipitation in their modeling of the basin to

555 address this. Discarding otherwise high-performing models based on a precipitation bias is

556 questionable when the bias will be removed before the data are used.

557 The following two approaches were explored to address this issue:

558 1) Forming separate precipitation, temperature, and atmospheric circulation indices from

559 the relevant metrics, then weighting those three indices equally to form the final model ranking.

560 This prevents the precipitation biases from dominating the overall model quality ranking, while

561 still allowing absolute rankings amongst the metrics in each class. We term this the "Index-3"

562 approach since it forms an overall index made up of three equally weighted subclasses of indices

563 (temperature, precipitation, and circulation).

564 2) Using a simple bias correction that removes the annual mean bias averaged over the

565 entire Upper Basin region (i.e., a single value) before calculating the metrics. Because only a

566 single annual value is removed for the entire region, this retains the models' simulation of the

567 annual cycle and spatial variability. However, it substantially reduces discrepancies between the

568 model fields and observations.

569 ***Index-3***

570 The Index-3 method forms an index that equally weights temperature, precipitation, and

571 circulation metrics. The temperature class includes all seasonal metrics of Upper Basin mean

572 temperature, the standard deviation of temperature averaged into 1-, 5-, and 10-year blocks, and

573 the amplitude and phase of the annual cycle of temperature. The precipitation class was formed

574 similarly. The circulation class includes all metrics based on ENSO and the PDO, including the

575 teleconnected responses of temperature and precipitation in the Upper Basin region, and the

576 metrics based on the wider-scale correlation maps of surface temperature and sea level pressure

577 with warm and cold season precipitation fluctuations in the Upper Basin. This equal-weighting

578 by class (temperature, precipitation, circulation) gives increased weight to the circulation-based

579 metrics since there are fewer circulation metrics than temperature or precipitation metrics.

580    After each metric was assigned to one of the classes, the values of all metrics that fell into
581    each class were averaged by model. This aggregation yields three quality values per model, one
582    each for temperature, precipitation, and circulation. For each class, the range of values across all
583    models was normalized to the range 0 (best) to 1 (worst), so that the three classes have the same
584    range, in keeping with the purpose of this exercise. The final Index-3 value for a model is the
585    average of the three normalized class values for that model.

586    The model quality scores using Index-3 are shown in the Supporting Information, Figure
587    S14, and show a substantial rearrangement of model rankings (cf. Figure 8). For instance, four
588    CMIP3 (black text) do relatively well in the regular rankings (low Dss), while in the Index-3
589    result the best CMIP3 model appears at rank 17.

590    Some insight into this behavior can be gained by examining changes over model
591    generations in the overall Index-3 and the individual temperature, precipitation, and circulation
592    indices (Figure 9). The difference between the means of pairs of CMIP distributions was
593    evaluated by a two-sample t-test, which indicates that the CMIP generation means are
594    significantly different for the circulation index, but not for the temperature, precipitation, or
595    overall Index-3 indices. In other words, progress across model generations has been dominated
596    by better depictions of large-scale atmospheric circulation, while regional biases (especially in
597    winter precipitation) have not fared as well. Bock et al. (2020) compared a variety of global
598    GCM fields to observations across the CMIP3, 5, and 6 generations, and generally found
599    improvement in the representation of global surface temperature and precipitation fields.
600    However their Figures 3 (temperature) and 4 (precipitation) show that biases across the Western
601    U.S., the focus of interest here but a small part of their global evaluation, show little
602    improvement across model generations. Fasullo (2020) likewise evaluated GCM simulations
603    across CMIP3, 5, and 6, and found that some of the biggest generational improvements were
604    found in aspects of the circulation and ENSO, generally being greater than the improvements
605    found for climatology or on seasonal timescales.

606    [FIGURE 9 GOES HERE]

607 ### *Simple bias correction*

608      The index-3 approach is useful in that it elucidates that the circulation metrics have been

609 the main improvement over model generations, but ultimately most applications bias correct

610 model data before using it. Often, this is done as part of a downscaling process. Accordingly,

611 metrics based on a simple bias correction scheme rather will now be explored, as it better reflects

612 how GCM data is generally used in UCRB studies.

613      Exactly how to evaluate a bias corrected model is still a research question. Common bias

614 correction approaches based on quantile mapping remap the entire model distribution to the

615 observed distribution at every point, which would obviate any comparison with observations in

616 terms of means or variability.

617      The approach taken here is to implement a very simple bias correction rather than a full

618 quantile mapping. The intent is to eliminate the mean biases but still evaluate the model's

619 simulation of the spatial variability and annual cycle of temperature and precipitation in the

620 Upper Basin. To do this, the mean model bias over all spatial points and times is removed, either

621 additively (for temperature) or multiplicatively (for precipitation). Following this simple bias

622 correction, the metrics are recalculated and the result analyzed as shown previously.

623      The portrait plot of metric values after simple bias correction is shown in Figure 10.

624 Comparing to the same result without the simple bias correction (Figure 5), it is clear that when

625 bias correction is added there is a significant overall improvement, as might be expected. A

626 version including the culled models is given in the Supporting Information, Figure S15.

627      Interestingly, some metric values become at least 1 standard deviation worse after simple

628 bias correction (Supporting information, Figure S16). The fact that some metrics degrade might

629 seem counterintuitive, but it happens due to offsetting errors. A model that has lower than

630 observed mean annual precipitation will be bias corrected by multiplying the precipitation fields

631 by a value greater than 1, so that the annual mean matches observations. If that model already

632 has too much precipitation variability, the variability will increase even more, and the skill score

633 will go down as a result.

634      [FIGURE 10 GOES HERE]

635         Figure 11 shows overall model quality rankings after the simple bias correction is

636 applied, with redundant information removed by forming the EOFs as described previously. One

637 striking aspect of Figure 11 is the high performance of the CMIP6 models (red text), which take

638 10 of the top 12 places (along with 2 CMIP5 models). Before the simple bias correction, the top

639 12 places had 3 CMIP3 models, 4 CMIP5 models, and 5 CMIP6 models—a much more equal

640 distribution (Figure 8). Like the Index-3 results, this again illustrates that while biases persist

641 across the model generations, correcting with even a single number (the annual and Upper Basin

642 regional average) reveals that the newer CMIP6 models, as a group, are clearly preferable.

643 Indeed, Figure 11 shows a strong preponderance of CMIP6 models in the top quarter of all

644 models. A similar plot but including all models (no culling) is given in the Supporting

645 Information Figure S17.

646         [FIGURE 11 GOES HERE]

647         The change in model quality from simple bias correction is quantified in Figure 12.

648 Before bias correction a two-sample t-test indicates no significant difference between the CMIP3

649 and CMIP6 means, but after the bias correction the difference is significant at the $p=0.01$ level.

650 Combined with our previous finding that the Index-3 circulation index shows significantly less

651 error in CMIP5 and CMIP6 models than CMIP3 models, this implies that the newer models still

652 struggle with systematic biases, but as a group they do a significantly better job than the older

653 CMIP3 generation in simulating spatial and temporal variability associated with atmospheric

654 circulation patterns.

655         [FIGURE 12 GOES HERE]

656         DISCUSSION

657 *Spatial Resolution in the Depiction of Climate Fields*

658         An important component of the metrics are spatial patterns of mean temperature,

659 precipitation, and variability, and the CMIP5 and CMIP6 models as a group have improved

660 spatial resolution compared with the older CMIP3 generation. This raises the question of whether

661 better results obtained from the CMIP5/6 models are simply due to a more resolved spatial

662 depiction of the fields. We can begin by examining the effect that degrading the spatial

663 resolution of the CMIP5/6 models to match the CMIP3 models has on the model scores. If the

664    primary reason CMIP5/6 models perform better is because they do not smear the spatial fields as

665    much as the lower-resolution CMIP3 models, then degrading the spatial resolution of the CMIP5

666    models might show less difference between the model generations than seen in Figure 12.

667        This is tested in Supporting Information Figure S18, which is the same as Figure 12

668    except that the 1-by-1-degree CMIP5 and CMIP6 data have been aggregated to the 2-by-2-

669    degree grid used by the CMIP3 models. The superiority of the CMIP5 and CMIP6 models

670    remains, and at the same level of significance. The better performance of the CMIP5/6 models in

671    the metrics is not due exclusively to a better resolved depiction of the surface temperature and

672    precipitation fields.

673        Another way to test the effects of model resolution is to stratify model performance by

674    the spatial resolution of the model, as shown in Figure 13. The final model performance is the

675    Dss value from Figure 11, and the spatial resolution is taken as the average of the latitudinal and

676    longitudinal resolutions from Table 1. The relationship between Dss scores and model resolution

677    is shown for each model generation individually (black, blue, and red least-squares best fit trend

678    lines for CMIP3, 5, and 6, respectively), and for all models taken together (purple trend lines).

679    Results both before (left panel) and after (right panel) the simple bias correction are shown. Only

680    one relationship between model quality and spatial resolution is significantly at the 95%

681    confidence interval: the *decrease* in model performance with higher resolution in CMIP3 with no

682    bias correction (left panel, black line). Otherwise, no statistically significant relationships

683    between final model score and the model resolution are found, either when all models are taken

684    together or when each model generation is considered individually, although the trends for the

685    CMIP5 models comes close. We do not argue that model resolution is immaterial to simulations

686    of the UCRB, but these results show that differences in spatial resolution are not the major factor

687    driving differences in performance across GCMs.

688        [FIGURE 13 GOES HERE]

689    *Relation of Model Quality to Projected Climate Changes*

690        It is natural to wonder whether the better-performing models have a systematically

691    different representation of future climate change than the worse-performing models. The

692    regression between model quality and model-projected precipitation trend for the SSP585

693    (CMIP6), RCP 8.5 (CMIP5), and SRES A2 (CMIP3) scenarios is shown in the Supporting

694    Information, Figure S19. No consistent relation between model quality and precipitation change

695    is found. In addition, few individual metrics were found to have any significant relationship with

696    the projected precipitation change. No combination of metrics identified in this way explained

697    more than 20 percent of the variability in projected precipitation change. By contrast, Rupp et al.

698    (2017) found that better performing models showed larger positive winter precipitation

699    projections in the Pacific Northwest. The difference here may be due to the Pacific Norwest

700    falling in the region where GCMs more consistently predict wetter conditions, while the UCRB

701    is close to the zero line where GCMs predict positive precipitation trends to the north and

702    negative trends to the south.

703    *Relation to Global Model Evaluations*

704    Our model evaluation has focused on the UCRB. How do our model rankings compare to

705    published global model rankings? Although a complete evaluation is beyond the scope of this

706    work, some interesting features are evident from the comparison. Many models are reasonably

707    consistent in their ranking across the evaluations, especially considering the uncertainties

708    involved (Figure 11). For example, GFDL-ESM4, FGOALS-f3-L, MPI-ESM1-2-HR, and

709    ACCESS-CM2 score well both here and in Brunner et al. 2020 (their Figure 4). However, there

710    are exceptions. For example, MIROC6 is one of the best models in Fasullo (2021, their Table 1),

711    average in Brunner et al. 2020, and one of the lowest-ranking models in this work. Conversely,

712    IPSL-CM6A-LR does poorly in Fasullo (2021), average in Brunner et al. 2020, and well here.

713    The existence of such models shows that both global and regional metrics should be consulted

714    before selecting a GCM to use in a regional study. Doing well on global metrics is not sufficient

715    to guarantee good performance in a regional setting, and doing well on the regional metrics is not

716    sufficient to guarantee good performance on the global metrics.

717    *Model Genealogy*

718    In this work we have not considered the commonality between models due to shared code

719    or parameterizations (e.g., Knutti et al. 2013, Brunner et al. 2020). However, this may be a

720    consideration when selecting GCMs for applications if a diverse set of models is desired. Our

721    purpose is to evaluate the models that met our data requirements and were available when the

722 analysis was undertaken; if desired, the model rankings develop here can be used to select which
723 one of a family of related set of models is best suited for a user's application.

724 # CONCLUSIONS

725     We have evaluated the ability of CMIP3, CMIP5, and CMIP6 GCMs to reproduce the
726 mean and variability of climate in the Upper Colorado River Basin, including at multi-year time
727 scales for drought applications (5- and 10-year averaging intervals) and teleconnections with the
728 wider hemispheric region. Using a set of 48 metrics, we have ranked 62 GCMs by overall quality
729 of their simulations of seasonal and annual temperature and precipitation, for both the mean and
730 variability in the region. The ranking included an initial culling of the GCMs, with 25% of each
731 generation of models discarded based on poor performance on global metrics. This reduces the
732 chance that a model does well in the limited region of the Upper Colorado River Basin for the
733 wrong dynamical reasons.

734     A key aspect of our approach is to evaluate the models after a simple bias correction has
735 been applied. This is motivated by the fact that stakeholders and impact studies in the region
736 generally use bias-corrected fields. However additional information was obtained from the
737 original (non-bias corrected) GCM fields, with the main finding that the CMIP3 models do
738 systematically worse than the CMIP5 and CMIP6 models on the metrics relating to global
739 atmospheric circulation. The CMIP6 models also do significantly better than the CMIP3 models
740 after the simple bias correction is applied. Nonetheless, it is worth noting that even after the
741 simple bias correction, the GCMs show appreciable residual biases in their depiction of the
742 climate in the Upper Basin, particularly in the interannual variability of winter precipitation.
743 Although GCMs are currently the best tools available for projecting future climate change over
744 broad regions, they have problems simulating relatively small regions with significant
745 topography, such as the Upper Basin. Our results show that these biases have persisted across
746 model generations, even as performance on metrics of atmospheric circulation has improved.

747 # SUPPORTING INFORMATION

748     Additional supporting information may be found online under the Supporting Information
749 tab for this article: Additional figures and illustrations.

750 # DATA AVAILABILITY

751      Metric values and model quality results are available for download at

752 http://cirrus.ucsd.edu/~pierce/pierce_et_al_2021_UCRB_GCM_selection. This will allow

753 individual practitioners to weight individual metrics or model results as needed for their

754 applications.

## LITERATURE CITED

764      Abatzoglou, J. T., and D. E. Rupp. 2017. "Evaluating climate model simulations of

765 drought for the northwestern United States." *Internat. J. of Climatol.* 37: 910-920.

766      Barnett, T., R. Malone, W. Pennell, D. Stammer, B. Semtner, and W. Washington. 2004.

767 "The Effects of Climate Change on Water Resources in the West: Introduction and Overview."

768 *Climatic Change* 62: 1-11.

769      Barnett, T. P., and D. W. Pierce. 2008. "Sustainable Water Deliveries from the Colorado

770 River in a Changing Climate." *Proceedings of the National Academy of Sciences* (PNAS) 106

771 (18): 7334-7338. Doi/10.1073/pnas.0812762106

772      Bennett, K., E., J. R. Urrego Blanco, A. Jonko, T. J. Bohn, A. L. Atchley, N. M. Urban,

773 and R. S. Middleton. 2018. "Global Sensitivity of Simulated Water Balance Indicators Under

774 Future Climate Change in the Colorado Basin." *Water Resources Research* (WRR) 54: 132-149.

775 https://doi.org/10.1002/2017WR020471

776      Bock, L, A. Lauer, M. Schlund, M. Barreiro, N. Bellouin, C. Jones, G. A. Meehl, V.

777 Predoi, M. J. Roberts, and V. Eyring. 2020. "Quantifying Progress Across Different CMIP

778    Phases with the ESMValTool." *J. Geophys. Res. Atmos*. 125, 28 pp.,

779    https://doi.org/10.1029/2019JD032321

780        Boone, A., P. de Rosnay, G. Balsamo, A. Beljaars, F. Chopin, et al. 2009. "The AMMA

781    Land Surface Model intercomparison project (ALMIP)." *Bull. Am. Met. Soc.* 90 (12): 1865-1880.

782    Doi: 10.1175/2009BAMS2786.1

783        Brunner, L., R. Lorenz, M. Zumwald, and R. Knutti. 2019. "Quantifying uncertainty in

784    European climate projections using combined performance-independence weighting." *Env. Res.*

785    *Lett*. (14): 124010. Doi: https://doi.org/10.1088/1748-9326/ab492f

786        Brunner, L., A. G. Pendergrass, F. Lehner, A. L. Merrifield, R. Lorenz, and R. Knutti.

787    2020. "Reduced Global Warming from CMIP6 Projections when Weighting Models by

788    Performance and Independence." *Earth Systems Dynamics* 11: 995-1012.

789    https://doi.org/10.5194/esd-11-995-2020

790        CA-DWR, California Department of Resources Climate Change Technical Advisory

791    Group. 2015. "Perspectives and Guidance for Climate Change Analysis."

792    https://www.water.ca.gov/LegacyFiles/climatechange/docs/2015/Perspectives_Guidance_Climat

793    e_Change_Analysis.pdf.

794        Cayan, D. R., T. Das, D. W. Pierce, T. P. Barnett, M. Tyree, and A. Gershunov. 2010.

795    "Future Dryness in the Southwest US and the Hydrology of the Early 21$^{st}$ Century Drought."

796    *Proceedings of the National Academy of Sciences* (PNAS) 107: 21271-21276.

797    Doi/10.1073/pnas.0912391107

798        Christensen, N. S., A. W. Wood, N. Voisin, D. P. Lettenmaier, and R. N. Palmer. 2004.

799    "The Effects of Climate Change on the Hydrology and Water Resources of the Colorado River

800    Basin." *Climatic Change* 62: 337-363.

801        Christensen, N. S., and D. P. Lettenmaier. 2007. "A Multimodel Ensemble Approach to

802    Assessment of Climate Change Impacts on the Hydrology and Water Resources of the Colorado

803    River Basin." *Hydrology and Earth System Sciences* (HESS) 11: 1417-1434.

804         Dawadi, S., and S. Ahmad. 2012. "Changing Climatic Conditions in the Colorado River

805 Basin: Implications for Water Resources Management." *J. Hydrology* 430-431: 127-141.

806 Doi:10.1016/j.jhydrol.2012.02.010

807         Dettinger, M., B. Udall, and A. Georgakakos. 2015. "Western Water and Climate

808 Change." *Ecological Applications* 25 (8): 2069-2093.

809         Dirmeyer, P. A., Y Jin, B. Singh, and XQ Yan. 2013. "Trends in land-atmosphere

810 interactions from CMIP5 simulations." J. *Hydromet.* 14(3): 829-849. Doi: 10.1175/JHM-D-12-

811 0107.1

812         Eyring, V., S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer, and K. E.

813 Taylor. 2016. "Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6)

814 Experimental Design and Organization." *Geosci. Model Dev.,* 9 (5): 1937-1958. DOI:

815 10.5194/gmd-9-1937-2016.

816         Fasullo, J. T. 2020. "Evaluating simulated climate patterns from the CMIP archives using

817 satellite and reanalysis datasets using the Climate Model Assessment Tool (CMATv1)." *Geosci.*

818 *Model Dev.* 13: 3627. https://doi.org/10.5194/gmd-13-3627-2020.

819         Ficklin, D. L., S. L. Letsinger, I. T. Stewart, and E. P. Maurer. 2016. "Assessing

820 Differences in Snowmelt-Dependent Hydrologic Projections using CMIP3 and CMIP5 Climate

821 Forcing Data for the Western United States." *Hydrology Research* 47 (2) 483-500. Doi:

822 10.2166/nh.2015.101

823         Ficklin, D. L., I. T. Stewart, and E. P. Maurer. 2013. "Climate Change Impacts on

824 Streamflow and Subbasin-Scale Hydrology in the Upper Colorado River Basin." *PLOS One* 8

825 (8): e71297. Doi:10.137/journal.pone.0071297

826         Flato, G., J. Marotzke, B. Abiodun, P. Braconnot, S.C. Chou, W. Collins, P. Cox, F.

827 Driouech, S. Emori, V. Eyring, C. Forest, P. Gleckler, E. Guilyardi, C. Jakob, V. Kattsov, C.

828 Reason, and M. Rummukainen. 2013. "Evaluation of climate models." In: *Climate Change*

829 *2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment*

830 *Report of the Intergovernmental Panel on Climate Change.* Edited by T.F. Stocker, D. Qin, G.-

831  K. Plattner, M. Tignor, S.K. Allen, J. Doschung, A. Nauels, Y. Xia, V. Bex, and P.M. Midgley,

832  Cambridge University Press, pp. 741-882, doi:10.1017/CBO9781107415324.020.

833  Gleckler, P. J., K. E. Taylor, and C. Doutriaux. 2008. "Performance metrics for climate

834  models." *Journal of Geophysical Research Atmospheres* 113: D06104. 10.1029/2007JD008972

835  Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J.,

836  Nicolas, J., Peubey, C., Radu, R., Rozum, I., Schepers, D., Simmons, A., Soci, C., Dee, D.,

837  Thépaut, J-N. 2018. "ERA5 hourly data on single levels from 1979 to present." *Copernicus*

838  *Climate Change Service (C3S) Climate Data Store (CDS)*. (Accessed on 22-JAN-2021),

839  10.24381/cds.adbb2d47.

840  Hidalgo, H. G., and J. A. Dracup, 2003. "ENSO and PDO Effects on Hydroclimatic

841  Variations of the Upper Colorado River Basin." *Journal of Hydrometeorlogy* 4: 5-23.

842  Hoerling, M., J. Barsugli, B. Livneh, J. Eischeid, X. Quan, and A. Badger. 2019. "Causes

843  for the Century-Long Decline in Colorado River Flow." *Journal of Climate* 32: 8181-8203.

844  Doi:10.1175/JCLI-D-19-0207.1

845  IPCC, 2013. "*Climate Change 2013: The Physical Science Basis. Contribution of*

846  *Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate*

847  *Change*." Cambridge: Cambridge University Press.

848  Jenkins, G. M., and Watts, D. G. 1968. *Spectral Analysis and its Applications*. Emerson-

849  Adams Press, Inc.

850  Jong, B-T., J. Ting, and R. Seager. 2021. "Assessing ENSO Summer Teleconnections,

851  Impacts, and Predictability in North America." *J. Climate* 34: 3629-3643.

852  https://doi.org/10.1175/JCLI-D-20-0761.1

853  Kalnay, E., M. Kanamitsu, R. Kistler, W. Collins, D. Deaven, et al. 1996. "The

854  NCEP/NCAR 40-year Reanalysis project." *Bulletin of the American Meteorological Society*

855  (BAMS) 77: 437-471.

856  Knutti, R., D. Masson, and A. Gettelman. 2013. "Climate model genealogy: Generation

857  CMIP5 and how we got there." *Geophysical Research Letters* 40: 1194-99.

858  doi:10.1002/grl.50256

859        Knutti, R., J. Sedláček, B. M. Sanderson, R. Lorenz, E. M. Fischer, and V. Eyring, 2017:
860    "A climate model projection weighting scheme accounting for performance and
861    interdependence." *Geophys. Res. Lett*. 44(4): 1909-1918. https://doi.org/10.1002/2016GL072012

862        Lanzante, J. R. 2005. "A cautionary note on the use of error bars." Journal of Climate 18:
863    3699-3703. https://doi.org/10.1175/JCLI3499.1

864        Li, JD., CY Miao, W. Wei, G. Zhang, LJ Hua, YL Chen, and XX Wang. 2021.
865    "Evaluation of CMIP6 Global Climate Models for Simulating Land Surface Energy and Water
866    Fluxes During 1979-2014." J. *Advances Modeling Earth Sys*. 13(6): art. e2021MS002515. Doi
867    10.1029/2021MS002515

868        Li, L, Y. Wang, V. K. Arora, D. Eamus, H. Shi, J. Li, L. Cheng et al. 2018. "Evaluating
869    Global Land Surface Models in CMIP5: Analysis of Ecosystem Water- and Light-Use
870    Efficiencies and Rainfall Partitioning." *J. Clim*. 31(8): 2995-3008. https://doi.org/10.1175/JCLI-
871    D-16-0177.1

872        Livneh, B., T. J. Bohn, D. W. Pierce, F. Munoz-Arriola, B. Nijssen, R. Vose, D. R.
873    Cayan, and L. Brekke. 2015. "A spatially comprehensive, hydrometeorological data set for
874    Mexico, the U.S., and Southern Canada 1950-2013." *Scientific Data* 2: article 150042 (2015).
875    doi:10.1038/sdata.2015.42.

876        Lorenz, R., N. Herger, J. Sedláček, V. Eyring, E. M. Fischer, and R. Knutti. 2018.
877    "Prospects and caveats of weighting climate models for summer maximum temperature
878    projections over North America." *Journal of Geophysical Research: Atmospheres* 123: 4509–
879    4526. Doi: https://doi.org/10.1029/2017JD027992

880        Massmann, C. 2020. "Evaluating the Suitability of Century-Long Gridded
881    Meteorological Datasets for Hydrological Modeling." *J. Hydrometeorology* 21(11): 2565-2580.
882    Doi: https://doi.org/10.1175/JHM-D-19-0113.1

883        McCabe, G. J, D. M. Wolock, G. T. Pederson, C. A. Woodhouse, and S. Mcafee. 2017.
884    "Evidence that Recent Warming is Reducing Upper Colorado River Flows." *Earth Interactions*
885    (21) 1-14.

886     McCabe, G. J., and D. M. Wolock. 2020. "The Water-Year Balance of the Colorado
887 River Basin." *Journal of the American Water Resources Association* (JAWRA) 56 (4): 724-737.
888 https://doi.org/10.1111/1752-1688.12848

889     Meehl, G. A., C. Covey, T. Delworth, M. Latif, B. McAvaney, J. F. B. Mitchell, R. J.
890 Stouffer, and K. E. Taylor. 2007. "The WCRP CMIP3 Multimodel Dataset: A New Era in
891 Climate Change Research." *Bulletin of the American Meteorological Society* (BAMS) Sept.
892 2007: 1383-1394. https://doi.org /10.1175/BAMS-88-9-1383

893     Moon, H., Gudmundsson, L., & Seneviratne, S. I., 2018. "Drought persistence errors in
894 global climate models." *J. Geophys. Res. Atmospheres* 123(7): 3483-3496.

895     Pierce, D. W., T. P. Barnett, B. D. Santer, and P. J. Gleckler. 2009. "Selecting global
896 climate models for regional climate change studies." *Proceedings of the National Academy of
897 Sciences* 106 (21): 8441-8446. doi:10.1073/pnas.0900094106

898     Pierce, D. W., L. Su, D. R. Cayan, M. D. Risser, B. Livneh, and D. P. Lettenmaier. 2021.
899 "An Extreme-Preserving Long-Term Gridded Daily Precipitation Dataset for the Conterminous
900 United States." *Journal of Hydrometeorology* 22(7): 1883-1895. DOI:
901 https://doi.org/10.1175/JHM-D-20-0212.1

902     Rajagopalan, B., K. Nowak, J. Prairie, M. Hoerling, B. Harding, J. Barsugli, A. Ray, and
903 B. Udall. 2009. "Water Supply Risk on the Colorado River: Can Management Mitigate?" *Water
904 Resources Research* (WRR) 45: W08201. Doi:10.1029/2008WR007652

905     Riahi, K., D. P. van Vuuren, E. Kriegler, J. Edmonds, B. C. O'Neill, et al. 2017. "The
906 Shared Socioeconomic Pathways and their Energy, Land Use, and Greenhouse Gas Emissions
907 Implications: An Overview." *Global Environmental Change* 42: 153-168. DOI:
908 10.1016/j.gloenvcha.2016.05.009

909     Rupp, D. W., J. T. Abatzoglou, K. C. Hegewisch, and P. W. Mote. 2013. "Evaluation of
910 CMIP5 20th century simulations for the Pacific Northwest USA." *Journal of Geophysical
911 Research Atmospheres* 118: 10,884-10,906. Doi:10.1002/jgrd.50843

912     Rupp, D. E., J. T. Abatzoglou, and P. W. Mote. 2017. "Projections of 21st century
913 climate of the Columbia River Basin." *Climate Dynamics* (49), 1783-1799.

914	Sanderson, B. M. and M. F. Wehner. 2017. "Model weighting strategy." In: *Climate*
915	*Science Special Report: Fourth National Climate Assessment, Volume I*, edited by D. J.
916	Wuebbles, D.W. Fahey, K.A. Hibbard, D.J. Dokken, B.C. Stewart, and T.K. Maycock, 436-442.
917	Washington D. C.: U.S. Global Change Research Program. doi: 10.7930/J06T0JS3.

918	Seager, R., M. Ting, C. Li, N. Naik, B. Cook, J. Nakamura, and H. Liu. 2012.
919	"Projections of Declining Surface-Water Availability for the Southwestern United States."
920	*Nature Climate Change* 3 (5) 482-486. Doi:10.1038/NCLIMATE1787

921	Tamaddun, K. A., A. Kalra, S. Kumar, and S. Ahmad. 2019. "CMIP5 Models' Ability to
922	Capture Observed Trends under the Influence of Shifts and Persistence: An In-Depth Study on
923	the Colorado River Basin." *Journal of Applied Meteorology and Climatology* (56) 1677-1688.
924	Doi: 10.1175/JAMC-D-18-0251.1

925	Taylor, K. E., R. J. Stouffer, and G. A. Meehl. 2012. "An Overview of CMIP5 and the
926	Experiment Design." *Bulletin of the American Meteorological Society* (BAMS) Apr. 2012: 485-
927	498. DOI: 10.1175/BAMS-D-11-00094.1

928	Tillman, F. D., S. Gangopadhyay, and T. Pruitt. "Understanding the Past to Interpret the
929	Future: Comparison of Simulated Groundwater Recharge in the Upper Colorado River Basin
930	(USA) Using Observed and General-Circulation-Model Historical Climate Data."
931	*Hydrogeological Journal* 25: 347-358. Doi: 10.1007/s100040-016-1481-0

932	Tokarska, K. B., M. B. Stolpe, S. Sippel, E. M. Fischer, C. J. Smith, F. Lehner, and R.
933	Knutti. 2020. "Past Warming Trend Constrains Future Warming in CMIP6 Models." *Science*
934	*Advances* 6 (12): eaaz9549. DOI: 10.1126/sciadv.aaz9549

935	Udall, B., and J. Overpeck. 2017. "The Twenty-First Century Colorado River Hot
936	Drought and Implications for the Future." *Water Resources Research* (WRR) 53: 2404-2418.
937	Doi:10.1002/2016WR019638

938	Ukkola, A. M., Pitman, A. J., De Kauwe, M. G., Abramowitz, G., Herger, N., Evans, J.
939	P., and Decker, M. (2018). "Evaluating CMIP5 model agreement for multiple drought metrics."
940	*Journal of Hydrometeorology* 19(6): 969-988.

941        USBR, 2012. "Colorado River Basin Water supply and demand study: Executive

942    summary." U.S. Dept. Interior, Bureau of Reclamation.

943    https://www.usbr.gov/watersmart/bsp/docs/finalreport/ColoradoRiver/CRBS_Executive_Summa

944    ry_FINAL.pdf (accessed 2021-02-22).

945        Vano, J. A., T. Das, and D. P. Lettenmaier. 2012. "Hydrologic Sensitivities of Colorado

946    River Runoff to Changes in Precipitation and Temperature." *Journal of Hydrometeorology* 13:

947    932-949. Doi: 10.1175/JHM-D-11-069.1

948        Vano, J. A., B. Udall, D. R. Cayan, J. T. Overpeck, et al. 2014. "Understanding

949    Uncertainties in Future Colorado River Streamflow." *Bulletin of the American Meteorological*

950    *Society* (BAMS) Jan. 2014: 59-78. Doi: 10.1175/BAMS-D-12-00228.1

951        Wilks, D. S. 2011. *Statistical Methods in the Atmospheric Sciences, Third Edition*.

952    Oxford: Academic Press.

953        Xiao, M., B. Udall, D. P. Lettenmaier. 2018. "On the Causes of Declining Colorado

954    River Streamflows." *Water Resources Research* (WRR) 54: 6739-6756.

955    https://doi.org/10.1029/2018WR023153

956    Table 1. Models used in this analysis, whether they belong to the CMIP-3, 5, or 6 generation, the

957    number of historical ensemble members analyzed, and the approximate resolution of the

958    atmospheric model, in degrees (longitude x latitude).

| Model | CMIP generation | Num. ensemble members | Approx. atmo resolution (deg) |
|---|---|---|---|
| 1. ACCESS-CM2 | 6 | 3 | 1.875 x 1.25 |
| 2. ACCESS-ESM1-5 | 6 | 5 | 1.875 x 1.25 |
| 3. AWI-CM-1-1-MR | 6 | 5 | 0.938 x 0.935 |
| 4. BCC-CSM2-MR | 6 | 1 | 1.125 x 1.121 |
| 5. BCC-ESM1 | 6 | 3 | 2.8 x 2.8 |
| 6. CNRM-CM6-1 | 6 | 1 | 1.4 x 1.4 |
| 7. CNRM-CM6-1-HR | 6 | 1 | 0.5 x 0.5 |
| 8. CNRM-ESM2-1 | 6 | 5 | 1.4 x 1.4 |
| 9. CanESM5 | 6 | 7 | 2.8 x 2.8 |
| 10. EC-Earth3 | 6 | 4 | 0.7 x 0.7 |
| 11. EC-Earth3-Veg | 6 | 5 | 0.7 x 0.7 |
| 12. FGOALS-f3-L | 6 | 3 | 1.25 x 1.0 |
| 13. FGOALS-g3 | 6 | 4 | 2.0 x 2.278 |
| 14. GFDL-CM4 | 6 | 1 | 1.25 x 1.0 |
| 15. GFDL-ESM4 | 6 | 1 | 1.25 x 1.0 |
| 16. HadGEM3-GC31-LL | 6 | 4 | 1.875 x 1.25 |
| 17. IPSL-CM6A-LR | 6 | 21 | 2.5 x 1.27 |
| 18. INM-CM4-8 | 6 | 1 | 2.0 x 1.5 |
| 19. INM-CM5-0 | 6 | 10 | 2.0 x 1.5 |
| 20. MIROC6 | 6 | 7 | 1.4 x 1.4 |
| 21. KACE-1-0-G | 6 | 3 | 1.875 x 1.25 |
| 22. MPI-ESM1-2-HR | 6 | 10 | 0.938 x 0.935 |
| 23. MPI-ESM-1-2-HAM | 6 | 2 | 1.875 x 1.865 |

| Model | CMIP generation | Num. ensemble members | Approx. atmo resolution (deg) |
|---|---|---|---|
| 24. MPI-ESM1-2-LR | 6 | 9 | 1.875 x 1.865 |
| 25. MRI-ESM2-0 | 6 | 6 | 1.125 x 1.121 |
| 26. NESM3 | 6 | 3 | 1.875 x 1.865 |
| 27. NorCPM1 | 6 | 30 | 2.5 x 1.9 |
| 28. NorESM2-LM | 6 | 3 | 2.5 x 1.9 |
| 29. NorESM2-MM | 6 | 3 | 1.25 x 0.942 |
| 30. TaiESM1 | 6 | 1 | 1.25 x 0.942 |
| 31. UKESM1-0-LL | 6 | 5 | 1.875 x 1.25 |
| 32. ACCESS1-0 | 5 | 1 | 1.875 x 1.25 |
| 33. ACCESS1-3 | 5 | 1 | 1.875 x 1.25 |
| 34. CanESM2 | 5 | 5 | 2.8 x 2.8 |
| 35. CCSM4 | 5 | 5 | 1.25 x 0.942 |
| 36. CESM1-BGC | 5 | 1 | 1.25 x 0.942 |
| 37. CESM1-CAM5 | 5 | 3 | 1.25 x 0.942 |
| 38. CMCC-CM | 5 | 1 | 0.75 x 0.75 |
| 39. CNRM-CM5 | 5 | 5 | 1.4 x 1.4 |
| 40. CSIRO-Mk3-6-0 | 5 | 10 | 1.875 x 1.865 |
| 41. EC-EARTH | 5 | 4 | 1.125 x 1.12 |
| 42. FGOALS-g2 | 5 | 1 | 2.8 x 3.0 |
| 43. FGOALS-s2 | 5 | 2 | 2.812 x 1.659 |
| 44. FIO-ESM | 5 | 3 | 2.8 x 2.8 |
| 45. GFDL-CM3 | 5 | 1 | 2.5 x 2.0 |
| 46. GFDL-ESM2G | 5 | 1 | 2.5 x 2.0 |
| 47. GFDL-ESM2M | 5 | 1 | 2.5 x 2.0 |
| 48. GISS-E2-R | 5 | 5 | 2.5 x 2.0 |
| 49. GISS-E2-R-CC | 5 | 1 | 2.5 x 2.0 |

| Model | CMIP generation | Num. ensemble members | Approx. atmo resolution (deg) |
|---|---|---|---|
| 50. HadGEM2-AO | 5 | 1 | 1.875 x 1.25 |
| 51. HadGEM2-CC | 5 | 1 | 1.875 x 1.25 |
| 52. HadGEM2-ES | 5 | 4 | 1.875 x 1.25 |
| 53. IPSL-CM5A-LR | 5 | 4 | 3.75 x 1.9 |
| 54. IPSL-CM5A-MR | 5 | 1 | 2.5 x 1.268 |
| 55. IPSL-CM5B-LR | 5 | 1 | 3.75 x 1.9 |
| 56. MIROC-ESM | 5 | 1 | 2.8 x 2.8 |
| 57. MIROC-ESM-CHEM | 5 | 1 | 2.8 x 2.8 |
| 58. MIROC5 | 5 | 1 | 1.4 x 1.4 |
| 59. MPI-ESM-LR | 5 | 3 | 1.875 x 1.865 |
| 60. MPI-ESM-MR | 5 | 1 | 1.875 x 1.865 |
| 61. MRI-CGCM3 | 5 | 1 | 1.125 x 1.121 |
| 62. NorESM1-M | 5 | 1 | 2.5 x 1.9 |
| 63. NorESM1-ME | 5 | 1 | 2.5 x 1.9 |
| 64. bcc-csm1-1 | 5 | 1 | 2.8 x 2.8 |
| 65. bcc-csm1-1-m | 5 | 1 | 1.125 x 1.121 |
| 66. inmcm4 | 5 | 1 | 2.0 x 1.5 |
| 67. bccr_bcm2_0 | 3 | 1 | 2.8 x 2.8 |
| 68. cccma_cgcm3_1 | 3 | 5 | 3.75 x 3.7 |
| 69. cnrm_cm3 | 3 | 1 | 2.8 x 2.8 |
| 70. csiro_mk3_0 | 3 | 1 | 1.875 x 1.865 |
| 71. gfdl_cm2_0 | 3 | 1 | 2.5 x 2.0 |
| 72. gfdl_cm2_1 | 3 | 1 | 2.5 x 2.0 |
| 73. giss_model_e_r | 3 | 1 | 5.0 x 3.95 |
| 74. inmcm3_0 | 3 | 1 | 5.0 x 4.0 |
| 75. ipsl_cm4 | 3 | 1 | 3.75 x 2.5 |

| Model | CMIP generation | Num. ensemble members | Approx. atmo resolution (deg) |
|---|---|---|---|
| 76. miroc3_2_medres | 3 | 3 | 2.8 x 2.8 |
| 77. miub_echo_g | 3 | 3 | 3.75 x 3.7 |
| 78. mpi_echam5 | 3 | 3 | 1.875 x 1.865 |
| 79. mri_cgcm2_3_2a | 3 | 5 | 2.8 x 2.8 |
| 80. ncar_ccsm3_0 | 3 | 7 | 1.4 x 1.4 |
| 81. ncar_pcm1 | 3 | 4 | 2.8 x 2.8 |
| 82. ukmo_hadcm3 | 3 | 1 | 3.75 x 2.5 |

959

**Figure Captions**

961	Figure 1. Histograms of model-projected surface temperature changes (°C; tas; red bars/left

962	column) and precipitation trends (mm/day per century; pr; green bars/right column) in the Upper

963	Colorado River Basin. Top row: for all models before the global culling. For models that have

964	multiple ensemble members, values are averaged across ensemble members before plotting.

965	Bottom row: after the global culling. Changes are for the sresa2, RCP 8.5, and SSP85 scenarios

966	for the CMIP3, CMIP5, and CMIP6 models, respectively. Temperature changes are evaluated as

967	the 2070-2099 mean minus the 1950-2005 mean. Precipitation changes are evaluated as a best-fit

968	least squares linear trend (millimeters per day per century) over the period 1950-2099. For

969	reference, the mean observed annual precipitation over the Upper Basin is approximately 1.09

970	millimeters per day.

971	Figure 2. The Upper Colorado River Basin (brown outline) and the evaluation domain used in

972	this work, as indicated by the centers of the 1-by-1 gridcells (black dots). Colors show elevation

973	in meters.

974	Figure 3. Example showing the fields for the calculation of winter (DJF) precipitation metric.

975	Top left: observed field (mm day$^{-1}$). Bottom left: Standard deviation of observations (mm day$^{-1}$).

976	Top center: A model that performs well, CanESM2. Bottom middle: the z-score for CanESM2,

977	i.e., the difference between the model and observations, divided by the observed standard

978	deviation. Right column: Same as the middle column, but for FIO-ESM, which performs poorly

979	on this metric.

980	Figure 4. Top left: the observed SST pattern (°C) associated with ENSO. Top right: the standard

981	deviation (°C) of the observed pattern. Middle row: for CESM1-CAM5, the model's observed

982	pattern of SST for ENSO (left) and the model's z-score (right, dimensionless). Bottom row:

983	same, for CSIRO-Mk3-6-0.

984	Figure 5. Portrait plot of the model skill scores. The metrics are along the X axis (orange/red

985	shows good skill, blues show poor skill), and the models along the Y axis. Metric labels, from

986	left to right, are: seasonal (DJF, MAM, JJA, SON) mean ($\bar{x}$) and standard deviation in 1-, 5-, and

987	10-year blocks ($\sigma_1$, $\sigma_5$, $\sigma_{10}$) of temperature (T) and precipitation (P); the seasonal cycle

988    evaluated via the amplitude (A) and phase ($\phi$) of temperature and precipitation; ENSO and the

989    PDO evaluated via the mean SST pattern in the tropical Pacific ($\bar{x}$), the spectrum (S), and the

990    teleconnected response over the western U.S. in temperature (T) and precipitation (P); and the

991    teleconnected (TCON) correlation maps of temperature (T) and sea level pressure (S) with Upper

992    Basin precipitation variability during the warm (Wrm; AMJJAS) and cold (Cld; ONDJFM)

993    seasons. Names of CMIP3 models are shown in black, CMIP5 models in blue, and CMIP6 in

994    red.

995    Figure 6. Distribution of skill scores sorted by mean metric value. Better simulated metrics have

996    higher skill scores (closer to the perfect value of 1) and are plotted at the top. Worse simulated

997    metrics are at the bottom. The whiskers and dots show the mean of the metric (center line),

998    interquartile range (box), 90 percent range (bars), and extreme values (dots).

999    Figure 7. The leading two EOFs (red) and associated PCs (blue) of the skill score matrix.

1000    CMIP3, 5, and 6 model names are in black, blue, and red, respectively. The first two EOFs

1001    explain 68.1% and 9.4% of the variance, respectively.

1002    Figure 8. Model quality rankings after the EOF process has been applied. Best models (lowest

1003    Dss values) are at the bottom, worst models at the top. The 95% confidence intervals (red lines)

1004    are estimates derived from an analysis of models with multiple ensemble members – see text for

1005    details. The number of realizations is shown by n along the right hand side. Black, blue, and red

1006    names indicate CMIP3, CMIP5, and CMIP6, respectively. The vertical black bars with diamonds

1007    illustrate, for a few example models chosen to span the results, the range of models whose

1008    rankings are statistically indistinguishable from the base model (indicated by the diamond) given

1009    the uncertainty in Dss. For example, the Dss value of TaiESM1 is not statistically distinct from

1010    the Dss values of models ranging from GFDL-ESM4 to GISS-E2-R-CC.

1011    Figure 9. Performance of the CMIP3, CMIP5, and CMIP6 (black, blue, and red, respectively)

1012    models on Index-3 (top left), and the individual components of index-3, the temperature (top

1013    right), precipitation (bottom left), and circulation (bottom right) indices. The diamonds show the

1014    mean of each CMIP's distribution; the dots show individual model values. The P value shown in

1015    the panels is the chance that the CMIP distributions have means that differ only due to sampling

1016    fluctuations.

1017    Figure 10. Portrait plot of the model skill scores for the case with simple bias correction. The
1018    format is the same as Figure 8; see that caption for figure details.

1019    Figure 11. Model quality rankings after the simple bias correction has been applied and
1020    redundant information removed via an EOF approach. The format is the same as Figure 8.

1021    Figure 12. Model errors (lower is better) for the CMIP3, CMIP5, and CMIP6 model generations,
1022    both before (light dots) and after (dark dots) the simple bias correction. The diamonds indicate
1023    the mean of the distributions. The P values shown in the lower left are the chance that the means
1024    of the CMIP3 and CMIP6 distributions differ only due to sampling fluctuations, as estimated by
1025    a two-sample t-test.

1026    Figure 13. Final model performance (Dss from Figure 11) as a function of model spatial
1027    resolution (column 4 of Table 1). Left: no bias correction. Right: With simple bias correction.
1028    The spatial resolution is taken as the average of the longitudinal and latitudinal resolutions, in
1029    degrees. Solid purple line: least-squares best-fit line using all models. Black, blue, and red lines:
1030    least-squares best fit lines using only CMIP 3, 5, and 6 models, respectively. None of the trends
1031    are significantly different from zero at the 95% confidence level except the CMIP3 trend in the
1032    No-BC case.

**tas change: CMIP3+5+6 ALL models**

mean: 5.97
stddev: 1.44

Count of models

proj change, degC fut−hist

**pr change: CMIP3+5+6 ALL models**

mean: 0.07
stddev: 0.15

Count of models

proj change, (mm/day)/century

**tas change: CMIP3+5+6 RETAINED models**

mean: 5.73
stddev: 1.03

Count of models

proj change, degC fut−hist

**pr change: CMIP3+5+6 RETAINED models**

mean: 0.05
stddev: 0.15

Count of models

proj change, (mm/day)/century

jawr_12974_f1.eps

jawr_12974_f2.eps

jawr_12974_f3.eps

**Observed ENSO sst anom [C]**

**Observed std dev [C]**

**CESM1−CAM5 ENSO sst anom [C]**

**CESM1−CAM5 z−score**

**CSIRO−Mk3−6−0 ENSO sst anom [C]**

**CSIRO−Mk3−6−0 z−score**

jawr_12974_f4.eps

Labels (top to bottom):
tas phase
tas MAM sd−10yr
PDO spectrum
tas SON sd−5yr
pr SON mean
tas SON sd−10yr
tas DJF sd−10yr
tas MAM sd−5yr
UCRB pr w/tas, warm
ENSO tas
ENSO spectrum
tas DJF sd−5yr
UCRB pr w/psl, warm
tas DJF mean
PDO tas
pr SON sd−10yr
pr phase
pr JJA mean
UCRB pr w/tas, cold
tas amp
ENSO pattern
tas JJA sd−10yr
pr SON sd−5yr
tas JJA sd−5yr
pr MAM sd−10yr
UCRB pr w/psl, cold
pr JJA sd−10yr
pr amp
tas SON sd−1yr
PDO pr
ENSO pr
PDO pattern
pr JJA sd−5yr
pr DJF sd−10yr
pr MAM mean
pr MAM sd−5yr
tas MAM sd−1yr
tas SON mean
tas MAM mean
tas DJF sd−1yr
pr DJF mean
pr DJF sd−5yr
pr SON sd−1yr
tas JJA sd−1yr
tas JJA mean
pr MAM sd−1yr
pr JJA sd−1yr
pr DJF sd−1yr

jawr_12974_f7.eps

jawr_12974_f8.eps

jawr_12974_f9.eps

jawr_12974_f11.eps

jawr_12974_f12.eps

jawr_12974_f13.eps