1 **An objective statistical downscaling technique for emulating WRF**

2

3 George Shu Heng Pau[1], Daniel B. Walton[*2], Samir Touzani[3]

4

5 [1]Lawrence Berkeley National Laboratory

6 Climate and Ecosystem Sciences Division

7 1 Cyclotron Road

8 Berkeley, CA 94720.

9 [2]University of California, Los Angeles

10 Institute of the Environment and Sustainability

11 La Kretz Hall, Suite 300

12 Los Angeles, CA 90095-1496.

13 [3]Lawrence Berkeley National Laboratory

14 Building Technology and Urban Systems Division

15 1 Cyclotron Road

16 Berkeley, CA 94720.

17

18 *Corresponding author email address: waltond@ucla.edu

19

20

21 **Abstract**

22

23 Accurate downscaling of global climate models (GCMs) is needed to quantify the local

24 impacts of climate change. Dynamical downscaling with a regional climate model has

25 been shown to capture important physical processes at fine scales, but it is too

26 computationally expensive to be used to downscale a large ensemble of GCMs or

27 multiple time periods and scenarios. Hybrid dynamical-statistical downscaling saves time

28 by using a statistical method to mimic the output of dynamical downscaling. Previous

29 applications of hybrid downscaling used a subjective statistical method to fit the region of

30 interest. It is preferable to use an objective, automated statistical technique that is easily

31 portable to any region. Here, Proper Orthogonal Decomposition Mapping (PODM) is

32 presented as a potential candidate. As a case study, PODM is used to mimic output from

33 the Weather Research and Forecasting (WRF) model used to project climate changes

34 over California's Sierra Nevada mountain range. The results show that PODM robustly

35 predicts WRF temperatures from coarse GCM output, with similar errors across different

36 GCM cases. PODM predictions have 15% lower error than the original hybrid model of

37 Walton et al. (2016). More importantly, PODM can be implemented using automated

38 procedures with limited manual tuning, allowing it to be deployed rapidly. PODM is also

39 shown to compare favorably to state-of-the-art machine-learning algorithms in the

40 context of hybrid downscaling. The use of an objective statistical technique like PODM

41 has the potential to streamline the application of hybrid downscaling for other regions.

42

2

## 1 Introduction

Predictions from global climate models (GCMs) are commonly used to study the impacts of climate change on various aspects of human activities (Gosling et al. 2011). However, existing climate models do not necessarily have the required resolution to accurately model relevant fine-scale features, such as complex topography in mountain ranges and urban heat island effects in cities (McCarthy et al. 2010). While ongoing efforts attempt to resolve these features directly within global climate models, downscaling procedures are practical approaches that allow us to obtain climate change predictions at the desired resolution of a particular impact study.

Downscaling techniques have been widely used to downscale climatic variables, typically precipitation and temperature, from global to regional scales; these techniques have been well-documented, e.g., in Benestad et al. (2008), Fowler et al. (2007), Gutmann et al. (2014), Maraun et al. (2010), and Wilby et al. (1998). There are two main approaches: dynamical downscaling and the statistical downscaling. Dynamical downscaling simulates the complex physical processes that underlie the local climate response using a regional climate model (RCM) forced at its boundaries by reanalysis or GCM output. Statistical downscaling uses a statistical model to map coarse GCM output to station observations or a gridded dataset. A wide variety of statistical downscaling methods are available. For example, Bias Correction with Spatial Disaggregation (Wood et al. 2004), has been particularly successful for downscaling of precipitation suitable as input to regional hydrological models. More complex regression models (Hanssen-Bauer et al.

66  2003; von Storch et al. 1993) have been used to directly model the relationship between

67  the predictors (e.g., sea level pressure) and climatic variables of interest. A similar

68  approach called pattern scaling (Tebaldi and Arblaster 2014) has been used within

69  integrated assessment models.

70

71  More recently, hybrid dynamical-statistical downscaling techniques have been developed

72  (referred to here as "hybrid downscaling"). Hybrid downscaling uses a statistical model

73  to extend the results of dynamical downscaling to multiple GCMs. Under this approach,

74  dynamical downscaling is applied to a small subset of GCMs. Then, a statistical model is

75  trained to mimic the dynamically downscaled results, and is applied to the remaining

76  GCMs.  This saves time when downscaling a large ensemble of GCMs, as applying a

77  statistical models is typically much faster than performing dynamical downscaling.

78  Hybrid downscaling may be valuable in situations where there are important features of

79  the climate change pattern that can only be captured through dynamical downscaling

80  (Berg et al. 2015; Sun et al. 2015a; Sun et al. 2015b).

81

82  The statistical models used in previous hybrid downscaling by Walton et al. (2015) and

83  Walton et al. (2016) require the user to manually investigate the dynamically downscaled

84  data and parameterize salient processes affecting the climate change signal in the region

85  of interest. An open question is whether an objective statistical method could be used

86  instead to minimize manual tuning and streamline the process. Furthermore, previous

87  statistical models were designed to downscale only changes in climatology, but it would

88  be desirable to be able to downscale time series as well. Here we investigate whether

4

89 Proper Orthogonal Decomposition Mapping (PODM; Pau et al., 2014) — a method that

90 was successfully used to downscale hydrological and biogeochemical quantities in Pau et

91 al. (2016) — could be used in hybrid downscaling. The aim of this paper is to determine

92 whether PODM can accurately and robustly emulate dynamically downscaled

93 temperatures when fed coarse GCM output. We also systematically investigate the effects

94 of using different predictors and predictands.

95

96 **2    Problem Setup**

97

98 In Walton et al. (2016), the authors downscaled GCM climate change projections for

99 California's Sierra Nevada mountain range. To capture the effects of complex topography

100 and snow albedo feedback (SAF) on the warming in the Sierra Nevada, high-resolution (3

101 km) simulations were performed with the Weather Research and Forecasting model

102 (WRF; Skamarock et al., 2008). Following a hybrid approach, five GCMs were

103 dynamically downscaled with WRF. Then, the WRF climate change patterns were used

104 to train a statistical model, called StatWRF, that was used to produce WRF-like climate

105 change patterns for an entire ensemble of 35 GCMs. Here we follow a similar procedure,

106 but using PODM and machine learning techniques instead of StatWRF to extend the

107 WRF results.

108

109 As we use the Walton et al. (2016) WRF simulations, it is necessary to briefly describe

110 them. WRF version 3.5 (Skamarock et al. 2008) is used in a configuration with three one-

111 way nested domains of 27, 9, and 3 km resolution, going from the outermost to innermost

112  domain (Figure 1). WRF was coupled to the community Noah land surface model with

113  multi-parameterization options (Noah-MP) (Niu et al. 2011). First, a historical simulation

114  was performed, with WRF forced by North American Regional Reanalysis (NARR;

115  (Mesinger et al. 2006) spanning the period 1991-2000.  (NARR data was provided by the

116  NOAA/OAR/ESRL PSD, Boulder, Colorado, USA, from their website at

117  https://www.esrl.noaa.gov/psd/data/gridded/data.narr.html.)  Next, five "future"

118  simulations were performed, each representing how the 1991-2000 period would have

119  transpired if the mean climate were altered by changes between the 2081-2100 and 1981-

120  2000 periods in a different CMIP5 GCM (Taylor et al. 2012) run under the RCP 8.5

121  scenario (Riahi et al. 2011). The five GCMs used are CNRM-CM5, GFDL-CM3,

122  INMCM4, IPSL-CM5A-LR, and MPI-ESM-LR (see acronym details at

123  http://www.ametsoc.org/PubsAcronymList). Each future simulation is forced with

124  boundary conditions created by adding the difference in GCM monthly climatology

125  (2081–2100 minus 1981–2000) to the 1991–2000 NARR data. This process was applied

126  to temperature, humidity, zonal and meridional winds, and geopotential height. Readers

127  should refer to Walton et al. (2016) for a full description of the WRF model and the

128  dynamical downscaling step.  WRF temperature data used in this study is available from

129  http://research.atmos.ucla.edu/csrl/pub.html.

130

131  In the statistical downscaling step, we employ PODM and machine learning techniques to

132  determine the high-resolution WRF monthly 2m air temperature ($T$) for the innermost

133  domain (D3) from low-resolution GCM output. Our predictand is $T_{\mathrm{WRF,fut}}$, the sequences

134  of monthly $T$ values in the "future" simulations. This is a more difficult task than using

135     the change in climatology, $\Delta\bar{T}_{\text{WRF}} = \bar{T}_{\text{WRF,fut}} - \bar{T}_{\text{WRF,hist}}$, as the predictand, as is considered

136     in Walton et al. (2016) because the statistical model must be able to explain inter-annual

137     variability, not just mean changes. In above, $\bar{T}_{\text{WRF,fut}}$ and $\bar{T}_{\text{WRF,hist}}$ are the 10-year average

138     of the monthly temperature in the 2091-2100 and 1991-2000 periods, respectively.

139

140     Four predictors are considered as input: the monthly NARR 2m air temperature from

141     1991-2000 ($T_{\text{NARR}}$), the monthly NARR surface temperature from the same period (

142     $T_{\text{S,NARR}}$), the difference in GCM monthly 2m air temperature climatology between the

143     2081-2100 and 1981-2000 periods ($\Delta\bar{T}_{\text{GCM}}$), and the difference in GCM monthly surface

144     temperature climatology between the 2081-2100 and 1981-2000 periods ($\Delta\bar{T}_{\text{S,GCM}}$).

145     [Note that the differences in monthly climatology have length 12, while the time series

146     for 1991-2000 has length 120 (10 years $\times$ 12 months/year). So, the sequence of

147     differences in monthly climatology are repeated 10 times when serving as a predictor.]

148     Since the resolutions between NARR and GCM are different, $\Delta\bar{T}_{\text{GCM}}$ and $\Delta\bar{T}_{\text{S,GCM}}$ are

149     interpolated to the resolution of NARR data.

150

151     We also consider an alternative way of preparing the predictors that is more similar to the

152     way the future WRF boundary conditions are constructed. Under this alternate

153     preparation, $T_{\text{NARR}}$ and $\Delta\bar{T}_{\text{GCM}}$ are combined into a single predictor

154 $$T_{\text{BC,fut}} = T_{\text{NARR}} + \Delta\bar{T}_{\text{GCM}}. \tag{1}$$

155  Similarly, we define $T_{S,BC,fut} = T_{S,NARR} + \Delta \bar{T}_{S,GCM}$. As part of the statistical model setup, we

156  determine which set of predictors (and which way of preparing them) minimizes the

157  error. The goal is to see whether surface temperature should be included along with 2m

158  air temperature $T$, and whether it is advantageous to use the combined predictors such as

159  $T_{BC,fut}$ that mimics the future WRF boundary conditions.

160

161  In addition, we consider direct and indirect approaches to obtaining $T_{WRF,fut}$ from the

162  above predictors. The direct approach is to use $T_{WRF,fut}$ as the predictand. The indirect

163  approach is to use the temperature change $\Delta T_{WRF} = T_{WRF,fut} - T_{WRF,hist}$ as the predictand

164  and then add the result to the historical sequence of temperatures $T_{WRF,hist}$ to determine

165  $T_{WRF,fut}$. The indirect approach could be useful since it matches the way the boundary

166  conditions are constructed, i.e. by adding the climate change signal to historical sequence.

167

168  To compare our results with StatWRF (which was designed to downscale climatological

169  changes), we also test PODM's skill in mimicking $\Delta \bar{T}_{WRF}$. We are interested to see if

170  PODM can improve on the accuracy of StatWRF while still capturing the temperature

171  sequences.

172

173  **3  Methods**

174

175    ## 3.1    Proper orthogonal decomposition mapping (PODM)

176

177    We give a summary of the PODM method, as formulated in Pau et al. (2016). This

178    method was first proposed by Robinson et al. (2006) and is derived from the Gappy

179    proper orthogonal decomposition (POD) method (Everson and Sirovich 1995). We first

180    consider a single multivariate predictor $\mathbf{p}$ (e.g. $T_{BC,fut}$ over the region D1) and a single

181    multivariate predictand $\mathbf{f}$ (e.g. $\Delta T_{WRF}$ over the region D3). The training dataset consists

182    of $N$ snapshots of $\mathbf{p}$ and $\mathbf{f}$, taken monthly over the 10-year simulation period using

183    different GCM outputs. For example, $N$ is 600 if output from the the five future

184    simulations is used (each simulation yields 120 monthly snapshots). Given $N$

185    corresponding sets of these $\mathbf{p}$ and $\mathbf{f}$ snapshots, we determine a set of POD bases that are

186    found through a singular value decomposition of the following data matrix:

187
$$\mathbf{W}^{\mathrm{PODM}} = \begin{bmatrix} \mathbf{p}_1 - \overline{\mathbf{p}} & \cdots & \mathbf{p}_N - \overline{\mathbf{p}} \\ \mathbf{f}_1 - \overline{\mathbf{f}} & \cdots & \mathbf{f}_N - \overline{\mathbf{f}} \end{bmatrix} \tag{2}$$

188    where $\overline{\mathbf{p}} = \frac{1}{N}\sum_{i=1}^{N}\mathbf{p}_i$ and $\overline{\mathbf{f}} = \frac{1}{N}\sum_{i=1}^{N}\mathbf{f}_i$. We determine $M$ right singular vectors,

189    $\mathbf{V} = \{\mathbf{v}_1,\ldots,\mathbf{v}_M\}$ corresponding to the $M$ largest singular values for the above data matrix.

190    The POD bases are then given by $\zeta_i = \mathbf{W}^{\mathrm{PODM}}\mathbf{v}_i, i = 1,\ldots,M$, and represent dominant

191    modes of variability in the snapshots within the data matrix $\mathbf{W}^{\mathrm{PODM}}$. By decomposing $\zeta_i$

192    into

193
$$\zeta_i = \begin{bmatrix} \zeta_i^{\mathbf{p}} \\ \zeta_i^{\mathbf{f}} \end{bmatrix} \tag{3}$$

194    where $\zeta_i^{\mathbf{p}}$ and $\zeta_i^{\mathbf{f}}$ are components of the POD basis vector associated with the predictor

195    and the predictand. A linear approximation of $\mathbf{f}$ in the vector space spanned by $\zeta_i$ is then

196    given by

197
$$\mathbf{f} \approx \mathbf{f}_{\mathrm{PODM}} = \overline{\mathbf{f}} + \sum_{i=1}^{M} \gamma_i \zeta_i^{\mathbf{f}} \ . \tag{4}$$

198    The PODM method determines $\gamma = \{\gamma_1, \ldots, \gamma_M\}$ that solves the following the least square

199    problem:

200
$$\gamma = \arg\min_{\alpha} \| \mathbf{p} - \overline{\mathbf{p}} - \sum_{i=1}^{M} \alpha_i \zeta_i^{\mathbf{p}} \|_2 \tag{5}$$

201    where $\|\cdot\|_2$ is the discrete $L^2$ measure. We expect the PODM model to be good if $\zeta_i^{\mathbf{p}}$ is

202    strongly correlated to $\zeta_i^{\mathbf{f}}$. We can augment the above procedure to consider multiple

203    predictors by stacking multiple predictors into a single vector, i.e.

204

205
$$\mathbf{p} = \begin{bmatrix} \mathbf{p}^1 \\ \mathbf{p}^2 \\ \vdots \\ \mathbf{p}^n \end{bmatrix}, \tag{6}$$

206    assuming we have $n$ different predictors. In this study, $n = 4$.

207

208    A key parameter of PODM model is the number of POD bases ($M$) used in the

209    approximation. Determining an appropriate $M$ is a balance of accuracy and stability. If $M$

210    is too large, the PODM approximation becomes unstable due to over-fitting (Everson and

211    Sirovich 1995; Pau et al. 2014). If $M$ is too small, then PODM will not capture all modes

212     of variability and accuracy will be diminished. To determine an appropriate $M$, we will

213     utilize a leave-one-out cross validation (LOOCV) procedure. There are six total cases that

214     we can use for training: results from the five future simulations and the historical case,

215     which represents the zero change scenario (i.e. $\Delta\bar{T}_{\mathrm{GCM}} = 0$ and $\Delta\bar{T}_{\mathrm{S,GCM}} = 0$). Under

216     LOOCV, we train different a PODM model using the historical case and four out of the

217     five GCM cases, with each of the five GCM cases taking a turn being left out. Each

218     PODM model is then used to predict the $T_{\mathrm{WRF}}$ for the case not included in the training

219     dataset. For each case, we determine the mean absolute error (MAE) for a range of $M$. An

220     "optimal" $M$, $M_{\mathrm{opt}}$, is given by one that leads to the lowest MAE, summed over the five

221     GCM cases. This procedure avoids over-fitting $M_{\mathrm{opt}}$ for to any particular GCM case.

222     Since there is seasonal variation in $T_{\mathrm{WRF}}$, we determine a different $M_{\mathrm{opt}}$ for each month

223     of the year, resulting in 12 different values of $M_{\mathrm{opt}}$.

224

225     Before applying the PODM method, it is a good practice to examine whether the POD

226     basis vectors generated from GCMs in the training dataset can be used to closely

227     approximate a new GCM. If a new GCM cannot be closely approximated by the POD

228     bases from the training dataset, then it's unlikely that PODM will be an effective method

229     for downscaling that GCM. The approximation accuracy can be quantified by

230     determining the *projection error* defined as

$$\varepsilon_{\mathrm{proj}}(\mathbf{p}) = \bar{\mathbf{p}} + \sum_{i=1}^{M}((\zeta_i^{\mathbf{p}})^T \mathbf{p})\zeta_i^{\mathbf{p}} - \mathbf{p}, \tag{7}$$

232    where $\bar{\mathbf{p}}$ and $\zeta^{\mathbf{p}}$ are obtained based on the POD procedure applied to snapshots of $\mathbf{p}$ in

233    the training dataset. For $\mathbf{p} = \Delta\bar{T}_{\mathrm{GCM}}$ of a GCM not in the training dataset, a large $\varepsilon_{\mathrm{proj}}(\mathbf{p})$

234    implies there are insufficient similarities between the solutions of that GCM and GCMs

235    used in the training data set. It is unlikely in that case that the predictand can be

236    accurately downscaled through the PODM method. In Section 4.1, we use this LOOCV

237    procedure to determine whether $\Delta\bar{T}_{\mathrm{GCM}}$ within the D1 region can be used as a predictor.

238

239    **3.2    Machine learning based regression approach**

240

241    In this study, we also consider machine learning (ML) based regression approaches in the

242    statistical downscaling step. State-of-the-art machine learning approaches have been used

243    in ecology (Elith et al. 2008; Maloney et al. 2012; Pittman and Brown 2011) and

244    hydrology (Erdal and Karakurt 2013; Nolan et al. 2015). They are also used to downscale

245    satellite images of land surface temperature (Keramitsoglou et al. 2013). Similar to

246    PODM method, machine-learning algorithms could be advantageous as they are typically

247    automatable and require limited manual tuning. We refer these models as ML-based

248    regression models.

249

250    In this paper, different ML methods are used to identify the relationship between the

251    predictors and the predictand. The predictors are defined to be the latitude, longitude,

252    elevation, $T_{\mathrm{NARR}}$, $T_{\mathrm{S,NARR}}$, $\Delta\bar{T}_{\mathrm{GCM}}$, and $\Delta\bar{T}_{\mathrm{S,GCM}}$. $\Delta T_{\mathrm{WRF}}$ is used as the predictand. This

253    amounts to solving the standard regression problem of relating the predictand to a

254    function of predictors:

12

255

$$\Delta T_{\text{WRF}} = f(latitude, longitude, elevation, T_{\text{NARR}}, T_{\text{S,NARR}}, \Delta \overline{T}_{\text{S,GCM}}, \Delta \overline{T}_{\text{S,GCM}})$$

257

258 where $f : \mathbb{R}^n \to \mathbb{R}$, with $n = 7$. Three different machine learning algorithm are used to

259 estimate the function $f$: gradient boosting machines (Freund and Schapire 1997; Friedman

260 et al. 2000; Friedman 2001), extremely randomized trees (Geurts et al. 2006) and elastic

261 net regression method (Zou and Hastie 2005). For further description of these algorithms,

262 see the appendix. Getting the GCM and NARR variables to the WRF resolution is done

263 in two steps. First, the GCM data were interpolated onto NARR grid using bivariate

264 spline approximation on a sphere. Then the GCM and NARR data were interpolated to

265 WRF grid with Gaussian process regression. This combination of preprocessing steps,

266 schematically shown in Figure 2, produced slightly better results than directly

267 interpolating GCM results onto the WRF grid using Gaussian process regression, but the

268 difference is small; our tests showed that the performance of the ML-based regression

269 models only depends weakly on the interpolation scheme.

270

271 **3.3    Evaluation procedure**

272

273 To evaluate how well our statistical downscaling methods (SDMs) emulate WRF, we

274 define an approximation error

275
$$\varepsilon_{\text{SDM}} = T_{\text{SDM}} - T_{\text{WRF,fut}},$$

276  where SDM can be any of the SDMs used in this study. We will primarily be looking at

277  the mean absolute error (MAE), $e_{\mathrm{MAE,SDM}}$ defined as the average absolute value of $\varepsilon_{\mathrm{SDM}}$ .

278  Unless otherwise noted, $e_{\mathrm{MAE,SDM}}$ is assumed to be the average over a 10-year simulation

279  while the monthly average $e_{\mathrm{MAE,SDM}}$ is evaluated for a particular month over a 10-year

280  simulation. The error $e_{\mathrm{MAE,SDM}}$ is used to cross-validate the SDMs based on the LOOCV

281  procedure described in Section 3.1.

282

283  **3.4  Data availability**

284

285  Temperature output generated by hybrid downscaling with StatWRF (Walton et al. 2016)

286  and with above PODM  and ML methods is available from the UCLA Climate Sensitivity

287  Research Lounge website (http://research.atmos.ucla.edu/csrl/pub.html).

288

289  **4   Results**

290

291  **4.1   Initial analysis of the GCM results**

292

293  Here do a preliminary check to see if the GCM patterns are similar enough that any

294  pattern can be approximated by POD bases generated from the remaining four GCMs.

295  We determine the MAE of the projection error, $e_{\mathrm{MAE,proj}}$ of $\Delta \bar{T}_{\mathrm{GCM}}$ , for each of the five

296  GCMs, when it is left out of the training dataset. We also examine how the projection

297  error depends on the number of POD bases. Figure 3 shows that the $e_{\mathrm{MAE,proj}}$ decreases

298    monotonically with the number of POD basis vectors, *M*. The GFDL-CM3 case has the

299    highest averaged $e_{\mathrm{MAE,proj}}$, indicating that the GFDL-CM3's $\Delta \overline{T}_{\mathrm{GCM}}$ patterns are least well

300    approximated from the other GCMs. However, for $M = 40$, the mean, and the standard

301    deviation of $e_{\mathrm{MAE,proj}}$ are 0.2 °C and 0.03 °C respectively. This indicates that for a large

302    number of POD bases, PODM can reasonably approximate the left-out GCM pattern,

303    regardless of which GCM is left out. This gives us confidence that PODM is suitable for

304    application for hybrid downscaling, where are the results of a small set of GCMs are

305    extended to a full ensemble.

306

307    **4.2    Dependence of model accuracy on the predictors and predictand**

308

309    The accuracy of a statistical model varies based on which combinations of predictors and

310    predictands are used. The two options for predictands are considered: predicting the

311    absolute future temperatures, $T_{\mathrm{WRF,fut}}$, or the difference in temperatures between the future

312    and historical simulations, $\Delta T_{\mathrm{WRF}}$. For the predictors there are multiple options: whether

313    to use $\Delta T_{\mathrm{GCM}}$ or $T_{\mathrm{BC,fut}}$, whether to include surface temperature $T_S$ along with 2m air

314    temperature *T*, and which domain over which the predictor is sampled. The domain

315    options are the innermost WRF domain (D3) covering the Sierra Nevada, the

316    intermediate WRF domain covering all of California (D2), and the largest WRF domain

317    covering the entire U.S. West Coast and part of the Pacific Ocean (D1; see Figure 1). In

318    each case, only NARR grid cells within the boundaries of the WRF domain are used as

319    the predictor. The GCM data is interpolated onto these NARR grid cells using a bivariate

320    spline interpolation method. PODM is applied to each different combination of

321    predictors, predictand, and domain to determine how each choice affects the resulting

322    statistical model error, $e_{\text{MAE, PODM}}$ .

323

324    Table 1 shows $e_{\text{MAE, PODM}}$ averaged over all five GCM cases for each combination of

325    choices identified above. In the first three columns, the values in parentheses represent

326    the average $e_{\text{MAE,PODM}}$ over a restricted set of factors. For example, when the predictors

327    are sampled over D1 and $\Delta T_{\text{WRF}}$ is the predictand, the average $e_{\text{MAE,PODM}}$ over different

328    combinations of remaining factors is 0.45 °C. Table 1 clearly shows that PODM can

329    more accurately predict $\Delta T_{\text{WRF}}$ than $T_{\text{WRF,fut}}$ : the average $e_{\text{MAE,PODM}}$ for $\Delta T_{\text{WRF}}$ is 0.53 °C,

330    while the average $e_{\text{MAE,PODM}}$ for $T_{\text{WRF,fut}}$ is 0.99 °C. This makes sense as we would expect

331    $T_{\text{WRF,fut}}$ to be a more difficult predictand to approximate using any method, because the

332    time series $T_{\text{WRF,fut}}$ has larger variability due to inclusion of the seasonal cycle, which is

333    not present in $\Delta T_{\text{WRF}}$ .

334

335    Using the largest domain (D1) results in universally greater accuracy compared with the

336    other domains. With $\Delta T_{\text{WRF}}$ as the predictand, D1 leads to an average $e_{\text{MAE,PODM}}$ that is

337    12%, and 30% lower than D2, and D3 respectively. This result shows that predictor

338    values outside the predictand domain contain valuable predictive information. However,

339    the predictive value decreases with increasing distance from the predictand domain: the

340    average $e_{\text{MAE,PODM}}$ improves 20% between D3 and D2, and only 12% between D2 and D1.

341

342 The inclusion of surface temperature as an auxiliary predictor (i.e. $T_{S,\mathrm{NARR}}$ and $\Delta \bar{T}_{S,\mathrm{GCM}}$ )

343 also universally improves the accuracy of the models. When $T_{\mathrm{WRF,fut}}$ is the predictand,

344 this inclusion reduces the error by 11-17%, depending on the domain size of the

345 predictor. However, when $\Delta T_{\mathrm{WRF}}$ is the predictand, the improvement is smaller: the error

346 is only reduced by 6.5% when the predictor domain is D1 and no reduction is observed

347 for the predictor domain of D2 and D3. Finally, the alternative formulation of the

348 predictors to mimic boundary conditions ( $T_{\mathrm{BC,fut}}$ ) improves accuracy when $T_{\mathrm{WRF,fut}}$ is the

349 predictand, and decreases accuracy when $\Delta T_{\mathrm{WRF}}$ is the predictand ( $e_{\mathrm{MAE,PODM}}$ is larger by

350 up to 6%). This makes sense, as one would generally expect best results when the form of

351 the predictor matches the form of the predictand.

352

353 The above analysis guides our formulation of PODM and other SDMs in the following

354 sections. First, results clearly indicate that the statistical model should be trained to

355 predict $\Delta T_{\mathrm{WRF}}$ , as opposed to predicting $T_{\mathrm{WRF,fut}}$ , regardless of whether the goal is to

356 predict $\Delta T_{\mathrm{WRF}}$ or $T_{\mathrm{WRF,fut}}$ . So, we formulate PODM and the SDMs with $\Delta T_{\mathrm{WRF}}$ as the

357 predictand. We use domain D1, as it leads to a higher accuracy model when compared to

358 the alternatives. $T_{S,\mathrm{NARR}}$ and $\Delta \bar{T}_{S,\mathrm{GCM}}$ are included as predictors since they slightly

359 improve PODM accuracy and the additional computational cost is minimal. Finally, we

360 use $\Delta \bar{T}_{\mathrm{GCM}}$ instead of $T_{\mathrm{BC,fut}}$ as predictor, since it better matches the form of our

361 predictand $\Delta T_{\mathrm{WRF}}$ .

362

363

**4.3    Comparison of statistical methods in approximating WRF**

365

366    The $e_{MAE,SDM}$ of the different SDMs for all the GCM cases are shown in Table 2. The

367    average error for PODM, $e_{MAE,PODM}$, over the five GCMs is 0.44 °C. The monthly average

368    $e_{MAE,PODM}$ varies with month, and the variation is different for each of the GCM cases

369    (Figure 4). However, when averaged over the five GCM cases, the monthly average

370    $e_{MAE,PODM}$ varies more smoothly, with higher errors in the summer months (reaching a

371    maximum in July) and lower errors in the winter months (reaching a minimum in

372    December). ML-based regression models are universally less accurate than the PODM

373    model (Table 2). The average $e_{MAE,SDM}$ of the different ML-based regression models are

374    50%–110% larger than for PODM model.

375

376    PODM also more accurately captures changes in climatology, $\Delta \bar{T}_{WRF}$ (Table 3). The

377    MAEs for the 10-year monthly temperature climatology ($e_{MAE,SDM}^{\Delta \bar{T}}$) of the ML-based

378    regression models are larger than that of PODM model by 45%–130% (Table 3).

379    StatWRF errors are 21% larger than PODM. PODM, like StatWRF, has much lower

380    errors compared to traditional statistical downscaling methods as well, including Bias

381    Correction and Constructed Analogs (BCCA; Maurer and Hidalgo 2008) and Bias

382    Correction with Spatial Disaggregation (BCSD; Wood et al. 2004). (Downscaled CMIP5

383 climate projections using BCCA and BCSD were obtained from http://gdo-

384 dcp.ucllnl.org/downscaled_cmip_projections/, Reclamation 2013.) Interestingly, ML-

385 based regression models do not perform any better than these well-established

386 downscaling techniques despite having higher degrees of complexity. Given the poor

387 performance of ML-based regression models, we focus only on PODM models in

388 subsequent sections.

389

390 **4.4   Comparing spatial distributions of the predictions**

391

392 Here, we compare the spatial patterns of $\Delta \overline{T}_{\text{PODM}}$, $\Delta \overline{T}_{\text{WRF}}$, and $\Delta \overline{T}_{\text{GCM}}$. Figure 5 shows

393 January and July as examples of months where $\Delta \overline{T}_{\text{PODM}}$ poorly and closely matches

394 $\Delta \overline{T}_{\text{WRF}}$, respectively. In January, there is a large disparity between the GCM temperature

395 changes and the WRF-downscaled temperature changes. Importantly, these biases are not

396 in the same direction for CNRM-CM5 and GFDL-CM3.  WRF-downscaled CNRM-CM5

397 has much less warming than CNRM-CM5 (about 1–2 °C).  Meanwhile, WRF-

398 downscaled GFDL-CM3 has much *more* warming than GFDL-CM3 (about 1–2 °C).

399 When relationship between the WRF-downscaled warming and GCM warming is

400 inconsistent between the cases, it is challenging for any statistical model to accurately

401 model it. Thus, PODM struggles to predict WRF-downscaled temperatures in January.

402 PODM better predicts WRF in July, when the relationship between the GCM warming

403 and the WRF-downscaled warming is more consistent.

404

405     We now compare changes in temperature climatology averaged over the five GCM cases,

406     denoted as $\left\langle \Delta \overline{T}_{\mathrm{PODM}} \right\rangle$, $\left\langle \Delta \overline{T}_{\mathrm{GCM}} \right\rangle$, and $\left\langle \Delta \overline{T}_{\mathrm{WRF}} \right\rangle$. Figure 6 shows that differences between

407     $\left\langle \Delta \overline{T}_{\mathrm{PODM}} \right\rangle$ and $\left\langle \Delta \overline{T}_{\mathrm{WRF}} \right\rangle$ are typically small (< 0.5 °C) except for the month of June.

408     PODM is able to capture the fine-scale details present in $\left\langle \Delta \overline{T}_{\mathrm{WRF}} \right\rangle$, such as snow albedo

409     feedback (Walton et al. 2016). This demonstrates that an automated, objective statistical

410     model can capture important features that had to be parameterized in previous hybrid

411     downscaling attempts.

412

413     **4.5    Downscaling 35 CMIP5 GCMs**

414

415     The purpose of hybrid downscaling is to enable rapid, high-quality downscaling of output

416     from a large number of GCMs. Here we demonstrate this capability by applying PODM

417     to 35 CMIP5 GCMs run under the RCP8.5 forcing scenario. Before applying PODM, we

418     check whether the original five GCMs are good representatives of the full ensemble of

419     GCMs. If so, then we can have more confidence that the PODM will have similar

420     accuracy in downscaling the new GCMs as it does in downscaling the original five. To do

421     this, we approximate the full ensemble of GCM warming patterns using POD bases

422     constructed from the five original GCM warming patterns, similar to the analysis in

423     Section 4.1. For the full ensemble, the mean and standard deviation of the approximation

424     errors $e_{\mathrm{MAE,proj}}$ for are 0.24 °C and 0.07 °C, respectively. In comparison, the mean and

425     standard deviation are 0.2 °C and 0.03 °C for the LOOCV errors obtained with the

426     original five GCMs. Since these values are of similar magnitudes, we expect PODM to

427     emulate WRF to a similar degree of accuracy as was found in Section 4.3 when

428     downscaling the entire ensemble.

429

430     To downscale the full ensemble, PODM is trained on data from all six WRF simulations

431     (1 historical + 5 future) as described in Section 2. For some GCMs, surface temperature

432     output is not available in the CMIP5 database. For these GCMs, surface temperature

433     changes are not included as a predictor. This should not significantly alter the accuracy as

434     including surface temperature resulted in only minimal gains (6.5% improvement).

435     Figure 7 shows that PODM captures spatial variations due to snow albedo feedback and

436     the complex topography of the Sierra Nevada that are visible in the WRF solution, but

437     not in the original GCM data.

438

439     **5     Discussion**

440

441     Our results show that PODM can emulate WRF in downscaling temperature changes with

442     errors less than 0.44 °C. This is slightly higher accuracy than the original StatWRF model

443     proposed by Walton et al. (2016). Additionally, PODM is objective and the training steps

444     can be automated. In contrast, StatWRF requires that the user parameterizes salient

445     physical processes affecting the climate change signal in the region of interest and to

446     manually determine the appropriate large-scale predictors. Thus PODM can be applied

447     quickly to any region, without the user needing expert knowledge about the region's

448     climate. We note that the skill of any statistical model — including  PODM — in

449     emulating WRF is likely to be region dependent, so hybrid downscaling users need to

450    verify PODM's skill when applying it elsewhere. It is also important to acknowledge that

451    while the PODM model can be valuable tool for making downscaled projections, it may

452    not enhance our understanding of the climate processes at play. For instance, the way

453    PODM model accounts for the additional warming in the Sierra Nevada due to snow

454    albedo feedback is part of PODM's internal workings that is opaque to the user.

455

456    The skill of PODM could be improved if more dynamically downscaled GCM

457    simulations are included in the training data. Indeed, the PODM model described in this

458    paper could achieve even higher accuracy if the dynamically downscaled GCMs are

459    chosen to best represent the full set of GCMs. The five GCM cases currently used in this

460    paper were chosen to represent the range of temperature and precipitation changes

461    predicted by the GCMs (Walton et al. 2016). However, if just temperature projections are

462    desired, a more representative set of GCMs could be selected by identifying the five

463    GCMs that minimize $e_{\text{MAE,proj}}$ when their POD bases are used to approximate the rest of

464    the ensemble of GCMs. We will study these training procedures in our future work.

465

466    It's important to acknowledge that accuracy results and optimal predictor/predictand

467    combinations for PODM might change in downscaling studies that do not use the pseudo-

468    global warming methodology (PGW). In our study, each future WRF simulation is

469    downscaling of historical NARR plus a change in GCM climatology. Thus, interannual

470    variability of each future simulation is nearly identical to the historical simulation. So,

471    when subtracting the future and historical sequences to determine $\Delta T_{\text{WRF}}$, the interannual

472    variability mostly cancels. In contrast, if raw GCM output were downscaled for the

473   future and historical simulations, as is often the case, then interannual variability between

474   the historical and future simulations will be unrelated.  In this case, $\Delta T_{\mathrm{WRF}}$ could have

475   considerably more variability that PODM needs to capture, and accuracy could be lower.

476   Furthermore, different formulation of the predictors will be needed.  The historical and

477   future GCM sequences $T_{\mathrm{GCM,hist}}$ and $T_{\mathrm{GCM,fut}}$ will probably need to be included, not just

478   the GCM difference in climatology, $\Delta \bar{T}_{\mathrm{GCM}}$.

479

480

481   **6   Conclusions**

482

483   In this study, we have demonstrated that an objective statistical downscaling method,

484   PODM, can approximate WRF temperatures changes with similar or better accuracy than

485   previous statistical methods in California's Sierra Nevada mountain range. ML-based

486   regression methods were also tested, but were found to be much less accurate than

487   PODM in emulating WRF. Our analysis shows that use of a large predictor domain

488   encompassing the entire U.S. West Coast yielded the highest accuracy, even though our

489   predictand domain is limited to the Sierra Nevada mountain range. The inclusion of

490   surface temperature as an additional predictor was found to moderately improve

491   accuracy. Our results also show that if the goal is to predict future temperatures, then

492   statistical models should be designed to predict them as anomalies from the current

493   climate (as opposed to directly predicting the future values).

494

495

508

509    **Appendix**

510

511    In this appendix, gradient boosting machine, Extra-Tree and Elastic Net are reviewed.

512    The first two are so-called ensemble machine learning algorithms based on regression

513    trees and the latter is a linear regression method using a combination of L1 and L2

514    penalties.  Hybrid downscaled output using these methods are available from the UCLA

515    Climate Sensitivity Research Lounge website

516    http://research.atmos.ucla.edu/csrl/pub.html.

517

518    **A.1 Gradient boosting machine (GBM)**

519

520     The gradient boosting machine (GBM) algorithm (Freund and Schapire 1997; Friedman

521     et al. 2000; Friedman 2001) combines iteratively several simple models, called "weak

522     learners", in order to obtain a "strong learner" with improved prediction accuracy. GBM

523     starts by initializing the model by a first guess of a regression tree model (Breiman et al.

524     1984) that maximally reduces the loss function (i.e. least squares). Then at each step a

525     new regression tree model is fitted to the current residual and added to the previous

526     model in order to update the residual, until the number of iteration $K$ is reached. By

527     fitting the regression tree model to the residuals the global model is improved in the

528     regions where it is not accurate.

529

530     GBM expresses the relationship between the scalar predictand (corresponding to $\Delta T_{\mathrm{WRF}}$)

531     and the $n$ scalar predictors ($\boldsymbol{p} = \{p^1,\ldots,p^n\}$, corresponding to the latitude, longitude,

532     elevation, $T_{\mathrm{NARR}}$, $T_{\mathrm{S,NARR}}$, $\Delta \bar{T}_{\mathrm{GCM}}$, and $\Delta \bar{T}_{\mathrm{S,GCM}}$) as an ensemble of $K$ additive functions:

$$f_{\mathrm{GBM}}(\boldsymbol{p}) = \sum_{k=1}^{K} \phi_k(\boldsymbol{p}), \tag{8}$$

534     where $\phi_k(\boldsymbol{p})$ is a regression tree model. Note that K is the number of GBM iteration

535     steps and $\boldsymbol{p}$ represents $m$ different scalar predictors at a particular grid block while **p**

536     represents a multivariate predictor.

537

538     As for any predictive machine learning algorithm, GBM has several parameters that need

539     to be tuned. These parameters are: 1) the complexity of the regression tree, which is

540     represented by the maximum number of split points of the decision tree; 2) K the number

541     of the algorithm iterations; 3) the learning rate, which is a relatively small positive value

542     between 0 and 1, and inversely proportional to K; and 4) the fraction of training data that

543     is used as a training subsample at each iterative step. To choose the combination of these

544     parameters that produce the best predictive GBM model the LOOCV procedure

545     (described in Section 3.1) combined with the so-called search grid method have been

546     used in this study. This method is based on predefining a grid of GBM parameters

547     combinations, then for each combination a GBM model is estimated. The best

548     combination is selected as the one that produce the most accurate model using the

549     LOOCV procedure. We used the minimization of the MAE as the accuracy criteria to

550     select the best combination.

551

552     The XGBoost python library (https://github.com/dmlc/xgboost), which is a relatively new

553     efficient implementation of GBM method, is used here. The performance of XGBoost has

554     been demonstrated in multiple data mining and machine learning challenges (Chen and

555     Guestrin 2016). We refer readers to Chen and Guestrin (2016) for details of the XGBoost

556     algorithm, especially the advanced features that have been implemented in it.

557

558     **A.2 Extremely Randomized Trees (Extra-Trees)**

559

560     Similar to GBM the Extra-Tree algorithm (Geurts et al. 2006) is based on a simple

561     averaging of the weak learner while the boosting algorithm of GBM is built upon a

562     constructive iterative strategy.  It builds a set of regression trees, which are trained by

563     selecting the decision trees splits points at random. In other words, instead of selecting

564    the splits points that are locally optimal, these splits points are selected randomly. The

565    predictions of each regression trees are simply averaged to create the final prediction.

566

567    The Extra-Tree procedure has two main parameters that need to be tuned. These

568    parameters are the maximum number of splits points of each regression tree and K the

569    number of regression trees of the ensemble. As for the GBM model the best combination

570    is selected using a search grid and the previously described LOOCV procedure. In this

571    work, we have used the Extra-Trees implementation of the scikit-learn python library

572    (Pedregosa et al. 2011).

573

574

575    **A.3 Elastic net linear regression**

576

577    We consider the standard linear regression model, which is defined for the given $n$ scalar

578    predictors $\boldsymbol{p} = \{p^1, \ldots, p^n\}$ (corresponding to the latitude, longitude, elevation, $T_{\mathrm{NARR}}$,

579    $T_{\mathrm{S,NARR}}$, $\Delta\overline{T}_{\mathrm{GCM}}$, and $\Delta\overline{T}_{\mathrm{S,GCM}}$) and the scalar predictand $f$ (corresponding to $\Delta T_{\mathrm{WRF}}$) by:

580
$$f_{\mathrm{linear}}(\boldsymbol{p}) = \sum_{i=1}^{n} \beta_i p^i,$$

581    The standard approach to estimate the regression coefficients $\beta = \{\beta_1, \ldots, \beta_n\}$ is to use the

582    ordinary least squares algorithm (OLS). However it is well known that the OLS often

583    underperforms in term of prediction accuracy compared to other linear techniques such as

584    ridge regression (Hoerl and Kennard 2000) and LASSO (Tibshirani 1996). The first one

585    applies an L2 penalty on the coefficients and the second one applies an L1 penalty.

586　Elastic Net regression method (Elastic Net, Zou and Hastie 2005) applies a convex

587　combination of L1 and L2 penalties on the regression coefficients. There are two

588　parameters to optimize: the ratio of L1 penalty to L2 penalty, and the magnitude of the

589　total penalty. These two parameters are tuned using the same methodology as for GBM

590　and Extra-Tree algorithms. The scikit-learn python library (Pedregosa et al. 2011)

591　implementation of the Elastic Net has been used in this work.

592

593　**Reference**

594

595　Benestad, R. E., I. Hanssen-Bauer, and D. Chen, 2008: *Empirical-statistical downscaling.*

596　　　World Scientific.

597　Berg, N., A. Hall, F. Sun, S. Capps, D. Walton, B. Langenbrunner, and D. Neelin, 2015:

598　　　Twenty-First-Century Precipitation Changes over the Los Angeles Region.

599　　　*Journal of Climate*, **28,** 401-421.

600　Breiman, L., J. Friedman, C. J. Stone, and R. A. Olshen, 1984: *Classification and*

601　　　*regression trees.* CRC press.

602　Chen, T., and C. Guestrin, 2016: XGBoost: A scalable tree boosting system. *arXiv*

603　　　*preprint arXiv:1603.02754.*

604　Elith, J., J. R. Leathwick, and T. Hastie, 2008: A working guide to boosted regression

605　　　trees. *Journal of Animal Ecology*, **77,** 802-813.

606　Erdal, H. I., and O. Karakurt, 2013: Advancing monthly streamflow prediction accuracy

607　　　of CART models using ensemble learning paradigms. *J Hydrol*, **477,** 119-128.

608      Everson, R., and L. Sirovich, 1995: Karhunen–Loeve procedure for gappy data. *Journal*

609          *of the Optical Society of America A*, **12,** 1657-1664.

610      Fowler, H. J., S. Blenkinsop, and C. Tebaldi, 2007: Linking climate change modelling to

611          impacts studies: recent advances in downscaling techniques for hydrological

612          modelling. *International Journal of Climatology*, **27,** 1547-1578.

613      Freund, Y., and R. E. Schapire, 1997: A Decision-Theoretic Generalization of On-Line

614          Learning and an Application to Boosting. *Journal of Computer and System*

615          *Sciences*, **55,** 119-139.

616      Friedman, J., T. Hastie, and R. Tibshirani, 2000: Additive logistic regression: a statistical

617          view of boosting (With discussion and a rejoinder by the authors). *The Annals of*

618          *Statistics*, **28,** 337-407.

619      Friedman, J. H., 2001: Greedy function approximation: A gradient boosting machine.

620          *Annals of Statistics*, **29,** 1189-1232.

621      Geurts, P., D. Ernst, and L. Wehenkel, 2006: Extremely randomized trees. *Mach Learn*,

622          **63,** 3-42.

623      Gosling, S. N., and Coauthors, 2011: A review of recent developments in climate change

624          science. Part II: The global-scale impacts of climate change. *Progress in Physical*

625          *Geography*, **35,** 443-464.

626      Gutmann, E., T. Pruitt, M. P. Clark, L. Brekke, J. R. Arnold, D. A. Raff, and R. M.

627          Rasmussen, 2014: An intercomparison of statistical downscaling methods used

628          for water resource assessments in the United States. *Water Resour Res*, **50,** 7167-

629          7186.

630     Hanssen-Bauer, I., E. J. Forland, J. E. Haugen, and O. E. Tveito, 2003: Temperature and

631         precipitation scenarios for Norway: comparison of results from dynamical and

632         empirical downscaling. *Climate Research*, **25,** 15-27.

633     Hoerl, A. E., and R. W. Kennard, 2000: Ridge regression: Biased estimation for

634         nonorthogonal problems. *Technometrics*, **42,** 80-86.

635     Keramitsoglou, I., C. T. Kiranoudis, and Q. Weng, 2013: Downscaling Geostationary

636         Land Surface Temperature Imagery for Urban Analysis. *IEEE Geoscience and*

637         *Remote Sensing Letters*, **10,** 1253-1257.

638     Maloney, K. O., M. Schmid, and D. E. Weller, 2012: Applying additive modelling and

639         gradient boosting to assess the effects of watershed and reach characteristics on

640         riverine assemblages. *Methods in Ecology and Evolution*, **3,** 116-128.

641     Maraun, D., and Coauthors, 2010: Precipitation downscaling under climate change:

642         Recent developments to bridge the gap between dynamical models and the end

643         user. *Reviews of Geophysics*, **48**.

644     Maurer, E. P., and H. G. Hidalgo, 2008: Utility of daily vs. monthly large-scale climate

645         data: an intercomparison of two statistical downscaling methods. *Hydrol. Earth*

646         *Syst. Sci.*, **12,** 551-563.

647     McCarthy, M. P., M. J. Best, and R. A. Betts, 2010: Climate change in cities due to

648         global warming and urban effects. *Geophys Res Lett*, **37,** n/a-n/a.

649     Mesinger, F., and Coauthors, 2006: North American Regional Reanalysis. *Bulletin of the*

650         *American Meteorological Society*, **87,** 343-360.

651    Niu, G. Y., and Coauthors, 2011: The community Noah land surface model with

652        multiparameterization options (Noah-MP): 1. Model description and evaluation

653        with local-scale measurements. *J Geophys Res-Atmos*, **116**.

654    Nolan, B. T., M. N. Fienen, and D. L. Lorenz, 2015: A statistical learning framework for

655        groundwater nitrate models of the Central Valley, California, USA. *J Hydrol*, **531,**

656        **Part 3,** 902-911.

657    Pau, G. S. H., G. Bisht, and W. J. Riley, 2014: A reduced-order modeling approach to

658        represent subgrid-scale hydrological dynamics for land-surface simulations:

659        application in a polygonal tundra landscape. *Geosci. Model Dev.*, **7,** 2091-2105.

660    Pau, G. S. H., C. Shen, W. J. Riley, and Y. Liu, 2016: Accurate and efficient prediction

661        of fine-resolution hydrologic and carbon dynamic simulations from coarse-

662        resolution models. *Water Resour Res*, **52,** 791-812.

663    Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine Learning in Python. *Journal*

664        *of Machine Learning Research*, **12,** 2825-2830.

665    Pittman, S. J., and K. A. Brown, 2011: Multi-scale approach for predicting fish species

666        distributions across coral reef seascapes. *Plos One*, **6,** e20583.

667    Reclamation, 2013: 'Downscaled CMIP3 and CMIP5 Climate and Hydrology

668        Projections: Release of Downscaled CMIP5 Climate Projections, Comparison

669        with preceding Information, and Summary of User Needs', prepared by the U.S.

670        Department of the Interior, Bureau of Reclamation, Technical Services Center,

671        Denver, Colorado. 47pp.

672    Riahi, K., and Coauthors, 2011: RCP 8.5—A scenario of comparatively high greenhouse

673        gas emissions. *Climatic Change*, **109,** 33.

674      Robinson, T., M. Eldred, K. Willcox, and R. Haimes, 2006: Strategies for Multifidelity

675          Optimization with Variable Dimensional Hierarchical Models. *47th*

676          *AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials*

677          *Conference*, American Institute of Aeronautics and Astronautics.

678      Skamarock, W. C., and Coauthors, 2008: A description of the advanced research WRF

679          Ver. 30. NCAR Technical Note.

680      Sun, A. Y., R. M. Miranda, and X. L. Xu, 2015a: Development of multi-metamodels to

681          support surface water quality management and decision making. *Environ Earth*

682          *Sci*, **73,** 423-434.

683      Sun, F., D. B. Walton, and A. Hall, 2015b: A Hybrid Dynamical–Statistical Downscaling

684          Technique. Part II: End-of-Century Warming Projections Predict a New Climate

685          State in the Los Angeles Region. *Journal of Climate*, **28,** 4618-4636.

686      Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An Overview of CMIP5 and the

687          Experiment Design. *Bulletin of the American Meteorological Society*, **93,** 485-

688          498.

689      Tebaldi, C., and J. Arblaster, 2014: Pattern scaling: Its strengths and limitations, and an

690          update on the latest model simulations. *Climatic Change*, **122,** 459-471.

691      Tibshirani, R., 1996: Regression shrinkage and selection via the Lasso. *J Roy Stat Soc B*

692          *Met*, **58,** 267-288.

693      von Storch, H., E. Zorita, and U. Cubasch, 1993: Downscaling of Global Climate Change

694          Estimates to Regional Scales: An Application to Iberian Rainfall in Wintertime.

695          *Journal of Climate*, **6,** 1161-1171.

696    Walton, D., A. Hall, F. Sun, M. Schwartz, and N. Berg, 2016: Incorporating Snow

697        Albedo Feedback into Downscaled Temperature and Snow Cover Projections for

698        California's Sierra Nevada. *Journal of Climate*, **30**, 1417-1438, doi:

699        http://dx.doi.org/10.1175/JCLI-D-16-0168.1.

700    Walton, D. B., F. Sun, A. Hall, and S. Capps, 2015: A Hybrid Dynamical–Statistical

701        Downscaling Technique. Part I: Development and Validation of the Technique.

702        *Journal of Climate*, **28,** 4597-4617.

703    Wilby, R. L., T. M. L. Wigley, D. Conway, P. D. Jones, B. C. Hewitson, J. Main, and D.

704        S. Wilks, 1998: Statistical downscaling of general circulation model output: A

705        comparison of methods. *Water Resour Res*, **34,** 2995-3008.

706    Wood, A. W., L. R. Leung, V. Sridhar, and D. P. Lettenmaier, 2004: Hydrologic

707        Implications of Dynamical and Statistical Approaches to Downscaling Climate

708        Model Outputs. *Climatic Change*, **62,** 189-216.

709    Zou, H., and T. Hastie, 2005: Regularization and variable selection via the elastic net. *J*

710        *Roy Stat Soc B*, **67,** 301-320.

711

712    Table 1: Averaged $e_{\text{MAE,PODM}}$ for different combinations predictors and predictands. The

713    value in parenthesis is the averaged $e_{\text{MAE,PODM}}$ over that particular parameter.

| Predictand (°C) | Predictor | | | $e_{\text{MAE,PODM}}$, °C |
| --- | --- | --- | --- | --- |
| | Domain (°C) | $T_{S,\text{NARR}}$ and $\Delta\bar{T}_{S,\text{GCM}}$ (°C) | $\Delta\bar{T}_{\text{GCM}}$ or $T_{\text{BC,fut}}$ | |
| $T_{\text{WRF}}$ (0.99) | D01 (0.68) | False (0.73) | $\Delta\bar{T}_{\text{GCM}}$ | 0.77 |
| | | | $T_{\text{BC,fut}}$ | 0.69 |
| | | True (0.62) | $\Delta\bar{T}_{\text{GCM}}$ | 0.63 |
| | | | $T_{\text{BC,fut}}$ | 0.61 |
| | D02 (0.96) | False (1.05) | $\Delta\bar{T}_{\text{GCM}}$ | 1.04 |
| | | | $T_{\text{BC,fut}}$ | 1.05 |
| | | True (0.87) | $\Delta\bar{T}_{\text{GCM}}$ | 0.90 |
| | | | $T_{\text{BC,fut}}$ | 0.84 |
| | D03 (1.34) | False (1.42) | $\Delta\bar{T}_{\text{GCM}}$ | 1.44 |
| | | | $T_{\text{BC,fut}}$ | 1.39 |
| | | True (1.26) | $\Delta\bar{T}_{\text{GCM}}$ | 1.37 |
| | | | $T_{\text{BC,fut}}$ | 1.15 |
| $\Delta T_{\text{WRF}}$ | D01 | False | $\Delta\bar{T}_{\text{GCM}}$ | 0.45 |

| (0.53) | (0.45) | (0.46) | $T_{\text{BC,fut}}$ | 0.46 |
|---|---|---|---|---|
| | | True (0.43) | $\Delta\bar{T}_{\text{GCM}}$ | 0.43 |
| | | | $T_{\text{BC,fut}}$ | 0.43 |
| | D02 (0.51) | False (0.51) | $\Delta\bar{T}_{\text{GCM}}$ | 0.49 |
| | | | $T_{\text{BC,fut}}$ | 0.52 |
| | | True (0.51) | $\Delta\bar{T}_{\text{GCM}}$ | 0.50 |
| | | | $T_{\text{BC,fut}}$ | 0.52 |
| | D03 (0.64) | False (0.64) | $\Delta\bar{T}_{\text{GCM}}$ | 0.64 |
| | | | $T_{\text{BC,fut}}$ | 0.65 |
| | | True (0.64) | $\Delta\bar{T}_{\text{GCM}}$ | 0.62 |
| | | | $T_{\text{BC,fut}}$ | 0.66 |

714

715

716  Table 2: The MAEs ($e_{\text{MAE,SDM}}$) of $\Delta T_{\text{SDM}}$ for the different SDMs, for each GCM cases,

717  and the ensemble averages.

| GCM | PODM | Elastic Net | GBM | Extra-Tree |
|---|---|---|---|---|
| CNRM-CM5 | 0.40 | 0.86 | 0.80 | 0.87 |
| GFDL-CM3 | 0.54 | 0.83 | 0.94 | 1.26 |

| | | | | |
|---|---|---|---|---|
| INMCM4 | 0.38 | 0.46 | 0.97 | 0.90 |
| IPSL-CM5A-LR | 0.48 | 0.72 | 1.14 | 1.15 |
| MPI-ESM-LR | 0.38 | 0.44 | 0.50 | 0.51 |
| **Average** | **0.44** | **0.66** | **0.87** | **0.94** |

718 Table 3: The MAEs ($e_{\text{MAE,SDM}}^{\Delta\bar{T}}$) of $\Delta\bar{T}_{\text{SDM}}$ for the different SDMs, for each GCM cases,

719 and the ensemble averages. BCCA stands for Bias Correction and Constructed Analogs

720 and BCSD stands for Bias Correction with Spatial Disaggregation. MAE data for

721 StatWRF, BCCA, BCSD, and linter interpolation are from Walton et al. (2016).

| GCM | PODM | Elastic Net | GBM | Extra-Tree | Stat-WRF | BCCA | BCSD | Linear inter-polation |
|---|---|---|---|---|---|---|---|---|
| CNRM-CM5 | 0.34 | 0.76 | 0.66 | 0.76 | 0.52 | 0.49 | 0.89 | 0.85 |
| GFDL-CM3 | 0.47 | 0.66 | 0.80 | 1.25 | 0.61 | 1.18 | 1.08 | 0.75 |
| INMCM4 | 0.31 | 0.34 | 0. 88 | 0.84 | 0.47 | 0.91 | 0.94 | 0.48 |
| IPSL-CM5A-LR | 0.41 | 0.68 | 1.10 | 1.11 | 0.31 | 0.78 | 0.56 | 0.43 |
| MPI-ESM-LR | 0.29 | 0.24 | 0. 28 | 0.30 | 0.35 | 0.63 | 0.43 | 0.44 |

| Average | 0.37 | 0.54 | 0. 74 | 0.85 | 0.45 | 0.80 | 0.78 | 0.59 |

722

723

724



725 Figure 1: (Taken from Walton et al., 2016) (a) Elevation (meters) and model setup with

726 three one-way nested WRF domains (D1, D2, and D3) at horizontal resolutions of 27, 9,

727 and 3 km. (b) Innermost domain elevation (meters).

728



730 Figure 2: Predictors data pre-processing steps for the downscaling approach based on

731 ML-based regression models.

732

38

733

Figure 3: The projection error $e_{\text{MAE,proj}}$ for approximating the left-out $\Delta \bar{T}_{\text{GCM}}$ pattern,

when $M$ bases are used. Here, $e_{\text{MAE,proj}}$ is an average over all 12 calendar months.

736



737

Figure 4: Monthly averaged $e_{\text{MAE,PODM}}$ versus month for different GCM cases. The values

in the parentheses are the averages for each of the GCM cases.

740

Figure 5: $\Delta\bar{T}_{\mathrm{GCM}}$, $\Delta\bar{T}_{\mathrm{WRF}}$, $\Delta\bar{T}_{\mathrm{PODM}}$, and $\varepsilon^{\Delta\bar{T}}_{\mathrm{PODM}} = \Delta\bar{T}_{\mathrm{PODM}} - \Delta\bar{T}_{\mathrm{WRF}}$ in January and July. Only results from the CNRM-CM5 and GFDL-CM3 cases are shown.

Figure 6: Changes in temperature climatology averaged over five GCM cases produced from three sources. (row 1) GCM changes $\langle \Delta \bar{T}_{\text{GCM}} \rangle$. (row 2) WRF changes $\langle \Delta \bar{T}_{\text{WRF}} \rangle$. (row 3) PODM changes produced via cross validation $\langle \Delta \bar{T}_{\text{PODM}} \rangle$. (row 4) differences between WRF and PODM $\varepsilon_{\text{PODM}}^{\langle \Delta \bar{T} \rangle}$ for each of month, averaged over the five GCMs cases.

754

Figure 7: The mean and the standard deviation of the monthly $\Delta\overline{T}_{\text{PODM}}$ compared to

756    $\Delta\overline{T}_{\text{GCM}}$.


757