1  Improving climate model accuracy by exploring parameter space with an O(10⁵) member

2  ensemble and emulator

3  Sihan Li[1,2], David E. Rupp[3], Linnia Hawkins[3,6], Philip W. Mote[3,6], Doug McNeall[4], Sarah

4  N. Sparrow[2], David C. H. Wallom[2], Richard A. Betts[4,5], Justin J. Wettstein[6,7,8]

5  [1]Environmental Change Institute, School of Geography and the Environment, University

6  of Oxford, Oxford, United Kingdom

7  [2]Oxford e-Research Centre, University of Oxford, Oxford, United Kingdom

8  [3]Oregon Climate Change Research Institute, College of Earth, Ocean, and Atmospheric

9  Science, Oregon State University, Corvallis, Oregon

10  [4]Met Office Hadley Centre, FitzRoy Road, Exeter, United Kingdom

11  [5]College of Life and Environmental Sciences, University of Exeter, Exeter, UK

12  [6]College of Earth, Ocean, and Atmospheric Science, Oregon State University, Corvallis,

13  Oregon

14  [7]Geophysical Institute, University of Bergen, Bergen, Norway

15  [8]Bjerknes Centre for Climate Change Research, Bergen, Norway

16  *Correspondence to*: Sihan Li (sihan.li@ouce.ox.ac.uk)

17

18

19

20

21

22

23

Geoscientific
Model Development
Discussions

Open Access

EGU

24  **Abstract**
25  Understanding the unfolding challenges of climate change relies on climate models, many

26  of which have large summer warm and dry biases over Northern Hemisphere continental

27  mid-latitudes. This work, using the example of the model used in the updated version of

28  the weather@home distributed climate model framework, shows the potential for

29  improving climate model simulations through a multi-phased parameter refinement

30  approach, particularly over northwestern United States (NWUS). Each phase consists of 1)

31  creating a perturbed physics ensemble with the coupled global - regional atmospheric

32  model, 2) building statistical emulators that estimate climate metrics as functions of

33  parameter values, 3) and using the emulators to further refine the parameter space. The

34  refinement process includes sensitivity analyses to identify the most influential parameters

35  for various model output metrics; results are then used to cull parameters with little

36  influence. Three phases of this iterative process are carried out before the results are

37  considered to be satisfactory; that is, a handful of parameter sets are identified that meet

38  acceptable bias reduction criteria. Results not only indicate that 74% of the NWUS regional

39  warm biases can be reduced by refining global atmospheric parameters that control

40  convection and hydrometeor transport, and land surface parameters that affect plant

41  photosynthesis, transpiration and evaporation, but also suggest that this iterative approach

42  to perturbed physics has an important role to play in the evolution of physical

43  parameterizations.

44

45    **Introduction**

46    Boreal summer (June-July-August, JJA) warm and dry biases over North Hemisphere (NH)

47    continental midlatitudes are common in many global and regional climate models (e.g.,

48    Boberg and Christensen, 2012; Mearns et al., 2012; Mueller and Seneviratne, 2014;

49    Kotlarski et al., 2014; Cheruy et al., 2014; Merrifield and Xie, 2016), including very high

50    resolution convection-permitting models (e.g. Liu et al., 2017). These biases can have non-

51    negligible impacts on climate change studies, particularly where relationships are non-

52    linear, such as is the case of surface latent heat flux as a function of water storage (e.g.

53    Rupp et al., 2017). Biases in present-day climate model simulations cast doubt on the

54    reliability of the future climate projections from those models. As shown by Boberg and

55    Christensen (2012), after applying a bias correction conditioned on temperature to account

56    for model deficiencies, the Mediterranean summer temperature projections were reduced

57    by up to 1°C. Cheruy et al. (2014) demonstrated that of the climate models contributing to

58    the Coupled Model Intercomparison Project Phase5 (CMIP5), the models that simulate a

59    higher-than-average warming overestimated the present climate net shortwave radiation

60    which increased more than multi-model average in the future; those models also showed a

61    higher-than-average reduction of evaporative faction in areas with soil moisture-limited

62    evaporation regimes. Both studies suggested that models with a larger warm bias in surface

63    temperature tend to overestimate the projected warming. The implication of the warm bias

64    goes beyond climate model simulations, as many impact modeling (e.g. hydrological, fire,

65    crop modeling) studies (e.g. Brown et al., 2004; Fowler et al., 2007; Hawkins et al., 2013;

66    Rosenzweig et al., 2014) use climate model simulation results as driving data. Recently,

67    there have been coordinated research efforts (Morcrette et al., 2018; van Weverberg et al.,

Geoscientific
Model Development
Discussions

68    2018; Ma et al., 2018; Zhang et al., 2018) to better understand the causes of the near-surface

69    atmospheric temperature biases through process level understanding and to identify the

70    model deficiencies that generate the bias. These studies suggest that biases in the net

71    shortwave and downward longwave fluxes as well as surface evaporative fraction are

72    contributors to surface temperature bias.

73

74    Older generation Hadley Centre coupled models (HadCM2 and HadCM3), and

75    atmospherere-only global (HadAM) and regional (HadRM) models have been used in

76    numerous attribution studies (e.g., Tett et al., 1996; Stott et al., 2004; Otto et al., 2012;

77    Rupp et al., 2017a; van Oldenborgh et al., 2016; Schaller et al., 2016; van Oldenborgh et

78    al., 2017; Uhe et al., 2018), and the same models have been used for future projections

79    (e.g., Rupp and Li, 2017; Rupp et al., 2017b; Guillod et al., 2018).  These model families

80    exhibit warm and dry biases during JJA over continental midlatitudes, biases that have

81    persisted over model generations and enhancements (e.g., Massey et al., 2015; Li et al.,

82    2015; Guillod et al., 2017). The more recent generations of Hadley Centre models –

83    HadGEMx (HadGEM1, Johns et al, 2016;  HadGEM2, Collins et al., 2008 ) also have the

84    same biases to some extent.

85

86    Many of the aforementioned studies using HadAM and HadRM generated simulations

87    through a distributed computing system known as climateprediction.net (CPDN, Allen et

88    al., 1999), within which a system called weather@home is used to dynamically downscale

89    global simulations using regional climate models (Massey et al., 2015; Mote et al., 2016;

90    Guillod et al., 2017).  As with the previous version of weather@home, the current

91   operational version of weather@home (version 2: weather@home2) uses the coupled

92   HadAM3P/HadRM3P with the atmosphere component based on HadCM3 (Gordon et al.,

93   2000), but updates the land surface scheme from the Met Office Surface Exchange Scheme

94   version 1 (MOSES1, Cox et al., 1999) to version 2(MOSES2, Essery et al., 2003).

95

96   Although the current model version in weather@home2 produces some global-scale

97   improvements in the global model's simulation of the seasonal mean climate, warm biases

98   in JJA increase over North America north of roughly 40° compared with the previous

99   version in weather@home1 (Fig. 2 in Guillod et al., 2017).  The warm and dry JJA biases

100  appear clearly in the regional model simulations over the northwestern US region (NWUS,

101  defined here as all the continental US land points west of 110° and between 40°N-49°N -

102  the grey bounding box in Fig.S1). These biases may be related to, among other things, an

103  imperfect parameterization of certain cloud processes, leading to excess downward solar

104  radiation at the surface, which in turn triggers warm and dry summer conditions that are

105  further amplified by biases in the surface energy and water balance in the land surface

106  model (Sippel et al., 2016; Guillod et al., 2017).  The fact that recent model enhancements

107  did not reduce biases over most of the northwest US motivates the present study, which

108  aims at reducing these warm/dry biases by way of adjusting parameter values, herein

109  referred to as 'parameter refinement'.

110

111  Many small-scale atmospheric processes have significant impacts on large-scale climate

112  states. Processes such as precipitation formation, radiative balance, and convection, occur

113  at scales smaller than the spatial resolution explicitly resolved by climate models, though

Geoscientific
Model Development
Discussions

114     very high resolution regional climate models are able to resolve or partially resolve some

115     of these processes (e.g., convection). These processes must be represented by

116     parameterizations that include parameters whose uncertainty are often high because: 1)

117     there are insufficient observations with which to constrain the parameters, 2) a single

118     parameter is inadequate to represent the different ways a process behaves across the globe,

119     and/or 3) there is incomplete understanding of the physical process (Hourdin et al., 2013).

120     Many studies have demonstrated the importance of considering parameterization

121     uncertainty in the simulation of present and future climates by perturbing single and

122     multiple model parameters within plausible parameter ranges usually established by expert

123     judgment (e.g., Murphy et al., 2004; Stainforth et al., 2005; Sanderson et al., 2008a, b,

124     2010, 2011; Collins et al., 2011; Bellprat et al., 2012a,b, 2016). These studies have argued

125     for careful tuning of models not only to reduce model parameter uncertainties by selecting

126     parameter values that result in a better match between model simulation results with

127     observations, but also to better understand relationships among physical processes within

128     the climate system via systematic experiments that alter individual parameter values or

129     combinations thereof, in order to assess model responses to perturbing parameters.

130

131     Improving a model by parameter refinement can be an iterative process of modifying

132     parameter values, running a climate simulation, comparing model output to observations,

133     and refining the parameter values again (Mauritsen et al., 2012; Schirber et al., 2013). This

134     iterative process can be both computationally expensive and labor-intensive. Any

135     parameter refinement process performed with the intent of improving the model also

136     involves unavoidably arbitrary decisions - though guided by expert judgement - about

137     which parameter(s) to adjust, which metric(s) to evaluate (i.e., which feature(s) of the

138     climate system to simulate at some level of accuracy), and which observational dataset(s)

139     to use as the basis for the evaluation metric(s). Nonetheless, model tuning through

140     parameter refinement is invariably needed to better match model simulations with

141     observations (Schirber et al., 2013).

142

143     One systematic, yet computationally demanding, approach to model tuning is through

144     perturbed physics experiments (Allen et al., 1999; Murphy et al., 2004). These experiments

145     use a perturbed physics ensemble (PPE) of simulations from a single model where a

146     handful of uncertain model parameters are varied systematically. Each set of perturbed

147     parameter (PP) values is considered to be a different model variant - a PP set refers to a

148     combination of parameter values from herein on. PPEs can be treated as a sparse sample

149     of behaviours from a vast, high-dimensional parameter space (Williamson et al., 2013). A

150     PPE directly informs us about model behaviour at those points in parameter space where

151     the model is run (the PP sets), and helps us infer model behavior in nearby parameter space

152     where the model has not been run.

153

154     Studies of climate model tuning using PPEs generally fall into three categories. The first

155     category makes only direct use of the ensemble itself (e.g., Murphy et al., 2004; Rowlands

156     et al., 2012) by screening out ensemble members that are deemed too far from the observed

157     target metrics. This is often referred to as ensemble filtering. However, this approach can

158     overlook certain critical parts of the parameter space not sampled by the PPE. One

159     promising improvement of this approach is to estimate the response of metric(s) in a

160     geophysical (e.g., atmospheric) model to parameter perturbations using a computationally

161     efficient statistical model (i.e. emulator) that is trained from the PPE results. The

162     emulator's skill is evaluated based on its metric prediction accuracy using independent

163     simulations of the model and, if deemed sufficiently skilful, can be used to estimate the

164     model's output metrics as a function of the model parameters in the parameter space not

165     sampled by the PPE.

166

167     The second category uses a PPE to train a statistical emulator, or establish some cost

168     function, which is then used to automatically search for optimal parameter values that

169     produce simulations closest to observations (e.g., Bellprat et al., 2012a, 2016; Zhang et al.,

170     2015; Tett et al., 2017). These studies advocated for this approach particularly because of

171     the efficiency and automation of available searching algorithms. However, as with any

172     model evaluation effort, the use of a cost function with multiple target metrics means that

173     optima for different metrics may occur at different parameter values. This approach

174     (automatically searching for optimal parameters) also runs the risk of being trapped into

175     local minima in the associated cost function; thus, searching results are heavily dependent

176     on the initial parameter values. Admittedly, the idea of automatic searching to obtain

177     optimal combinations of model parameters is appealing, but in reality there is still a high

178     level of subjectivity, e.g. selecting which model performance metrics and observation(s) to

179     use in evaluating the model, and the methods of optimization and searching algorithm.

180

181     Unlike the second category, which searches for the optimal parameter values that result in

182     the closest match to observations, the third category, named 'history matching' (McNeall

Geoscientific
Model Development
Discussions

183    et al., 2013, 2016; Williamson et al., 2013, 2015, 2017), seeks to rule out parameter choices

184    that do not adequately reproduce observations. History matching uses PPEs to train

185    statistical emulators that predict key metrics from the model output, and then uses the

186    emulators to rule out parameter space that is implausible. Williamson et al. (2017)

187    demonstrated that this method is more powerful when iterative steps are taken to rule out

188    implausible parameter space, where each step helps refine the parameter space containing

189    potentially better performing model variants. A drawback is that iterative history matching

190    requires more model runs in the not-ruled-out-yet parameter space for later iterations. The

191    method we adopted in this study fits in the third category, where the parameter values were

192    refined through phases of experiments.

193

194    All three approaches begin with an initial PPE, which can be computationally expensive

195    even with a modest number of free parameters. To cope with the computational demand,

196    many previous studies have generated PPEs from a global climate model (GCM) using

197    CPDN. The studies span a range of topics, from the earlier studies focusing on climate

198    sensitivity (e.g., Murphy et al., 2004; Stainforth et al., 2005; Sanderson et al., 2008a,b,

199    2010, 2011), to later ones attempting to generate plausible representations of the climate

200    without flux adjustments (e.g. Irvine et al., 2013; Yamazaki et al., 2013) and using history

201    matching to reduce parameter space uncertainty (Williamson et al., 2013). More recently,

202    Mulholland et al. (2016) demonstrated the potential of using PPEs to improve the skill of

203    initialized climate model forecasts of 1 month lead time, and Sparrow et al. (2018) showed

204    that large PPE can be used to identify subgrid scale parameter settings that are capable of

205    best simulating the ocean state over the recent past (1980-2010). However, very little has

206    been published on using PPEs for parameter refinement with the aim of improving regional

207    climate models (RCMs).

208

209    The goals of this study were to: 1) identify model parameters that most strongly control the

210    annual cycle of near-surface temperature and precipitation over the NWUS in

211    weather@home2, and 2) select model parameterizations that reduce the warm/dry summer

212    biases without introducing or unduly increasing other biases. We acknowledge that

213    changing a model in any way inevitably involves making sequences of choices that

214    influence the behaviour of the model. Some of the model behavioural changes are targeted

215    and desirable, but parameter refinement may have unintended negative consequences.

216    There is a general concern that 'improved' performance arises because of compensation

217    among model errors, and an 'accurate' climate simulation may very well be achieved by

218    compensating errors in different processes, rather than by best simulating every physical

219    process. This concern motivated us to select multiple parameter sets from the tuning

220    exercise rather than seek an "optimal" set. Though having multiple parameter sets does not

221    eliminate the problem, to the degree that each parameter set compensates for errors

222    uniquely, obtaining a similar model response to some change in forcing across parameter

223    sets may provide more confidence in that response.

224

225    It is worth noting that this study looks mainly at atmospheric parameters because we

226    intended to focus this study on larger-scale atmospherics dynamics that influence the

227    boundary conditions of the regional model, especially how much moisture and heat is

228    advected to the regional model, while local land surface/atmosphere interactions are being

Geoscientific
Model Development
Discussions

229    examined in a subsequent study that perturbs a suite of atmospheric and land surface

230    parameters in the regional model.

231

232    **2. Methodology**

233    Throughout this paper we use 'simulated' to refer to outputs from climate models, and

234    'emulated' results to refer to estimated/predicted outputs from statistical emulators.

235

236    **2.1. Overview of the parameter refinement process**

237    This study carried out an iterative parameter refinement exercise, or an 'iterative

238    refocusing' procedure to use a term coined in Williamson et al. (2017). The multi-

239    dimensional parameter space is reduced in phases, where each phase includes the following

240    steps:

241    1) Randomly sample the initially defined parameter space (defined by the bounds of the 17

242    parameters listed in Table1) to generate sets of parameter combinations;

243    2) generate a PPE with the parameters sets from step (1) through weather@home;

244    3) train statistical emulators for multiple climate metrics using the PPE from step (2);

245    4) reduce the parameter space (i.e., narrow the ranges of acceptable values for parameters)

246    such that the space excludes ensemble parameter sets that are 'too far away' from target

247    metrics;

248    5) randomly sample the reduced parameter space to design a new set of parameter

249    combinations;

250    6) use the trained emulators to filter the sample from step (5), and reject a parameter set if

251    the emulator prediction is too far away from a target value;

252  7) repeat steps (2) through (6) until the desired outcome is achieved.

253  Detailed descriptions of the parameter refinement process throughout three phases is

254  presented in Appendix A, including decisions on what key climate metrics to use in each

255  phase, and the stopping point of this iterative exercise - after three phases.

256

257  Here we briefly summarize the objective of each phase. The objective of Phase 1 was to

258  eliminate regions of parameter space that led to top-of-atmosphere (TOA) radiative fluxes

259  that are too far out of balance. The objective of Phase 2 was to reduce biases in the

260  simulated regional climate of NWUS, while not straying too far away from TOA radiative

261  (near-) balance. Lastly, the objective of Phase 3 was to further refine parameter space,

262  specifically to reduce the JJA warm and dry bias over the NWUS.

263

264  The principle climate metrics used to access the effect of parameter perturbation are: Phase

265  1) TOA radiative fluxes, where we considered outgoing (reflected) shortwave radiation

266  (SW) and outgoing longwave radiation (LW) separately; Phase 2) NWUS regional surface

267  metrics - the mean magnitude of the annual cycle of temperature (MAC-T), and mean

268  temperature (T) and precipitation (Pr) in December-January-February (DJF) and (JJA),

269  while still being mindful of SW and LW; and Phase 3) same as Phase 2, except for selecting

270  model parameterizations that reduce the JJA warm and dry biases over the NWUS.

271

272  **2.2. Climate simulations with weather@home**

273  The climate simulations used in this study were generated through the weather@home

274  climate modelling system (Massey et al., 2015; Mote et al., 2016) with updates (Guillod et

Geoscientific
Model Development
Discussions

275    al., 2017) that includes MOSES2. MOSES2 simulates the fluxes of $CO_2$, water, heat, and

276    momentum at the interface of the land and atmospheric boundary layer, and is capable of

277    representing a number of sub-grid tiles within each grid box, allowing a degree of sub-grid

278    heterogeneity in surface characteristics to be modeled (Williams et al., 2012).

279

280    The western North America application of weather@home (weather@home-WNA)

281    consists of HadRM3P (0.22° × 0.22°) nested within HadAM3P (1.875° longitude ×1.25°

282    latitude). Weather@home-WNA prior to recent enhancements was evaluated for how well

283    it reproduced various aspects of the recent historical climate of the western US by Li et al.

284    (2015), Mote et al. (2016), Rupp and Li (2016), and Rupp et al. (2017).  Notable warm/dry

285    biases in JJA were present over the NWUS and these biases persist with MOSES2 (Fig.

286    S1), with a temperature bias of 3.9 °C and a precipitation biases of -8.5 mm/month (-32%)

287    in JJA over Washington, Oregon, Idaho and western Montana, as compared with the

288    PRISM gridded observational dataset (Daly et al., 2008). Note these were biases using

289    default, i.e. standard physics (SP), model parameter values.

290

291    Each simulation in the PPE spanned 2 years, with the first year serving as spin-up and only

292    the second year used in the analysis. Simulations began on 1 December of each year for

293    the years 1995 to 2005, except for Phase 1 (see description of Phases in Appendix A).

294    Climate metrics were averaged over December 1996 to November 2007 (except Phase 1).

295    This time period was chosen because it contained a wide range of SST anomaly patterns -

296    including the very strong 1997-98 El Niño – which helps reduce the influence that any

297   particular SST anomaly pattern may have on the sensitivities of chosen climate metrics to

298   parameters.

299

300   **2.3. Perturbed parameters**

301   In our PPE, we initially selected 17 model parameters to perturb simultaneously, 16 in the

302   atmospheric model, and one in the land surface model (Table 1). The atmospheric

303   parameters are a subset of those perturbed in Murphy et al. (2004) and Yamazaki et al.

304   (2013); both studies also perturbed ocean parameters, and Yamazaki et al. (2013) perturbed

305   forcing parameters (e.g., scaling factor for emission from volcanic emissions) as well. Our

306   selection of parameters was constrained to those available to be perturbed using

307   weather@home at the time. Ranges for most parameter perturbations were 1/3 to 3 times

308   the default value, but for certain parameters (e.g., empirically adjusted cloud fraction,

309   EACF), only values greater than the default value were used (Table 1).  We intentionally

310   began with ranges generally wider than those used in previous studies (Murphy et al. 2004;

311   Yamazaki et al. 2013) because we intended to refine the ranges through multiple phases of

312   PPEs.

313

314   Though a principal objective was to evaluate sensitivity of the regional climate to

315   atmospheric parameters, sensitivities may be a function of land-atmosphere exchanges

316   (Sippel et al., 2016; Guillod et al., 2017).  While many parameters influence land-

317   atmosphere energy and water exchanges in MOSES2, one (V_CRIT_ALPHA) has been

318   shown to be particularly important (Booth et al., 2012) so was included in our tuning

319   exercise.   V_CRIT_ALPHA defines the soil water content below which transpiration

320   begins being limited by soil water availability and not solely the evaporative demand.

321

322   **2.4 Observational data**

323   The regional biases in MAC-T, JJA-T, JJA-Pr, DJF-T and DJF-Pr  - were all calculated

324   with respect to the 4-km resolution monthly PRISM dataset, after regridding the PRISM

325   data to the HadRM3P grid. To consider observational uncertainty, we also compared JJA-

326   T biases using four other observational datasets: 1) NCEP/NCAR Reanalysis 1 (NCEP,

327   Kalnay et al., 1996), 2) the Climate Forecast System Reanalysis and Reforecast (CFSR,

328   Saha et al., 2010), 3) the Modern-Era Retrospective Analysis for Research and

329   Applications Version2 (MERRA2, Gelaro et al., 2017), and 4) Climatic Research Unit

330   temperature dataset v4.00 (CRU, Harris et al., 2014).  The four datasets are not shown here

331   for the regional analysis because the maximum regionally averaged difference (0.71 °C)

332   among the datasets is less than 1/5 of  the regionally averaged JJA-T bias. Throughout this

333   paper, regional biases are calculated with respect to PRISM.

334

335   The biases in global temperature were calculated with respect to CRU, MERRA2, CSFR,

336   NCEP, and the Climate Prediction Centre global land surface temperature data; the latter

337   is a combination of the station observations collected from Global Historical Climatology

338   Network version 2 and the Climate Anomaly Monitoring System (GHCN-CAMS, Fan and

339   van den Dool, 2008).  The biases in global precipitation were calculated with respect to

340   CRU, MERRA2, CFSR, Global Precipitation Climatology Project monthly precipitation

341   (GPCP, Adler et al., 2003), Global Precipitation Climatology Centre monthly precipitation

342    (GPCC, Schneider et al., 2013), ERA-Interim reanalysis dataset (ERAI, Dee et al., 2011),

343    Japanese 55-year Reanalysis (JRA-55, Onogi et al., 2007) and NOAA-CIRES 20th Century

344    Reanalysis version 2c (20CRv2c, Compo et al.. 2011). All the datasets were regridded to

345    the HadAM3P grid before biases were calculated.

346

347    For all the observational datasets, data from December 1996 to November 2007 (the same

348    time period the model simulations cover as shown in Table2) was used to calculate model

349    biases.

350

351    **2.5 Sensitivity Analysis**

352    The response of the climate model to perturbations in the multidimensional parameter

353    space can be non-linear. In order to isolate the influence of each parameter on key climate

354    metrics and eliminate parameters that do not have a strong control on those metrics, we

355    performed two types of sensitivity analysis. One determines the sensitivity of a single

356    parameter by perturbing one parameter with all other parameters fixed, i.e. one-at-a-time

357    (OAAT) sensitivity analysis. Following Carslaw et al. (2013) and McNeall et al. (2016),

358    we also used a global sensitivity analysis using Fourier Amplitude sensitivity test (FAST)

359    for qualitative sensitivity analysis to validate the results of OAAT and to estimate

360    interactions among parameters. FAST allows the computation of the total contribution of

361    each input parameter to the output's variance, where total includes the factor's main effect,

362    as well as the interaction terms involving that input parameter. The computational aspects

363    and advantages of FAST are described in Satelli et al. (1999).

364

Geoscientific
Model Development
Discussions

## 3. Results and Discussion

Top-of-atmosphere (TOA) radiative balance is an emergent property in GCMs (Irvine et al., 2013), and the fact that the models of the IPCC Assessment Report 4 did not need flux-adjustment was seen as an improvement over earlier models (Solomon et al., 2007). Although climate models approximately balance the net absorption of solar radiation with the outward emission of longwave radiation (OLR) at the TOA, the details of how solar absorption and terrestrial emission are distributed in space and time depend on global atmospheric and oceanic circulation, clouds, ice, and other aspects of model behaviour. The surface expression of those global processes is also important given that a primary and practical purpose of climate modelling is to understand how (surface) climate will change. We describe the responses of both global TOA and regional surface climate to parameter refinement.

## 3.1. TOA radiative fluxes

In Fig. 1, we show the TOA energy flux components from the PPEs from each of the three phases. In Phase 1, many parameter sets (72%) resulted in TOA energy fluxes that vastly exceeded our ranges of acceptability (as defined in Appendix A). In Phase 2, most of the parameter sets resulted in TOA energy fluxes that fell within the ranges of acceptability; the 20% that did not reveal the error in our predictions using the emulator since the parameter sets were chosen to specifically achieve TOA fluxes within the region of acceptability. In Phase 3, nearly all (97%) the parameter sets yielded acceptable results. It is worth mentioning again that in Phase 3, selection of parameter sets was based only

Geoscientific
Model Development
Discussions

387  secondarily on TOA fluxes and primarily on regional climate metrics (see detailed

388  description of Phase 3 in Appendix A).

389

390  Rowlands et al. (2012) discarded any ensemble member that required a global annual mean

391  flux adjustment of absolute magnitude greater than 5 W m-2 (see red lines in Fig. 1) and

392  Yamazki et al. (2013) defined a confidence region of (SW, LW) that corresponded to a

393  TOA imbalance of less than $5\ W\ m^{-2}$ as one that did 'not drift significantly' from a realistic

394  TOA state.  Although the ranges of acceptability (Fig.1) permits net TOA imbalance

395  greater than $5\ W\ m^{-2}$, more than half (55.8%) of the Phase 3 parameter sets generated a

396  TOA imbalance less than $5\ W\ m^{-2}$, and the smallest TOA imbalance was less than 0.1 W

397  $m^{-2}$.

398

399  The entrainment coefficient (ENTCOEF) and the ice fall speed (VF1) were the dominant

400  controls on the TOA outgoing SW and LW fluxes, respectively (see SW and LW response

401  to these two parameters shown in the bottom two rows of Fig. S2).  Why these parameters

402  are important becomes clear from understanding their respective roles in the climate model,

403  especially with respect to convection and hydrometeor transport.

404

405  The atmospheric model simulates a statistical ensemble of air plumes inside each

406  convectively unstable grid cell. On each model layer, a proportion of rising air is allowed

407  to mix with surrounding air and vice-versa, representing the process of turbulent

408  entrainment of air into convection and detrainment of air out of the convective plumes

409  (Gregory and Rowntree, 1990). The rate at which these processes occur in the model is

Geoscientific
Model Development
Discussions

410     proportional to ENTCOEF, which is a parameter in the model convection component

411     (Table1). The implication of perturbing ENTCOEF has been investigated by (Sanderson et

412     al, 2008b) using single perturbation experiments, and they showed that a low ENTCOEF

413     leads to a drier middle troposphere and moister upper troposphere. Conversely, increasing

414     ENTCOEF results in increased low level moisture (more low level clouds) and decreased

415     high level moisture (less high level clouds). Because the albedo effects of low clouds

416     dominate their effects on emitted thermal radiation (Hartmann et al., 1992; Stephens,

417     2005), increasing ENTCOEF increases the outgoing SW fluxes.

418

419     VF1 is the speed at which ice particles may fall in clouds. A larger ice fall speed is

420     associated with larger particle sizes and increased precipitation. Wu (2002) studied ice fall

421     speed parameterization in radiative convective equilibrium models, and found that a

422     smaller ice fall speed leads to a warmer, moister atmosphere, more cloudiness, weak

423     convection and less precipitation, which could lead decreased outgoing LW TOA flux due

424     to absorption in the cloud itself and/or in the moist air. Higher ice fall speeds produce the

425     opposite - a cooler, clearer, less cloudiness, strong convection and more precipitation,

426     which increases the outgoing LW flux.

427

428     **3.2. Regional climate improvements**

429     A primary and practical purpose of climate modelling is to understand how (surface)

430     climate will change, but model biases can have non-negligible impacts on projections. In

431     Phase 2 and 3 we evaluate the response of regional surface climate to parameter

Geoscientific
Model Development
Discussions

Open Access

EGU

432    perturbations, and refine the parameter space to reduce biases in regional temperature and

433    precipitation.

434

435    In Phase 2, we identified ENTCOEF and VF1 as distinct from the other 15 parameters with

436    respect to their influence on the overall suite of climate metrics to a first order

437    approximation (Fig. S3). Recall the regional surface metrics considered were MAC-T, JJA-

438    T, JJA-Pr, DJF-T, and DJF-Pr. Though MAC-T is our principal metric (section2.1), MAC-

439    T co-varies with JJA-T, JJA-Pr, and DJF-T (Fig. S3), so moving in parameter space toward

440    lower bias in MAC-T reduces biases in JJA-T, JJA-Pr, and DJF-T. MAC-T does not co-

441    vary strongly with DJF-Pr.

442

443    Each OAAT relationship in Fig. 2 depends on the initial ranges of the input parameters

444    from the ensemble design, and is computed while holding all other parameters at their

445    ensemble mean values.   Because sensitivity can change as one moves through the

446    parameter space (e.g. CW_LAND and ENTCOEF in Fig. 2), these relationships must be

447    interpreted with care. Within the refined parameter space in Phase 2, ENTCOEF and the

448    parameter that limits photosynthesis (and thereby latent heat flux via transpiration) as a

449    function of soil water (V_CRIT_ALPHA) were the most influential individual parameters

450    and counter each other when both increased (Fig. 2 and Fig. S3).  The parameter that

451    controls the cloud droplet to rain threshold over land (CW_LAND) also had strong

452    influence on MAC-T across the lower end of the parameter perturbation range (up to

453    0.004). The other parameters had little to effectively no influence on MAC-T. The results

Geoscientific
Model Development
Discussions

454     of OAAT sensitivity analysis for the other output metrics considered in Phase 2 are

455     presented in Fig. S6-S11.

456

457     The global sensitivities of the simulated outputs (the ones considered in Phase 2) due to

458     each input, as both a main effect and total effect, including interaction terms, are presented

459     in Fig. 3. ENTCOEF was the most important parameter for all three surface temperature

460     metrics, with a total sensitivity index of ~0.7, 0.5, and 0.4 for MAC-T, JJA-T, and DJF-T

461     respectively , where maximum sensitivity is 1 (see Satelli et al. 1999). For the metrics

462     MAC-T and JJA-T, V_CRIT_ALPHA was the next most important, with a total sensitivity

463     index of ~0.3 for both metrics. For JJA-Pr, the most important parameter was VF1,

464     followed by ENTCOEF; for DJF-Pr, the most important parameter was ENTCOEF, closely

465     followed by the parameter that controls the roughness length for free heat and moisture

466     transport over the sea (Z0FSEA).

467

468     The interaction terms were relatively small, accounting for a few percent of the variance,

469     except for the effect of ENTCOEF on DJF-Pr, where the interaction with other parameters

470     accounts for ~ 1/3 of the variance.  In a study constraining carbon cycle parameters by

471     comparing emulator output with forest observations, McNeall et al. (2016) also found the

472     importance of the interaction terms negligible.  In contrast, Bellprat et al. (2012b) used

473     quadratic emulator to objectively calibrate a regional climate model, and found non-

474     negligible interaction terms. They showed that excluding the interactions in the emulator

475     increased the error of the emulated temperature and precipitation results by almost 20%.

476    Further work could be done to assess the magnitude and functional form (i.e. linear or

477    nonlinear) of the interaction terms, but is beyond the scope this study.

478

479    Only the parameters with a total sensitivity index larger than ~0.1 for MAC-T, JJA-T, DJF-

480    T, JJA-Pr, or DJF-Pr were retained for perturbation in Phase 3: CW_LAND, VF1,

481    ENTCOEFF, V_CRIT_ALPHA, ASYM_LAMBDA, G0, and Z0FSEA. Although the

482    parameter that controls the rate at which cloud liquid water is converted to precipitation

483    (CT) had a total sensitivity index of ~0.1 for SW, it was excluded from further perturbation

484    because the primary interest in Phase 2 was in regional surface metrics, not TOA radiative

485    fluxes.

486

487    Phase 3 demonstrated the power of our approach for reducing regional mean biases in

488    MAC-T, JJA-T and JJA-Pr. Simulations from Phase 3 resulted in MAC-T biases 1- 3°C

489    lower than SP (Fig.4 middle row). All Phase 3 parameter sets improved the JJA-Pr dry bias

490    with several eliminating the bias entirely. Many parameter sets reduced the bias in JJA-T

491    to less than 1.5°C, a dramatic improvement (~63%) over the 4°C SP bias. However, these

492    improvements come at a small price, namely a larger regional (NWUS) dry bias in DJF-Pr

493    (about -15% compared with PRISM in the worst case). Because our primary goal was to

494    reduce JJA warm and dry biases, any model variant from Phase 3 is preferable to SP. Any

495    subset of parameterizations from phase 3 can now be used in subsequent experiments.

496

497    V_CRIT_ALPHA plays an important role in controlling JJA-T and MAC-T (as shown in

498    Fig. 2 and Fig. S6) due to its role in the surface hydrological budget. V_CRIT_ALPHA

499   defines the critical point as a fraction of the difference between the wilting soil water

500   content and the saturated soil water content (as described in Appendix C). The critical

501   point is the soil moisture content below which plant photosynthesis becomes limited by

502   soil water availability. When V_CRIT_ALPHA is zero, transpiration starts to be limited as

503   soon as the soil is not completely saturated, whereas when it is one, transpiration continues

504   unlimited until soil moisture reaches wilting point at which point transpiration switches

505   off. Lower values of V_CRIT_ALPHA reduce the critical point allowing plant

506   photosynthesis to continue unabated at lower soil moisture levels, i.e. plants are not water-

507   limited. As plants photosynthesize water is extracted from soil layers and transpired,

508   increasing the local atmospheric humidity and lowering the local temperature through

509   latent cooling. Our results are consistent with previous findings by Seneviratne et al.

510   (2006), who also show reducing the temperature and increasing humidity can feedback

511   onto the regional temperature and precipitation during the summer months.

512

513   The only apparent constraints on ranges of parameter values through three phases of

514   parameter refinement were seen for V_CRIT_ALPHA and ENTCOEF. Values of

515   V_CRIT_ALPHA lower than 0.7 were required to keep the bias of MAC-T under 3 °C.

516   For ENTCOEF, the range between 3 and 5 contains the best candidates to reduce regional

517   warm/dry biases. The range of ENTCOEF identified here is consistent with findings of

518   Irvine et al. (2013), which also show that low values of ENTCOEF tend to give warmer

519   conditions. However, results from other previous studies varies. Williamson et al. (2015)

520   found that low values of ENTCOEF are implausible, and that there are more plausible

521   model variants at the upper end of its perturbed range, whereas Sexton et al. (2011) and

522 Rowlands et al. (2012) consider the range between 2 and 4 to contain the best model

523 variants. The discrepancy in optimal ranges for ENTCOEF are to be expected given that

524 the primary metrics used to evaluate the effect of parameter refinement are different, with

525 ours being JJA warm/dry biases over the NWUS, William et al. (2015) being the behaviour

526 of Antarctic Circumpolar Current, and other previous studies being climate sensitivities.

527 This demonstrates that any parameter refinement process is tailored to a specific objective,

528 and choices regarding metrics (e.g., variables, validation dataset(s), and / or cost functions)

529 may determine which part of parameter space is ultimately accepted.

530

531 **3.3. Effects on global scale climate**

532 To avoid introducing or increasing biases over other parts of the globe by our regionally-

533 focused model improvement effort,  we investigated the large-scale effects of the selected

534 10 'good' (least biased in MAC-T) sets of global parameter values. We focused on surface

535 temperature and precipitation because they are key variables of the climate system and are

536 of high interest for impact studies.

537

538 Figure 5 shows the meridional distribution of Northern Hemisphere (NH) mid-latitude

539 temperature (over land) and precipitation in DJF and JJA. Because of the wide range of

540 parameter values in the PPEs of Phase 1 and Phase 2, the spread for these PPEs is quite

541 large, whereas the ensemble spread in Phase 3 is substantially smaller. Compared with the

542 SP ensemble, the new parameter values (final 10 sets) reduced the zonal mean JJA

543 temperature throughout the NH mid-latitudes (30 °N -60 °N), by ~1 °C – 4 °C (depending

544 on the particular combination of parameters), and increased JJA precipitation over the same

545    latitude bands, except for latitudes south of 33 °N and north of 58 °N. In DJF, the effects

546    are not as large nor are the changes consistent in sign across the NH mid-latitude region

547    (though south of ~38 °N all 10 parameter sets give increasing precipitation).

548

549    To examine how parameter refinements affect spatial patterns of biases, we compare the

550    seasonal mean biases of temperature (Fig. 6) and precipitation (Fig. 7) under SP and the

551    selected PP settings, against CRU data. The SP simulations have large warm biases in JJA

552    (and to a lesser extent in MAM and SON, Fig. 6 b-d) over the NH mid-latitude land region,

553    that are substantially lower in the PP simulations (Fig. 6 f-h and Fig.6 j-l).  In the tropics,

554    the SP simulations have cold biases over northern South America, central Africa and

555    southern Asia in most seasons that are ameliorated in the PP simulations in some cases

556    (e.g. central Africa in DJF and SON) - even though the focus of the PP simulations was

557    improving the climate of the NWUS. The SP simulations also have cold biases over most

558    of the Southern Hemisphere continents in mid-latitudes in most seasons. A large fraction

559    of the JJA temperature biases were reduced in the PP simulations, as shown in Fig. 6c, g

560    and k. These salient features in JJA temperature biases under SP and PP are not particular

561    to the selection of observational dataset (see Fig. S12-S15 for comparison with other

562    datasets). In the other three seasons, however, the spatial patterns of temperature biases are

563    not consistent across observational datasets.

564

565    The reduction of JJA temperature from SP to PP (Fig. 6k) and the resulting reduction in

566    bias are accompanied by reduction in precipitation in the equatorial regions; increased

567    precipitation over northern North America, northern Africa, and Europe (Fig. 7k); and

Geoscientific
Model Development
Discussions

568   decreased incoming shortwave radiation at the surface and increased evaporation (Fig.

569   S16). Stronger evaporative cooling and reduced surface radiation lead to a cooling of the

570   JJA climate, which roughly agrees with the geographical pattern of reduced mean JJA

571   temperature, consistent with findings in Zhang et al. (2018) that both overestimated surface

572   shortwave radiation and underestimated evaporation contribute to the warm biases in JJA

573   in CMIP5 climate models.

574

575   For precipitation, the largest biases in SP are over Amazonia in DJF and MAM (Fig. 7a

576   and b), and northern South America, equatorial Africa, and south Asia in JJA (Fig. 7c).

577   These summer biases are increased in the PP simulations (Fig. 7k). However, it is difficult

578   to know whether we are improving the model's global precipitation patterns because of the

579   large uncertainty in historical precipitation observational datasets. Still, it is worth

580   comparing the PP simulations with both a variety of observational-based datasets and other

581   GCMs (Fig. 8). The precipitation amounts differ substantially across different

582   observational datasets, as well as across climate models. In the tropics, Phase 3 PP

583   simulated precipitation is mostly lower (except DJF just north of the equator) and has

584   narrower range than the observations or other climate models, but is higher in DJF and JJA

585   (up to 25% higher) than the SP simulation results. Outside the tropics, the precipitation

586   distributions in PP remain similar to those of SP, and differences from observational

587   datasets and other GCMs are less affected by the use of PP. The tropical precipitation

588   improvements in JJA can be taken as a general improvement, though not with high

589   confidence due to the variability across observational datasets. To further highlight the

590   uncertainties in precipitation, global maps of differences in biases between SP and our

591     selected parameter settings, in comparison with other observational-based datasets, are

592     presented in Fig. S17-23.

593

594     The fact that the large JJA warm bias (shared with many other GCMs and RCMs; see e.g.

595     Mearns et al., 2012; Kotlarski et al., 2014) could be reduced substantially through the use

596     of PP is a notable result, especially since the bias persisted through initial tuning efforts

597     and through the recent updates from version 1 to version 2 of weather@home. We

598     demonstrated here that significant improvements in the simulation of JJA temperature can

599     be made through parameter refinements, and that these JJA temperature biases are not

600     necessarily structural issues of the climate model. These improvements in simulating JJA

601     temperature generally did not overall improve JJA precipitation patterns across the globe,

602     and even worsened the bias in some places (e.g. South America).

603

604     **4. Conclusions**

605     Through an iterative parameter refinement approach to improve model performance, we

606     identified a region of climate model parameter space in which HadAM3P outperforms the

607     SP variant in simulating summer climate over the NWUS specifically, and over NH mid-

608     latitude land in general, while approximately maintaining TOA radiative (near-) balance.

609     Improving the northwest US climate comes with tradeoffs, e.g. larger JJA dry bias over

610     Amazonia. However, it is important to note that there are large uncertainties in observed

611     precipitation climatology, especially outside of the North American and European mid-

612     latitudes, so both apparent increases and decreases in biases should be treated with caution,

613     and compared against the range across observational datasets. In the end, we consider the

614    cost of increasing biases in parts of the globe acceptable for the purposes of selecting

615    multiple global model variants to drive the regional model with reduced JJA biases over

616    NWUS. The fact that improvements can be made at all (for a substantial area of the world)

617    through targeted PPE is encouraging.

618

619    Our parameter refinement yielded important improvements in the representation of the

620    summer climate over the NWUS, and it follows that biases in other models may also be

621    reduced by refining certain parameters that, although may not be identical to those in

622    HadAM3/RM3P, influence the same physical processes similarly. We found ENTCOEF

623    and V_CRIT_ALPHA to be the dominant parameters in reducing JJA biases. These

624    parameters control cloud formation and latent heat flux, respectively.  Bellprat et al. (2016)

625    found the key parameter responsible for reduction of JJA biases is increased hydraulic

626    conductivity, which increases the water availability at the land surface and leads to

627    increased evaporative cooling, stronger low cloud formation, and associated reduced

628    incoming shortwave radiation. We only perturbed one land surface parameter, but the

629    effects of additional land surface parameters are being explored in a subsequent study.

630    Given that land model parameters such as V_CRIT_ALPHA could reasonably be expected

631    to interact with sensitive atmospheric parameters like ENTCOEF, it is particularly

632    interesting to consider the multivariate sensitivity of a range of parameters that span across

633    component models (e.g., land, ice, atmosphere, ocean). We argue that this frontier of

634    parameter sensitivity exploration should be done in a transparent and systematic manner,

635    and we have demonstrated that statistical emulators can be effectively leveraged to reduce

636    computational expense.

Geoscientific
Model Development
Discussions

637

638    The fact that V_CRIT_ALPHA (which is a parameter in the land surface scheme MOSES2)

639    was found to be an important parameter on regional MAT-C and JJA-T, has much further

640    implications beyond this study. MOSES2 is the land surface scheme used in HadGEM1

641    and HadGEM2 family, which were used in CMIP4 and CMIP5. Moreover, the Joint UK

642    Land Environment Simulator (JULES) model (which is the land surface scheme of the

643    CMIP6 generation Hadley Centre models HadGEM3 family, https://www.wcrp-

644    climate.org/wgcm-cmip/wgcm-cmip6) is a development of MOSES2. What we have

645    learned about the atmosphere-land surface interactions here is relevant to even the most

646    recent HadGEM model generation and the in-progress CMIP6.

647

648    The reduction of JJA biases that we achieved in our multi-phase parameter refinement is

649    notable. However, despite out efforts, the 'best' performing parameter set still simulates a

650    MAC-T bias of 1.5 °C, and a JJA-T bias of 1 °C, over the NWUS. Future work could be

651    done to determine whether the model can be further improved by tuning additional land-

652    surface scheme parameters, and/or to what extent the remaining biases are due to structural

653    errors of the model for which we cannot (nor even should not) compensate by refining

654    parameter values. However, with the reduction in JJA temperature bias, future projections

655    using the new parameter settings over the SP should be at less risk of overestimating

656    projected warming in summer (as discussed in the introduction).

657

658    It is also worth noting that we restricted our analysis to seasonal and annual mean climate

659    metrics. Given the use of weather@home for attribution studies of many extreme weather

660    events (e.g., Otto et al., 2012; Rupp et al., 2017a) as well as their impacts, such as flooding-

661    related property damages (Schaller et al., 2016) and heat-related mortality (Mitchell et al.,

662    2016), an important next step would be to investigate how the tails of distributions of

663    weather variables respond to parameter perturbations.

664

665    Another important next step would be to apply the selected PPE over the weather@home

666    - European domain, given the non-trivial JJA warm bias identified over Europe by previous

667    studies (Massey et al., 2014; Sippel et al., 2016; Guillod et al., 2017). Bellprat et al. (2016)

668    showed that regional parameters tuned over Europe domain also produced similar

669    promising results over North America domain but the same model parameterization yielded

670    larger overall biases over North America than for Europe. One could test the transferability

671    of parameter values over different regional domains in the weather@home framework,

672    given weather@home currently uses the same GCM to drive several RCMs over different

673    parts of the world, all using the  same parameter values.

674

675    The methodology presented in this study could be applied to other models in the evolution

676    of physical parameterizations, and we advocate that parameter refinement process should

677    be more explicit and transparent as done here. Choices and compromises made during the

678    refinement process may significantly affect model results and influence evaluations against

679    observed climate, hence should be taken into account in any interpretation of model results,

680    especially in intercomparison of multimodel analyses to help understanding of model

681    differences.

682

Geoscientific
Model Development
Discussions

**Code availability**

HadRM3P is available from the UK Met Office as part of the Providing REgional

Climates for Impacts Studies (PRECIS) program. Access to the source code is dependent

on attendance at a PRECIS training workshop

(http://www.metoffice.gov.uk/research/applied/international-development/precis/obtain).

The code to embed the Met Office models within weather@home is proprietary and not

within the scope of this publication.

**Data availability**

The model output data for the experiment used in this study will be freely available at the

Centre for Environmental Data Analysis (http://www.ceda.ac.uk) in the next few months.

Until the point of publication within the CEDA archive, please contact the corresponding

author to access the relevant data.

**Appendix A: Detailed experimental process**

The overarching goal is to refine parameter values to reduce warm and dry summer bias in

the NWUS. In total four ensembles were generated, one using the SP values and one for

each of 3 PPE phases.  Details of each ensemble are listed in Table 2.

Internal variability of the atmospheric circulation can confound the relationship between

parameters values and the response being sought (i.e. result in a low signal-to-noise ratio).

Averaging over multiple ensemble members with the same parameter values but different

atmospheric initial conditions (ICs) can clarify the true sensitivity to parameters by

706  increasing the signal-to-noise ratio. We set up multiple ICs for each parameter set, but the

707  numbers of ICs applied was not consistent throughout the experiment. The IC applied in

708  each phase was determined somewhat subjectively, trying to strike a balance between

709  running a large enough PPE to probe as many processes and interactions between

710  parameters as possible, having multiple ICs so that the results were representative of the

711  parameter perturbations instead of reflecting the influence of any particular IC, while under

712  the practical limitation of data transfer, storage, and analysis. The actual IC ensemble size

713  used in the final analysis was also constrained by the number of successfully completed

714  returns from the distributed computing network.

715

716  The four ensembles are summarized below:

717  **SP:** A preliminary "standard physics" (SP) ensemble with 10 ICs that used only the default

718  model parameters was generated to provide a benchmark to access the effects of parameter

719  perturbations.

720

721  **Phase 1:** The objective of this phase was to eliminate regions of parameter space that led

722  to top-of-atmosphere (TOA) radiative fluxes that are strongly out of balance. Exclusion

723  criteria were deliberately lenient, to avoid eliminating regions of the parameter space that

724  could potentially reproduce the observed temperature and precipitation over the western

725  US. We perturbed 17 parameters simultaneously, using space-filling Latin hypercube

726  sampling (McKay et el., 1979) - maximizing the minimum distance between points - to

727  generate 340 sets of parameterizations across the range of parameter values described in

728  Table 1. To generate enough ensemble members for a statistical emulator, Loeppky et al.

729    (2009) suggested that the number of sets of parameter values be 10 times the number of

730    parameters ($p$). We used more than $10p$ sets of parameter values in this, and subsequent

731    phases of PPE.  A total of 2040 simulations (340 sets of parameter values x 6 ICs) were

732    submitted to the volunteer computing network.  This phase was considered finalized when

733    simulations with 220 sets of parameter values and 3 IC ensemble members per set were

734    returned from the computing network.

735

736    Model results were used to train a statistical emulator which maps the relationship between

737    parameter values and key climate metrics. In this phase, the metrics were outgoing LW and

738    (reflected) SW TOA radiative fluxes. We considered these two metrics separately because

739    the total net radiation could mask deficiencies in both types of radiation through

740    cancellation of errors.

741

742    For the emulator, a 2-layer feed-forward Artificial Neural Network (ANN, Knutti et al.,

743    2003; Sanderson et al., 2008; Mulholland et al., 2016) was used. Although other machine-

744    learning algorithms could be suitable (Rougier et al., 2009; Neelin et al., 2010; Bellprat et

745    al., 2012a,b, 2016), we chose ANN because it permits multiple simultaneous emulator

746    targets (i.e., TOA SW and LW at the same time). We used an ellipse (Fig. 1) to define the

747    space of acceptability for SW and LW, starting with the observational uncertainty ranges

748    given in Stephens et al. (2012), but tripling them (deliberately setting a lenient elimination

749    criteria), and then expanding both the negative and positive thresholds by an additional 1

750    W m$^{-2}$ to account for internal variability as estimated from SP (Fig. S5).  Sets of parameter

Geoscientific
Model Development
Discussions

751    values that fall within our range of acceptability were retained, and the ranges of these

752    refined/restricted parameter values defined the remaining  parameter space.

753

754    A new set of 1,000 parameter configurations was generated from the remaining parameter

755    space using space-filling Latin hypercube sampling. With this new ensemble we increased

756    the sample density within the refined parameter space. The statistical emulator was used to

757    predict SW and LW for each of these 1,000 new sets of parameters, and 41% fell within

758    our range of acceptability, reflecting the deficiency of the emulator to some extent.

759    Parameter sets that fell within the acceptable range were used in Phase 2.

760

761    **Phase 2:** The objective of this phase was to reduce biases in the simulated climate of the

762    NWUS, where the warm summer biases were the most obvious (Fig. S1), while not straying

763    far from TOA radiative (near-) balance. The climate metrics considered were the mean

764    magnitude of the annual cycle of temperature (MAC-T), and mean temperature (T) and

765    precipitation (Pr) in December-January-February (DJF) and June-July-August (JJA).

766    Although a primary motivation for this study was to investigate and reduce the warm and

767    dry bias in JJA over NWUS, MAC-T was treated as the primary metric in Phase 2 because

768    it is a comprehensive measure of climate feedbacks in response to a large change in forcing,

769    e.g., solar SW (Hall and Qu 2006).  MAC-T is also strongly correlated to the other regional

770    metrics (particularly JJA-T) as evident in Fig. S3 – MAC-T against other metrics. We chose

771    a NWUS average MAC-T of +/-3 °C as the bias threshold over which parameter space

772    would be eliminated.  Though this threshold is arbitrary, falling below it would mean

773    reducing the MAC-T bias for the NWUS by about 50%.

774

775    We did not treat all metrics as equally important. The order of importance in this second

776    phase was MAC-T > JJA-T, JJA-Pr, DJF-T, and DJF-Pr > SW and LW.

777

778    The 410 sets of new PPE from Phase 1 became the starting point for Phase 2. A total of

779    27,060 simulations (410 sets of parameter values x 6 ICs x 11 years) was submitted to the

780    computing network. This phase was considered finalized when simulations with 170 sets

781    of parameter values and 3 IC ensemble members per set and per year were completed.

782    These 5,610 simulations were used to train a suite of statistical emulators for various

783    climate metrics. An additional 94 sets of parameters with 3 IC ensemble members per set

784    and per year completed after starting Phase 3 and were used to validate the emulators

785    trained within Phase 2 (see Appendix B).

786

787    Separate statistical emulators were trained for MAC-T, JJA-T, JJA-Pr, DJF-T, DJF-Pr,

788    SW, and LW. Although ANN has the advantage of using multiple metrics as targets

789    simultaneously, the underlying emulator structure remains obscure, because an ANN is a

790    network of simple elements called neutrons which are organized in multilayer, and

791    different layers may perform different kinds of transformations on the inputs. For the sake

792    of simplicity and transparency,  in Phase 2 we used kriging instead - which is similar to a

793    Gaussian process regression emulator -  following McNeall et al. (2016) as coded in the

794    package DiceKriging (Roustant et al., 2012) in the statistical programming environment R.

795    We used universal kriging, with no 'nugget' term, meaning that the uncertainty on model

796    outputs shrinks to zero at the parameter input points that have already been run through our

797    climate model (Roustant et al., 2012). To validate if the emulators were adequate to predict

798    outputs at unseen parameter inputs, we needed to assure that it predicted relatively well

799    across our designed parameter inputs. For each emulator, we performed 'leave-one-out'

800    cross validation. The cross validation results showed no significant deviations in prediction

801    of the outputs (results not shown).

802

803    In addition to reducing parameter space in Phase 2, we also looked for parameters that

804    consistently showed little influence on our metrics of interest, as any reduction in

805    parameters could benefit subsequent experiments by reducing the overall dimensionality.

806    To identify which parameters have the most influence over the metrics of interest, we

807    performed two types of sensitivity analyses as described in Section 2.5. In the end, the 7

808    most influential parameters were retained after parameter reduction in Phase 2; these are

809    the bold-faced parameters in Table 1.

810

811    After eliminating parameter space resulting in MAC-T biases larger than 3°C, and reducing

812    the number of perturbed parameters to 7, we continued the parameter refinement process,

813    and randomly selected 100 parameter sets that emulated MAC-T biases less than 3°C and

814    had large spread in ENTCOEF and VIF1 (within the refined ranges of Phase 2). 100 was

815    subjectively chosen as a cut off number of new PPE sets to run through weather@home in

816    the next phase, mainly due to concern of not knowing how many more phases would be

817    required to reach our goal, while recognizing the practical constraints posed by the large

818    datasets that would potentially be generated in the following phases.

819

820    **Phase 3:** This objective of this phase was to further refine parameter space to reach the

821    target of northwest US regional bias in MAC-T less than 3°C, and then select 10 sets of

822    parameter values that met this criterion. The results in this phase satisfied our target, so we

823    stopped the iterative process here.

824

825    We were aware that our approach of regionally targeted parameter refinements might

826    degrade model performance elsewhere. Upon achieving our regional target, we

827    investigated the effects of our model tuning on global model metrics.

828

829    **Appendix B: Emulated vs. simulated results**

830    We used 94 additional ensemble members returned from Phase 2 (the 94 simulations that

831    completed after building the emulators from the Phase 2 PPE and starting Phase 3) to

832    provide out-of-sample validations of the emulators trained in Phase 2. In Fig. B1, we show

833    predictions from emulators against model-simulated values for all the output metrics. In all

834    cases, the linear relationship between the emulated and simulated is very strong (regression

835    coefficient regcoef>0.9), while the emulated results can predict the simulated results

836    relative well, with coefficient of determination $R2 > 0.9$ in the best cases (SW, LW and

837    JJA-T). It is not surprising that R2 for DJF-Pr is the smallest, considering precipitation in

838    DJF over NWUS is dominated by larger-scale atmospheric features such as the polar jet

839    stream, the Pacific subtropical high, and storm tracks (e.g.,Mock, 1996; Neelin et al., 2013;

840    Seager et al., 2014; Langenbrunner et al., 2015), and the internal variability of this metric

841    is the highest among those considered.

842

843    In Fig. B2, we present the emulated vs. simulated results in Phase 3 for the 95 PP sets that

844    were returned in Phase3. These 95 PP sets were run through the emulators from Phase 2 to

845    predict the climate metrics, then the emulated results were compared with the simulated

846    results returned from weather@home simulations. In most cases, r and R2 are lower than

847    the Phase 2 results (Fig. B1), except for LW and DJF-T, where R2 increases by a few

848    percent. This decrease in emulator prediction accuracy could be due to the fact that in Phase

849    3, only 7 parameters were perturbed simultaneously while keeping the rest at their default

850    values, so we have eliminated parts of the parameter space, which are no longer available

851    to the emulators.

852

853    The comparisons between simulated and emulated results from Phase 2 to Phase 3 highlight

854    the necessity of doing parameter refinement exercise in phases. Training a statistical

855    emulator once, then using it to search for optimal parameter settings may not always yield

856    optimum results.  An emulator may not fully capture the behaviour of the climate model in

857    every aspect, especially when the number of parameters perturbed was changed during the

858    process, such as in our case.

859

860    **Appendix C: Soil moisture control on plant photosynthesis in MOSES**

861    The critical point $\theta_{crit}$ (m$^3$ of water per m$^3$ of soil) is the soil moisture content below which

862    plant photosynthesis becomes limited by soil water availability and is calculated by:

863                    $$\theta_{crit} = \theta_{wilt} + V\_CRIT\_ALPHA\ (\theta_{sat}-\theta_{wilt})$$

864    where $\theta_{sat}$ is the saturation point, i.e. the soil moisture content at the point of saturation;

865    and $\theta_{wilt}$ is the wilting point, below which leaf stomata close. V_CRIT_ALPHA varies

866    between zero and one, meaning that $\theta_{crit}$ varies between $\theta_{wilt}$ and $\theta_{sat}$ (Cox et al., 1999).

867

868    **Author contributions**

869    The model simulations were designed by S. Li ,D. E. Rupp, L. Hawkins, with inputs from

870    P. W. Mote, and D. McNeall. All the results were analysed and plotted by S. Li. The paper

871    was written by S. Li, with edits from all co-authors.

872

873    **Competing interests**

874    The authors declare that they have no conflict of interest.

875

876    **Acknowledgements**

883

884

885

886

887     **Reference:**

888     Adler, R.F., Huffman, G.J., Chang, A., Ferraro, R., Xie, P.P., Janowiak, J., Rudolf, B.,

889     Schneider, U., Curtis, S., Bolvin, D., Gruber, A., Susskind, J., Arkin, P., and Nelkin, E.:

890     The Version 2 Global Precipitation Climatology Project (GPCP) Monthly Precipitation

891     Analysis (1979-Present), J. Hydrometeor., 4,1147-1167, https://doi.org/10.1175/1525-

892     7541(2003)004<1147:TVGPCP>2.0.CO;2, 2003.

893     Allen, M.: Do-it-yourself climate prediction, Nature, 401, 642, doi:10.1038/44266, 1999.

894     Bellprat, O., Kotlarski, S., Lüthi, D., and Schär, C.: Exploring perturbed physics ensembles

895     in a regional climate model, Journal of Climate, 25(13), 4582-4599,

896     https://doi.org/10.1175/JCLI-D-11-00275.1, 2002a.

897     Bellprat, O., Kotlarski, S., Lüthi, D., and Schär, C.: Objective calibration of regional

898     climate models, Journal of Geophysical Research: Atmospheres, 117(D23),

899     https://doi.org/10.1029/2012JD018262, 2012b.

900     Bellprat, O., Kotlarski, S., Lüthi, D., De Elía, R., Frigon, A., Laprise, R., and Schär, C.:

901     Objective calibration of regional climate models: application over Europe and North

902     America, Journal of Climate, 29(2), 819-838, https://doi.org/10.1175/JCLI-D-15-0302.1,

903     2016.

904     Boberg, F. and Christensen, J. H.: Overestimation of Mediterranean summer temperature

905     projections due to model deficiencies, Nat. Climate Change, 2, 433–436, doi:10.1038/

906     nclimate1454, 2012.

907     Booth, B. B. B., Jones, C. D., Collins, M., Totterdell, I. J., Cox, P. M., Sitch, S.,

908     Huntingford, C., Betts, R. A., Harris, G. R., and Lloyd, J.: High sensitivity of future global

909    warming to land carbon cycle processes, Environ. Res. Lett., 7, 024002, doi:10.1088/1748-

910    9326/7/2/024002, 2012.

911    Brown, T. J., Hall, B. L., and Westerling, A. L.: The impact of twenty-first century climate

912    change on wildland fire danger in the western United States: an applications perspective,

913    Climatic change, *62*(1-3), 365-388, 2004.

914    Carslaw, K. S., Lee, L. A., Reddington, C. L., Pringle, K. J., Rap, A., Forster, P. M., Mann,

915    G. W. , Spracklen, D. V. , Woodhouse, M. T. ,  Regayre, L. A., and Pierce, J. R.: Large

916    contribution of natural aerosols to uncertainty in indirect forcing, Nature, 503(7474), 67,

917    2013.

918    Cheruy, F., Dufresne, J. L., Hourdin, F., and Ducharne, A. :Role of clouds and land-

919    atmosphere coupling in midlatitude continental summer warm biases and climate change

920    amplification in CMIP5 simulations, Geophys. Res. Lett., 41, 6493–6500,

921    doi:10.1002/2014GL061145, 2014.

922    Collins, W.J., Bellouin, N., Doutriaux-Boucher, M., Gedney, N., Hinton, T., Jones, C.D.,

923    Liddicoat, S., Martin, G., O'Connor, F., Rae, J., Senior, C., Totterdell, I., Woodward, S.,

924    Reichler, T., and Kim, J.: Evaluation of the HadGEM2 model, Hadley Cent. Tech. Note,

925    74, 2008.

926    Collins, M., Booth, B. B., Bhaskaran, B., Harris, G. R., Murphy, J. M., Sexton, D. M., and

927    Webb, M. J.: Climate model errors, feedbacks and forcings: a comparison of perturbed

928    physics and multi-model ensembles, Climate Dynamics, 36(9-10), 1737-1766, 2011.

929    Compo, G.P., Whitaker, J.S., Sardeshmukh, P.D., Matsui, N., Allan, R.J., Yin, X., Gleason,

930    B.E., Vose, R.S., Rutledge, G., Bessemoulin, P., Brönnimann, S., Brunet, M., Crouthamel,

931    R.I., Grant, A.N., Groisman, P.Y., Jones, P.D., Kruk, M.C., Kruger, A.C., Marshall, G.J.,

932    Maugeri, M., Mok, H.Y., Nordli, Ø., Ross, T.F., Trigo, R.M., Wang, X.L., Woodruff, S.D.,

933    Worley, S.J.: The Twentieth Century Reanalysis Project. Quart. J. Roy. Meteor. Soc., 137,

934    1-28, https://doi.org/10.1002/qj.776, 2011.

935    Liu, C., Ikeda, K., Rasmussen, R., Barlage, M., Newman, A.J., Prein, A.F., Chen, F., Chen,

936    L., Clark, M., Dai, A. Dudhia, J., Eidhammer, T., Gochis, D., Gutmann, E., Kurkute, S.,Li,

937    Y., Thompson, G., and Yates, D.: Continental-scale convection-permitting modeling of the

938    current and future climate of North America, Clim Dyn (2017) 49, 71-95,

939    https://doi.org/10.1007/s00382-016-3327-9, 2017.

940    Cox, P. M.: Description of the TRIFFID dynamic global vegetation model. Hadley Centre

941    technical note, 24, 1-16, 2001.

942    Cox, P. M., Betts, R. A., Bunton, C. B., Essery, R. L. H., Rowntree, P. R., and Smith, J. :

943    The impact of new land surface physics on the GCM simulation of climate and climate

944    sensitivity, Climate Dynamics, 15(3), 183-203, 1999.

945    Daly, C., Halbleib, M., Smith, J.I., Gibson, W.P., Doggett, M.K., Taylor, G.H., Curtis, J.

946    and Pasteris, P.P.: Physiographically sensitive mapping of climatological temperature and

947    precipitation across the conterminous United States, International Journal of Climatology:

948    a Journal of the Royal Meteorological Society, 28(15), 2031-2064, 2008.

949    Dee, D.P., Uppala, S.M., Simmons, A.J., Berrisford, P., Poli, P., Kobayashi, S., Andrae,

950    U., Balmaseda, M.A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A.C.M., van de Berg,

951    L., Bidlot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A.J., Haimberger,

952    L., Healy, S.B., Hersbach, H., Hólm, E.V., Isaksen, L., Kållberg, P., Köhler, M.,

953    Matricardi, M., McNally, A.P., Monge-Sanz, B.M., Morcrette, J.J., Park, B.K., Peubey, C.,

954    de Rosnay, P., Tavolato, C., Thépaut, J.N., Vitart, F. :The ERA-Interim reanalysis:

955    configuration and performance of the data assimilation system, Quarterly Journal of the

956    Royal Meteorological Society 137656 553–597, DOI: 10.1002/qj.828, 2011.

957    Essery, R. L. H., Best, M. J., Betts, R. A., Cox, P. M., and Taylor, C. M.: Explicit

958    Representation of Subgrid Heterogeneity in a GCM Land Surface Scheme, J.

959    Hydrometeorol.,              4,              530–543,              doi:10.1175/1525-

960    7541(2003)004<0530:EROSHI>2.0.CO;2, 2003.

961    Fan, Y. and van den Dool, H.: A Global Monthly Land Surface Air Temperature Analysis

962    for 1948-Present,  J. Geophys. Res., 113, doi: 10.1029/2007JD008470, 2008.

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

964    impacts studies: recent advances in downscaling techniques for hydrological

965    modelling, International journal of climatology, 27(12), 1547-1578, 2007.

966    Gelaro, R., McCarty, W., Suárez, M.J., Todling, R., Molod, A., Takacs, L., Randles, C.A.,

967    Darmenov, A., Bosilovich, M.G., Reichle, R. and Wargan, K.: The modern-era

968    retrospective analysis for research and applications, version 2 (MERRA-2).,Journal of

969    Climate, 30(14), 5419-5454., J. Clim., doi: 10.1175/JCLI-D-16-0758.1, 2017.

970    Gregory, D., and Rowntree, P. R. :A mass flux convection scheme with representation of

971    cloud ensemble characteristics and stability-dependent closure, Monthly Weather Review,

972    118(7), 1483-1506, 1990.

973    Guillod, B. P., Jones, R. G., Bowery, A., Haustein, K., Massey, N. R., Mitchell, D. M.,

974    Otto, F. E. L., Sparrow, S. N., Uhe, P., Wallom, D. C. H., Wilson, S., and Allen, M. R.:

975    weather@home 2: validation of an improved global–regional climate modelling system,

976    Geosci. Model Dev., 10, 1849-1872, DOI:10.5194/gmd-10-1849-2017, 2017.

977   Guillod, B.P., Jones, R.G., Dadson, S.J., Coxon, G., Bussi, G., Freer, J., Kay, A.L., Massey,

978   N.R., Sparrow, S.N., Wallom, D.C. and Allen, M.R. :A large set of potential past, present

979   and future hydro-meteorological time series for the UK. Hydrology and Earth System

980   Sciences, 22(1), 611-634, https://doi.org/10.5194/hess-22-611-2018, 2018.

981   Hall, A., and Qu, X. (2006). Using the current seasonal cycle to constrain snow albedo

982   feedback in future climate change, Geophysical Research Letters, 33(3), 2006.

983   Harris, I.P.D.J., Jones, P.D., Osborn, T.J. and Lister, D.H.: Updated high-resolution grids

984   of monthly climatic observations–the CRU TS3. 10 Dataset, International journal of

985   climatology, 34(3), 623-642, doi:10.1002/joc.3711, 2014.

986   Hartmann, D.L., Ockert-Bell, M.E., and Michelsen, M.L.,:The effect of cloud type on

987   Earth's energy balance: Global analysis, Journal of Climate, 5(11),1281-1304,

988   https://doi.org/10.1175/1520-0442(1992)005<1281:TEOCTO>2.0.CO;2, 1992.

989   Hawkins, E., Osborne, T. M., Ho, C. K., and Challinor, A. J. : Calibration and bias

990   correction of climate projections for crop modelling: an idealised case study over Europe,

991   Agricultural and Forest Meteorology, 170, 19-31, 2013.

992   Hourdin, F., Grandpeix, J.-Y., Rio, C., Bony, S., Jam, A., Cheruy, F., Rochetin, N.,

993   Fairhead, L., Idelkadi, A., Musat, I., Dufresne, J.L., Lahellec, A., Lefebvre, M.-P., and

994   Roehrig, R.: LMDZ5B: the atmospheric component of the IPSL climate model with

995   revisited parameterizations for clouds and convection, Climate Dynamics, 40, 2193–2222,

996   doi:10.1007/s00382-012-1343-y, http://dx.doi.org/10. 1007/s00382-012-1343-y, 2013.

997   Irvine, P. J., Gregoire, L. J., Lunt, D. J., and Valdes, P. J.: An efficient method to generate

998   a perturbed parameter ensemble of a fully coupled AOGCM without flux-adjustment,

999   Geoscientific Model Development, 6(5), 1447-1462, 2013.

1000   Johns, T.C., Durman, C.F., Banks, H.T., Roberts, M.J., McLaren, A.J., Ridley, J.K., Senior,

1001   C.A., Williams, K.D., Jones, A., Rickard, G.J., Cusack, S., Ingram, W.I., Crucifix, M.,

1002   Sexton, D. M. H., Joshi, M.M., Dong, B.-W., Spencer, H., Hill, R. S. R., Gregory, J.M.,

1003   Keen, A.B., Pardaens, A.K., Lowe, J.A., Bodas-Salcedo, A., Stark, S., and Searl, Y. : The

1004   new Hadley Centre climate model (HadGEM1): Evaluation of coupled simulations,

1005   Journal of Climate, 19(7), 1327-1353, https://doi.org/10.1175/JCLI3712.1,  2006.

1006   Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M.,

1007   Saha, S., White, G., Woollen, J., Zhu, Y., Chelliah, M., Ebisuzaki, W., Higgins, W.,

1008   Janowiak, J., Mo, K. C., Ropelewski, C., Wang, J., Leetmaa, A., Reynolds, R., Jenne, R.,

1009   and Joseph, D.: The NCEP/NCAR 40-Year Reanalysis Project, Bull. Am. Meteorol. Soc.,

1010   77, 437–471, 1996.

1011   Knutti, R., Stocker, T. F., Joos, F., and Plattner, G. K.: Probabilistic climate change

1012   projections using neural networks, Climate Dynamics, 21(3-4), 257-272, 2003.

1013   Kotlarski, S., Keuler, K., Christensen, O. B., Colette, A., Déqué, M., Gobiet, A., Goergen,

1014   K., Jacob, D., Lüthi, D., van Meijgaard, E., Nikulin, G., Schär, C., Teichmann, C., Vautard,

1015   R., Warrach-Sagi, K., and V. Wulfmeyer: Regional climate modeling on European scales:

1016   a joint standard evaluation of the EURO-CORDEX RCM ensemble, Geoscientific Model

1017   Development, 7(4), 1297-1333, 2004.

1018   Langenbrunner, B., Neelin, J.D., Lintner, B.R., and Anderson, B.T.: Patterns of

1019   precipitation change and climatological uncertainty among CMIP5 models, with a focus

1020   on the midlatitude Pacific storm track, Journal of Climate, 28(19), 7857-7872, 2015.

1021    Li, S., Mote, P. W., Rupp, D. E., Vickers, D., Mera, R., and Allen, M.R.: Evaluation of a

1022    regional climate modeling effort for the western United States using a superensemble from

1023    weather@ home, Journal of Climate, 28(19), 7470-7488, 2015.

1024    Loeppky J.L., Sacks J., and Welch W.J.: Choosing the Sample Size of a Computer

1025    Experiment: a Practical Guide, Technometrics, 51(4):366–376, 2009.

1026    Ma, H.Y., Klein, S.A., Xie, S., Zhang, C., Tang, S., Tang, Q., Morcrette, C.J., Van

1027    Weverberg, K., Petch, J., Ahlgrimm, M., Berg, L.K., Cheruy, F., Cole, J., Forbes, R.,

1028    Gustafson Jr, W. I., Huang, M., Liu, Y., Merryfield, W., Qian, Y., Roehrig, R., and Wang,

1029    Y.-C.: CAUSES: On the role of surface energy budget errors to the warm surface air

1030    temperature error over the Central United States, Journal of Geophysical Research:

1031    Atmospheres, 123, 2888–2909, https://doi.org/10.1002/2017JD027194, 2018.

1032    Massey, N., Jones, R., Otto, F. E. L., Aina, T., Wilson, S., Murphy, J. M., Hassell, D.,

1033    Yamazaki, Y. H., and Allen, M. R.: weather@home—development and validation of a

1034    very large ensemble modelling system for probabilistic event attribution, Quarterly Journal

1035    of the Royal Meteorological Society, 141, 1528–1545, doi:10.1002/qj.2455, 2015.

1036    Mauritsen, T., Stevens, B., Roeckner, E., Crueger, T., Esch, M., Giorgetta, M., Haak, H.,

1037    Jungclaus, J., Klocke, D., Matei, D., Mikolajewicz, U., Notz, D., Pincus, R., Schmidt, H.,

1038    and Tomassini, L.: Tuning the climate of a global model, Journal of Advances in Modeling

1039    Earth Systems, 4, M00A01, doi:doi:10.1029/2012MS000154, 2012.

1040    McKay, M. D., Beckman, R. J., and Conover, W. J.: Comparison of three methods for

1041    selecting values of input variables in the analysis of output from a computer code,

1042    Technometrics, 21(2), 239-245., 1979.

1043    McNeall, D. J., Challenor, P. G., Gattiker, J. R., and Stone, E. J.: The potential of an

1044    observational data set for calibration of a computationally expensive computer model,

1045    Geosci. Model Dev., 6, 1715–1728, doi: 10.5194, 2013.

1046    McNeall, D., Williams, J., Booth, B., Betts, R., Challenor, P., Wiltshire, A., and Sexton,

1047    D.: The impact of structural error on parameter constraint in a climate model, Earth System

1048    Dynamics, 7(4), 917-935, 2016.

1049    Mearns, L.O., Arritt, R., Biner, S., Bukovsky, M.S., McGinnis, S., Sain, S., Caya, D.,

1050    Correia Jr, J., Flory, D., Gutowski, W., Takle, E.S., Jones, R., Leung, R., Moufouma-Okia,

1051    W., McDaniel, L., Nues, A.M.B., Qian, Y., Roads-*, J., Sloan., L., and Snyder, M.: The

1052    North American regional climate change assessment program: overview of phase I results,

1053    Bulletin of the American Meteorological Society, 93(9), 1337-1362, 2012.

1054    Merrifield, A. L., and  Xie, S. P.: Summer US surface air temperature variability:

1055    Controlling factors and AMIP simulation biases, Journal of Climate, 29(14), 5123–5139.

1056    https://doi.org/10.1175/JCLI-D-15-0705.1, 2016.

1057    Mitchell, D., Heaviside, C., Vardoulakis, S., Huntingford, C., Masato, G., Guillod, B. P.,

1058    Frumhoff, P., Bowery, A., Wallom, D., and Allen, M.: Attributing human mortality during

1059    extreme heat waves to anthropogenic climate change, Environ. Res. Lett., 11, 074006,

1060    doi:10.1088/1748-9326/11/7/074006, 2016.

1061    Mock, C. J.: Climatic Controls and Spatial Variations of Precipitation in the Western

1062    United States, J. Climate, 9(5), 1111–1125, 1996.

1063    Morcrette, C.J., Van Weverberg, K., Ma, H.Y., Ahlgrimm, M., Bazile, E., Berg, L.K.,

1064    Cheng, A., Cheruy, F., Cole, J., Forbes, R. and Gustafson Jr, W.I.: Introduction to

1065    CAUSES: Description of weather and climate models and their near-surface temperature

Geoscientific
Model Development
Discussions

1066    errors in 5 day hindcasts near the Southern Great Plains, Journal of Geophysical Research:

1067    Atmospheres, 123(5), pp.2655-2683. https://doi.org/10.1002/2017JD027199, 2018.

1068    Mote, P.W., Allen, M.R., Jones, R.G., Li, S., Mera, R., Rupp, D.E., Salahuddin, A. and

1069    Vickers, D.: Superensemble regional climate modeling for the western United States,

1070    Bulletin of the American Meteorological Society (97), 203-215, doi: 10.1175/BAMS-D-

1071    14-00090.1, 2016.

1072    Mueller, B., and Seneviratne S. I.: Systematic land climate and evapotranspiration biases

1073    in CMIP5 simulations, Geophys. Res. Lett., 41, 128–134, doi:10.1002/2013GL058055,

1074    2014.

1075    Mulholland, D. P., Haines, K., Sparrow, S. N., and Wallom, D.: Climate model forecast

1076    biases assessed with a perturbed physics ensemble, Climate Dynamics, 49(5-6), 1729-

1077    1746, 2017.

1078    Murphy, J. M., Sexton, D. M. H., Barnett, D. N., Jones, G. S., Webb, M. J., Collins, M.,

1079    and Stainforth, D. A.: Quantification of modelling uncertainties in a large ensemble of

1080    climate change simulations, Nature, 430, 768–772, doi:10.1038/nature02771, 2004.

1081    Neelin, J. D., Bracco, A., Luo, H., McWilliams, J.C., and Meyerson, J. E.: Considerations

1082    for parameter optimization and sensitivity in climate models. Proc. Natl. Acad. Sci. USA,

1083    107(50), 21349–21354, doi:10.1073/pnas.1015473107, 2010.

1084    Neelin, J.D., Langenbrunner, B., Meyerson, J.E., Hall, A. and Berg, N.: California winter

1085    precipitation change under global warming in the Coupled Model Intercomparison Project

1086    phase 5 ensemble, Journal of Climate, 26(17), 6238-6256, 2013.

1087    oceanic quasi-equilibrium states, J. Atmos. Sci., 59, 1885– 1897.

1088    oceanic quasi-equilibrium states, J. Atmos. Sci., 59, 1885– 1897Bellprat, O., Kotlarski, S.,

1089    Lüthi, D., De Elía, R., Frigon, A., Laprise, R., and Schär, C.: Objective Calibration of

1090    Regional Climate Models: Application over Europe and North America, Journal of

1091    Climate, 29, 819–838, doi:10.1175/jcli-d-15-0302.1, 2016.Williamson, D., Goldstein, M.,

1092    Allison, L., Blaker, A., Challenor, P., Jackson, L., and Yamazaki, K.: History matching for

1093    exploring and reducing climate model parameter space using observations and a large

1094    perturbed physics ensemble, Climate Dynamics, 41, 1703–1729, doi:10.1007/s00382-013-

1095    1896-4, http://dx.doi.org/10.1007/s00382-013-1896-4, 2013.

1096    Onogi, K., Tsutsui, J., Koide, H., Sakamoto, M., Kobayashi, S., Hatsushika, H.,

1097    Matsumoto, T., Yamazaki, N., Kamahori, H., Takahashi, K., Kadokura, S., Wada, K., Kato,

1098    K., Oyama, R., Ose, T., Mannoji, N., and Taira, R.: The JRA-25 Reanalysis, J. Met. Soc.

1099    Jap., 85(3), 369-432, doi: 10.2151/jmsj.85.369, 2007.

1100    Otto, F.E.L., Massey, N., van Oldenborgh, G.J., Jones, R.G. and Allen, M.R.:Reconciling

1101    Two Approaches to Attribution of the 2010 Russian Heat Wave, Geophysical Research

1102    Letters, 39(L04702), 2012.

1103    Rosenzweig, C., Elliott, J., Deryng, D., Ruane, A.C., Müller, C., Arneth, A., Boote, K.J.,

1104    Folberth, C., Glotter, M., Khabarov, N., Neumann, K., Piontek, F., Pugh, T.A.M., Schmid,

1105    E., Stehfest, E., Yang, H., and Jones, J.W. :Assessing agricultural risks of climate change

1106    in the 21st century in a global gridded crop model intercomparison, Proceedings of the

1107    National Academy of Sciences, 111(9), 3268-3273, 2014.

1108    Rougier, J.C., Sexton, D.M.H., Murphy, J.M., and Stainforth, D. : Analyzing the climate

1109    sensitivity of the HadSM3 climate model using ensembles from different but related

1110    experiments, J Clim 22(13):3540–3557, https://doi.org/10.1175/2008JCLI2533.1, 2009.

Roustant, O., Ginsbourger, D., and Deville, Y.: DiceKriging, DiceOptim: Two R Packages for the Analysis of Computer Experiments by Kriging-Based Metamodeling and Optimization, Journal of Statistical Software, 51(i01), 2012.

Rowlands, D. J., Frame, D. J., Ackerley, D., Aina, T., Booth, B. B., Christensen, C., and Gryspeerdt, E.: Broad range of 2050 warming from an observationally constrained large climate model ensemble, Nature Geoscience, 5(4), 256, 2012.

Rupp, D.E., Li, S., Mote, P.W., Massey, N., Sparrow, S.N., and Wallom, D.C.: Influence of the Ocean and Greenhouse Gases on Severe Drought Likelihood in the Central US in 2012, Journal of Climate (30), 1789-1806, doi: 10.1175/JCLI-D-16-0294.1, 2017a.

Rupp, D. E., Li, S., Mote, P. W., Shell, K.M., Massey, N., Sparrow, S. N., Wallom, D. C. H., and Allen, M. R.: Seasonal Spatial Patterns of Projected Anthropogenic Warming in Complex Terrain: A Modeling Study of the Western USA, Climate Dynamics (48), 2191-2213, doi: 10.1007/s00382-016-3200-x, 2017b.

Rupp, D. E. and Li, S.: Less warming projected during heavy winter precipitation in the Cascades and Sierra Nevada, Int. J. Climatol., 37(10): 3984–3990. doi:10.1002/joc.4963, 2017.

Saha, S., Moorthi, S., Pan, H.L., Wu, X., Wang, J., Nadiga, S., Tripp, P., Kistler, R., Woollen, J., Behringer, D. and Liu, H.: The NCEP climate forecast system reanalysis, Bulletin of the American Meteorological Society, 91(8), 1015-1058, 2010.

Saltelli, A., Tarantola, S., and Chan, K. S.: A quantitative model-independent method for global sensitivity analysis of model output, Technometrics, 41(1), 39-56, 1999.

1132    Sanderson, B. M.: A multimodel study of parametric uncertainty in predictions of climate

1133    response to rising greenhouse gas concentrations, Journal of Climate, 24(5), 1362-1377,

1134    2011.

1135    Sanderson, B. M., Knutti, R., Aina, T., Christensen, C., Faull, N., Frame, D. J., Ingram,

1136    W.J., Piani, C., Stainforth, D.A., Stone, D.A., and Allen, M. R.: Constraints on model

1137    response to greenhouse gas forcing and the role of subgrid-scale processes, Journal of

1138    Climate, 21(11), 2384-2400, 2008a--ANN.

1139    Sanderson, B. M., Piani, C., Ingram, W. J., Stone, D. A., and Allen, M. R.: Towards

1140    constraining climate sensitivity by linear analysis of feedback patterns in thousands of

1141    perturbed-physics GCM simulations, Climate Dynamics, 30(2-3), 175-190, 2008b--

1142    entcoef.

1143    Sanderson, B. M., Shell, K. M., and Ingram, W.: Climate feedbacks determined using

1144    radiative kernels in a multi-thousand member ensemble of AOGCMs, Climate dynamics,

1145    35(7-8), 1219-1236, 2010.

1146    Schaller, N., Kay, A.L., Lamb, R., Massey, N.R., van Oldenborgh, G.J., Otto, F.E.L.,

1147    Sparrow, S.N., Vautard, R., Yiou, P., Ashpole, I., Bowery, A., Crooks, S.M., Haustein, K.,

1148    Huntingford, C., Ingram, W.J., Jones, R.G., Legg, T., Miller, J., Skeggs, J., Wallom, D.,

1149    Weisheimer, A., Wilson, S., Stott, P.A. and Allen, M.R. : Human Influence on Climate in

1150    the 2014 Southern England Winter Floods and Their Impacts, Nature Climate Change, 6:

1151    627-634, 2016.

1152    Schirber, S., Klocke, D., Pincus, R., Quaas, J., and Anderson, J. L.: Parameter estimation

1153    using data assimilation in an atmospheric general circulation model: From a perfect toward

1154    the real world, Journal of Advances in Modeling Earth Systems,5(1), 58–70,

1155    doi:10.1029/2012MS000167, 2013.

1156    Schneider, U., Becker, A., Finger, P., Meyer-Christoffer, A., Rudolf, B., and Ziese, M.:

1157    GPCC Full Data Reanalysis Version 6.0 at 0.5°: Monthly Land-Surface Precipitation from

1158    Rain-Gauges        built        on        GTS-based        and        Historic        Data,

1159    doi:10.5676/DWD_GPCC/FD_M_V6_050, 2011.

1160    Seager, R., Neelin, D., Simpson, I., Liu, H., Henderson, N., Shaw, T., Kushnir, Y., Ting,

1161    M. and Cook, B.: Dynamical and thermodynamical causes of large-scale changes in the

1162    hydrological cycle over North America in response to global warming, Journal of Climate,

1163    27(20), 7921-7948, 2014.

1164    Seneviratne, S. I., Lüthi, D., Litschi, M., and Schär, C.: Land–atmosphere coupling and

1165    climate change in Europe, Nature, 443(7108), 205, 2006.

1166    Sexton, D. M., Murphy, J. M., Collins, M., and Webb, M. J.: Multivariate probabilistic

1167    projections using imperfect climate models part I: outline of methodology, Climate

1168    dynamics, 38(11-12), 2513-2542, 2012.

1169    Sippel, S., Otto, F.E., Forkel, M., Allen, M.R., Guillod, B.P., Heimann, M., Reichstein, M.,

1170    Seneviratne, S.I., Thonicke, K. and Mahecha, M.D.: A novel bias correction methodology

1171    for climate impact simulations, Earth System Dynamics, 7(1), 71-88, 2016.

1172    Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K. B., Tignor, M., and

1173    Miller, H. L.: Climate change 2007: The physical science basis, in Contribution of Working

1174    Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate

1175    Change, 2007 Cambridge University Press, Cambridge, United Kingdom and

1176    New York, NY, USA, 2007.

1177 Sparrow, S., Wallom, D., Mulholland, D. P., and Haines, K.: Climate model forecast biases

1178 assessed with a perturbed physics ensemble, Climate Dynamics, 2016.

1179 Stainforth, D. A., Aina, T., Christensen, C., Collins, M., Faull, N., Frame, D. J.,

1180 Kettleborough, J. A., Knight, S., Martin, A., Murphy, J., Piani, C., Sexton, D., Smith, L.

1181 A., Spicer, R. A., Thorpe, A. J., and Allen, M. R.: Uncertainty in predictions of the climate

1182 response to rising levels of greenhouse gases, Nature, 433(7024), 403–406, 2005.

1183 Stephens, G.L.:Cloud feedbacks in the climate system: A critical review, Journal of

1184 climate, 18(2), 237-273, https://doi.org/10.1175/JCLI-3243.1, 2005.

1185 Stephens, G. L., Li, J., Wild, M., Clayson, C. A., Loeb, N., Kato, S., L'ecuyer, T.,

1186 Stackhouse Jr, P.W., Lebsock, M. and Andrews, T. : An update on Earth's energy balance

1187 in light of the latest global observations, Nature Geoscience, 5(10), 691, 2012.

1188 Stott, P. A., Stone, D. A., and Allen, M. R.: Human contribution to the European heatwave

1189 of 2003, Nature, 432(7017), 610, 2004.

1190 Tett, S. F., Mitchell, J. F., Parker, D. E., and Allen, M. R.: Human influence on the

1191 atmospheric vertical temperature structure: Detection and observations, Science,

1192 274(5290), 1170-1173, 1996.

1193 Tett, S. F., Yamazaki, K., Mineter, M. J., Cartis, C., and Eizenberg, N.: Calibrating climate

1194 models using inverse methods: case studies with HadAM3, HadAM3P and HadCM3,

1195 Geoscientific Model Development, 10(9), 3567-3589, 2017.

1196 Uhe, P., Philip, S., Kew, S., Shah, K., Kimutai, J., Mwangi, E., van Oldenborgh, G.J.,

1197 Singh, R., Arrighi, J., Jjemba, E., Cullen, H. and Otto, F.E.L.: Attributing Drivers of the

1198 2016 Kenyan Drought, International Journal of Climatology, 38, e554-e568, 2018.

1199   van Oldenborgh, G.J., Otto, F.E.L., Haustein, K. and AchutaRao, K.: The Heavy

1200   Precipitation Rvent of December 2015 in Chennai, India, In Explaining Extremes of 2015

1201   from a Climate Perspective. Bulletin of the American Meteorological Society, 97(12), S87-

1202   S91, 2016.

1203   van Oldenborgh, G.J., van der Wiel, K., Sebastian, A., Singh, R., Arrighi, J., Otto, F. E.L.,

1204   Haustein, K., Li, S., Vecchi, G. and Cullen, H. : Attribution of Extreme Rainfall from

1205   Hurricane Harvey, August 2017, Environmental Research Letters,12(12),124009, 2017.

1206   Van Weverberg, K., Morcrette, C.J., Petch, J., Klein, S.A., Ma, H.Y., Zhang, C., Xie, S.,

1207   Tang, Q., Gustafson Jr, W.I., Qian, Y. and Berg, L.K., : CAUSES: Attribution of surface

1208   radiation biases in NWP and climate models near the U.S. Southern Great Plains, Journal

1209   of      Geophysical      Research:      Atmospheres,      123,      3612–3644.

1210   https://doi.org/10.1002/2017JD027188, 2018.

1211   Williams, J. H. T., Smith, R. S., Valdes, P. J., Booth, B. B. B., and Osprey, A.: Optimising

1212   the FAMOUS climate model: inclusion of global carbon cycling, Geoscientific Model

1213   Development, 5, 3089-3129, 2013.

1214   Williamson, D., Goldstein, M., Allison, L., Blaker, A., Challenor, P., Jackson, L. and

1215   Yamazaki, K.: History matching for exploring and reducing climate model parameter space

1216   using observations and a large perturbed physics ensemble. Climate dynamics, 41(7-8),

1217   1703-1729, 2013.

1218   Williamson, D., Blaker, A. T., Hampton, C., and Salter, J.: Identifying and removing

1219   structural biases in climate models with history matching, Climate dynamics, 45(5-6),

1220   1299-1324, 2015.

1221    Williamson, D. B., Blaker, A. T., and Sinha, B.: Tuning without over-tuning: parametric

1222    uncertainty quantification for the NEMO ocean model, Geoscientific Model Development,

1223    10(4), 1789-1816, doi:10.5194/gmd-10-1789-2017,2017.

1224    Wu, X. : Effects of ice mircrophysics on tropical radiative-convective- oceanic quasi-

1225    equilibrium states, J. Atmos. Sci., 59, 1885– 1897, 2002.

1226    Wu, X. (2002), Effects of ice microphysics on tropical radiative-convective-

1227    Wu, X. (2002), Effects of ice microphysics on tropical radiative-convective-

1228    Yamazaki, K., Rowlands, D. J., Aina, T., Blaker, A., Bowery, A., Massey, N., Miller, J.,

1229    Rye, C., Tett, S. F. B., Williamson, D., Yamazaki, Y. H., and Allen, M. R.: Obtaining

1230    diverse behaviors in a climate model without the use of flux adjustments, JGR-

1231    Atmospheres, 118, 2781–2793, doi:10.1002/jgrd.50304, 2013.

1232    Zhang, C., Xie, S., Klein, S.A., Ma, H.Y., Tang, S., Van Weverberg, K., Morcrette, C.J.

1233    and Petch, J.: CAUSES: Diagnosis of the summertime warm bias in CMIP5 climate models

1234    at the ARM Southern Great Plains site, Journal of Geophysical Research: Atmospheres,

1235    123, 2968–2992. https:// doi.org/10.1002/2017JD027200, 2018.

1236    Zhang, T., Li, L., Lin, Y., Xue, W., Xie, F., Xu, H., and Huang, X.: An automatic and

1237    effective parameter optimization method for model tuning, Geoscientific Model

1238    Development, 8, 3579–3591, doi:10.5194/gmd-8-3579-2015, http://www.geosci-model-

1239    dev.net/8/3579/2015/, 2015.

1240

Geoscientific
Model Development
Discussions



1241

**Figure 1.** Global mean top-of-atmosphere (TOA) outgoing (reflected) shortwave radiation

(SW) and outgoing longwave radiation (LW) from the four ensembles run through

weather@home2. Horizontal and vertical dashed lines denote the reference values for SW

and LW taken from Stephens et al. (2012). The filled brown circle denotes our SP. The

ellipse indicates the uncertainty ranges we are willing to accept for SW and LW

respectively, which includes the observational uncertainty range taken from Stephens et al.

(2012), but tripled, plus the uncertainty range due to initial condition perturbations

estimated from our SP reference ensemble. The red solid lines highlight net TOA energy

flux of +/- 5 $Wm^{-2}$.

1251

**Figure 2.** One-at-a-time sensitivity analysis of magnitude of annual cycle of temperature
(MAC-T) over Northwest to each input parameter in turn, with all other parameters held at
mean value of all the designed points. Heavy lines represent the emulator mean, and shaded
areas represent the estimate of emulator uncertainty, at the ±1 SD level.

1268

**Figure 3.** Sensitivity analysis of model output metrics in Phase 2 via the FAST algorithm

of Saltelli et el. (1999).



**Figure 4.** Phase 3 PPE parameter inputs and summary model output metrics evaluated. 95

parameter sets are shown. The parameter values and model outputs under SP are marked

in red. The horizontal and vertical red lines mark the transition from parameter inputs and

model output metrics.

**Figure 5.** Comparison between three PPEs and SP zonal mean HadAM3P simulated North

Hemisphere mid-latitude (30°N-60°N) a) DJF mean temperature over land, b) JJA mean

temperature over land, c) DJF mean precipitation, and d) JJA mean precipitation. Output

from the selected 10 parameter sets selected, based on NWUS MAC-T, are shown in blue.

Note that the plotting order is the same as the legend, so most Phase 1 curves are obscured

by subsequent phases.

1287

**Figure 6.** Biases of SP temperature over land in a) DJF, b) MAM, c) JJA, and d) SON,

compared with CRU over December 1996 through November 2007. Biases of selected PP

compared with CRU are shown in e)-h), while the differences between selected PP and SP,

i.e. the absolute increase or decrease of biases in PP with respect to the SP values, are

shown in i) - l). The PP results are the composites of the 10 selected sets, 6 IC per set.

Geoscientific
Model Development
Discussions

Open Access

EGU



1293

1294    **Figure 7.** Same as Fig. 6, but for precipitation.

**Figure 8.** Annual (a,d), DJF (b,e) and JJA (c,f) meridional distributions of precipitation from Phase 3 and SP (all panels), reanalysis datasets MERRA2, JRA-55, CFSR, ERAI and 20CRv2c shown (a - c) and GCMs CanAM4-AMIP, CESM1-CAM5, and HadGEM2-A, shown in (d - f ).

1300

**Figure B1.** Emulator predicted results vs. model simulated results in Phase 2 for different

model output metrics based on 94 parameter sets not used to train the emulator (the 94 sets

that finished after starting Phase3). The regression coefficient (regcoef) and coefficient of

determination ($R^2$) by emulated results are shown in each panel. The dashed line in each

panel denotes the 1:1 line.

Geoscientific
Model Development
Discussions

Open Access



1306

**Figure B2.** Same as Fig. B1, but for the 95 parameter sets in Phase 3. Note the ranges of

x- and y-axis are set to be the same as in Fig. B1.

1307

1308

1309

1310

1311

1312

1313

1314

**Table 1.** Parameters perturbed in our tuning exercise with the

post-culling parameters highlighted in bold.

| Parameter | Default | Low | High | Description | Model component |
|---|---|---|---|---|---|
| CT ($s^{-1}$) | $6 \times 10^{-4}$ | $0.5 \times 10^{-4}$ | $1.2 \times 10^{-3}$ | Rate at which cloud liquid water is converted to precipitation | Cloud |
| CW_SEA (kg $m^{-3}$) | $2.0 \times 10^{-5}$ | $0.5 \times 10^{-5}$ | $2.0 \times 10^{-4}$ | Threshold cloud liquid water content over sea | Cloud |
| **CW_LAND** (kg $m^{-3}$) | $1.0 \times 10^{-3}$ | $0.5 \times 10^{-3}$ | $1.0 \times 10^{-2}$ | Threshold cloud liquid water content over land | Cloud |
| EACF | 0.5 | 0.5 | 0.6 | Empirically adjusted cloud fraction | Cloud |
| **VF1** (m $s^{-1}$) | 2 | 0.5 | 4 | Ice fall speed | Cloud |
| **ENTCOEF** | 3 | 0.3 | 9.5 | Entrainment rate coefficient | Convection |
| ALPHAM | 0.5 | 0.45 | 0.65 | Albedo at melting point of sea ice | Radiation |
| DTICE (°C) | 10 | 2 | 11 | Temperature range over which ice albedo varies | Radiation |
| ICE_SIZE (m) | $3.0 \times 10^{-5}$ | $2.5 \times 10^{-5}$ | $4.0 \times 10^{-5}$ | Ice particle size | Radiation |
| KAY_GWAVE (m) | $1.8 \times 10^{4}$ | $1.0 \times 10^{4}$ | $2.0 \times 10^{4}$ | Surface gravity wave drag: typical wavelength | Dynamics |
| KAY_LEE_GWAVE ($m^{-3/2}$) | $2.7 \times 10^{5}$ | $1.5 \times 10^{5}$ | $3.0 \times 10^{5}$ | Surface gravity wave trapped lee wave constant | Dynamics |
| START_LEVEL_GWDRAG | 3 | 3 | 5 | Lowest model level for gravity wave drag | Dynamics |
| **V_CRIT_ALPHA** | 0.5 | 0.01 | 0.99 | Control of photosynthesis with soil moisture | Land surface |
| **ASYM_LAMBDA** | 0.15 | 0.05 | 0.5 | Vertical distance over which air parcels travel before mixing | Boundary layer |

| | | | | with their surroundings | |
|---|---|---|---|---|---|
| CHARNOCK | 0.012 | 0.009 | 0.020 | Constant in Charnock formula for calculating roughness length for momentum transport over sea | Boundary layer |
| **G0** | 10 | 5 | 20 | Used in calculation of stability function for heat, moisture, and momentum transport | Boundary layer |
| **Z0FSEA** (m) | $1.3\times10^{-3}$ | $2.0\times10^{-4}$ | $5\times10^{-3}$ | Roughness length for free heat and moisture transport over the sea | Boundary layer |

1315

1316 **Table 2.** The specifics of four ensembles used in this study.

1317

| Experiment | Start dates | Number of parameters | Number of parameter sets in PPE | IC per parameter set per year used in the analysis |
|---|---|---|---|---|
| SP | 1 Dec 1995, 1996, …, 2005 | 1 | 1 | 6 |
| PPE Phase 1 | 1 Dec 1995 | 17 | 220 | 3 |
| PPE Phase 2 | 1 Dec 1995, 1996, …, 2005 | 17 | 264 | 3 |
| PPE Phase 3 | 1 Dec 1995, 1996, …, 2005 | 7 | 95 | 6 |

1318