

1 **The International Land Model Benchmarking (ILAMB) System:**
2 **Design, Theory, and Implementation**

3 **Nathan Collier¹, Forrest M. Hoffman^{1,2}, David M. Lawrence³, Gretchen Keppel-Aleks⁴,**
4 **Charles D. Koven⁵, William J. Riley⁵, Mingquan Mu⁶, James T. Randerson⁶**

5 ¹Climate Change Science Institute, Oak Ridge National Laboratory, Oak Ridge, TN, USA

6 ²Department of Civil & Environmental Engineering, University of Tennessee, Knoxville, TN, USA

7 ³Climate & Global Dynamics Division, National Center for Atmospheric Research, Boulder, CO, USA

8 ⁴Department of Climate and Space Sciences and Engineering, University of Michigan, Ann Arbor, MI, USA

9 ⁵Climate Sciences Department, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

10 ⁶Department of Earth System Science, University of California, Irvine, CA, USA

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as doi: [10.1029/2018MS001354](https://doi.org/10.1029/2018MS001354)

Corresponding author: Nathan Collier, nathaniel.collier@gmail.com

Abstract

The increasing complexity of Earth system models (ESMs) has inspired efforts to quantitatively assess model fidelity through rigorous comparison with best-available measurements and observational data products. ESMs exhibit a high degree of spread in predictions of land biogeochemistry, biogeophysics, and hydrology, which are sensitive to forcing from other model components. Based on insights from prior land model evaluation studies and community workshops, the authors developed an open source model benchmarking software package that generates graphical diagnostics and scores model performance in support of the International Land Model Benchmarking (ILAMB) project. Employing a suite of *in situ*, remote sensing, and reanalysis datasets, the ILAMB package performs comprehensive model assessment across a wide range of land variables and generates a hierarchical set of webpages containing statistical analyses and figures designed to provide the user insights into strengths and weaknesses of multiple models or model versions. Described here is the benchmarking philosophy and mathematical methodology embodied in the most recent implementation of the ILAMB package. Comparison methods unique to a few specific datasets are presented, and guidelines for configuring an ILAMB analysis and interpreting resulting model performance scores are discussed. ILAMB is being adopted by modeling teams and centers during model development and for model intercomparison projects, and community engagement is sought for extending evaluation metrics and adding new observational datasets to the benchmarking framework.

1 Introduction

As Earth system models (ESMs) become increasingly complex and observational data volumes rapidly expand, there is a growing need for comprehensive and multi-faceted evaluation of model fidelity. Process-rich ESMs pose challenges to developers implementing new parameterizations or tuning process representations, and to the broader community seeking information about the skill of model predictions. Model developers and software engineers require a systematic means for evaluating changes in model results to ensure that developments improve the scientific performance of target process representations while not adversely affecting results in other, possibly less familiar, parts of the model. To advance understanding and predictability of terrestrial biogeochemical processes and their interactions with hydrology and climate under conditions of increasing atmospheric carbon dioxide, rigorous analysis methods, employing best-available observational data, are

43 required to objectively assess and constrain model predictions, inform model development,
44 and identify needed measurements and field experiments (*Hoffman et al., 2017*).

45 Building upon past model evaluation work (*Randerson et al., 2009*), we developed an
46 extensible model benchmarking package in support of the goals of the International Land
47 Model Benchmarking (ILAMB; <https://www.ilamb.org/>) activity. ILAMB's goals are
48 to

- 49 1. develop internationally accepted benchmarks for land model performance by draw-
50 ing upon international expertise and collaboration;
- 51 2. promote the use of these benchmarks by the international community for model
52 intercomparison and development;
- 53 3. strengthen linkages among experimental, remote sensing, and climate modeling
54 communities in the design of new model tests, benchmarks, and measurement pro-
55 grams; and
- 56 4. support the design and development of a new, open source, benchmarking software
57 system for use by the international community.

58 Three ILAMB workshops have been held—in Exeter, United Kingdom, in 2009; Irvine,
59 California, United States, in 2011 (*Luo et al., 2012*); and Washington, DC, United States,
60 in 2016 (*Hoffman et al., 2017*)—to engage the modeling, measurements, and remote sens-
61 ing communities in the identification of observational datasets and the design of model
62 evaluation metrics. In this way, community consensus was sought for the curation of ob-
63 servational data and the methodology of model evaluation and scoring, which are de-
64 scribed below.

65 Recognition that the capacities of the terrestrial and marine biosphere to store an-
66 thropogenic carbon will weaken under climate warming (*Cox et al., 2000; Friedlingstein*
67 *et al., 2001; Fung et al., 2005; Denman et al., 2007; Randerson et al., 2015; Mahowald*
68 *et al., 2017; Moore et al., 2018*) and that uncertainties in carbon cycle feedbacks must
69 be quantified and reduced to improve projections of future climate change (*Friedling-*
70 *stein et al., 2006; Gregory et al., 2009; Arora et al., 2013; Ciais et al., 2013; Friedlingstein*
71 *et al., 2014; Hoffman et al., 2014*) has inspired efforts to quantitatively evaluate model per-
72 formance through comparison with *in situ* and remote sensing observations (*Anav et al.,*
73 *2013; Eyring et al., 2016*). Multi-model simulation results from the third Coupled Model

74 Intercomparison Project (CMIP3; *Meehl et al.*, 2007) and fifth Coupled Model Intercom-
75 parison Project (CMIP5; *Taylor et al.*, 2012), which informed the Intergovernmental Panel
76 on Climate Change (IPCC) Fourth and Fifth Assessment Reports (AR4 and AR5), pro-
77 vided opportunities for developing and testing model evaluation diagnostics, formal met-
78 rics, and exploration of benchmarking concepts and techniques. Early work on coupled
79 model evaluation and establishing formal metrics focused primarily on atmospheric vari-
80 ables (*Reichler and Kim*, 2008; *Gleckler et al.*, 2008). Following the first two ILAMB
81 workshops, the land modeling community began exploring standardized and comprehen-
82 sive benchmarking for terrestrial carbon cycle models (*Cadule et al.*, 2010; *Blyth et al.*,
83 2011; *Abramowitz*, 2012; *Kelley et al.*, 2013; *Dalmonech and Zaehle*, 2013; *Piao et al.*,
84 2013; *Anav et al.*, 2013; *Bouskill et al.*, 2014; *Ghimire et al.*, 2016). While some researchers
85 define benchmarking as a series of model tests based on a pre-defined expected level of
86 performance (*Abramowitz*, 2005; *Best et al.*, 2015), most of the systematic benchmarking
87 strategies explored by the land modeling community to date do not depend upon the estab-
88 lishment of an expected level of performance.

89 The ILAMB software package, hereafter referred to as ILAMB, shares some of the
90 same goals as existing model diagnostic and evaluation tools, such as the Protocol for the
91 Analysis for Land Surface models (PALS; *Abramowitz*, 2012), the Program for Climate
92 Model Diagnosis and Intercomparison (PCMDI) Metrics Package (PMP; *Gleckler et al.*,
93 2016), the Earth System Model Evaluation Tool (ESMValTool; *Eyring et al.*, 2016), the
94 Land surface Verification Toolkit (LVT; *Kumar et al.*, 2012), and a wide variety of often
95 custom-developed diagnostic packages in use at international modeling centers. Some of
96 these tools provide model-to-model comparisons, a large collection of standalone graph-
97 ical diagnostics, or workflow infrastructure that allows one to regenerate analysis results
98 from previously published studies but with new model outputs. In contrast, ILAMB was
99 designed to compare multiple models or model versions with observations simultaneously,
100 assess functional relationships between prognostic variables and one or more forcing vari-
101 ables through variable-to-variable comparisons (e.g., gross primary production vs. precip-
102 itation), and score model performance across a suite of metrics, variables, and datasets.
103 Model performance is evaluated for variables in categories of biogeochemistry (Table 2),
104 hydrology (Table 3), radiation and energy (Table 4), and climate forcing (Table 5).

105 For every variable, ILAMB generates graphical diagnostics (spatial contour maps,
106 time series line plots, and Taylor diagrams (*Taylor*, 2001)) and scores model performance

107 for the period mean, bias, root mean squared error (RMSE), spatial distribution, inter-
108 annual coefficient of variation, seasonal cycle, and long-term trend. Model performance
109 scores are calculated for each metric and variable and are scaled based on the degree of
110 certainty of the observational dataset, the scale appropriateness, and the overall impor-
111 tance of the constraint or process to model predictions, following a customizable rubric
112 described below (Table 1). Scores are aggregated across metrics and datasets, producing a
113 single scalar score for each variable for every model or model version. As shown in Fig-
114 ure 1, these scalar scores are presented graphically. On the left side we use a stop-light
115 color scheme to indicate aggregate performance for each model by variable. On the right,
116 we show relative performance (i.e., Z-score), indicating which models or model versions
117 perform better with respect to others contained in the overall analysis.

118 We do not view these aggregate absolute scores as a determinant of ‘good’ or ‘bad’
119 models. We envision the scores as a tool to more quickly identify relative differences
120 among models and model versions which the scientist must then interpret. As in any eval-
121 uation methodology, many of our choices are subjective and must be considered as the
122 scores are interpreted. Where possible, the ILAMB implementation allows for users to
123 customize weights and diagnostics in order to incorporate aspects of model performance
124 relevant to their scientific goals. ILAMB may be thought of as a framework which may be
125 expanded to incorporate community ideas regarding model benchmarking. Thus while our
126 choices are subjective, they are informed by the preferences of a larger community and
127 can be considered as an initial suggestion.

128 The remainder of this paper describes the ILAMB methodology used to compute
129 aggregate absolute scores. First we describe how we compare an individual observational
130 dataset to model output (Section 2). Then we explain how scores are aggregated across
131 datasets for each variable and present the datasets used in the land model evaluation (Sec-
132 tion 3). In Section 4 we present some salient points about how the ILAMB software is
133 designed. Finally, in Section 5 we discuss what ILAMB scores mean and how they should
134 be used.

135 **2 Methodology**

136 In this section we describe the methodology used to assess how well a model cap-
137 tures information contained in a reference (e.g. observational) dataset. For the purposes

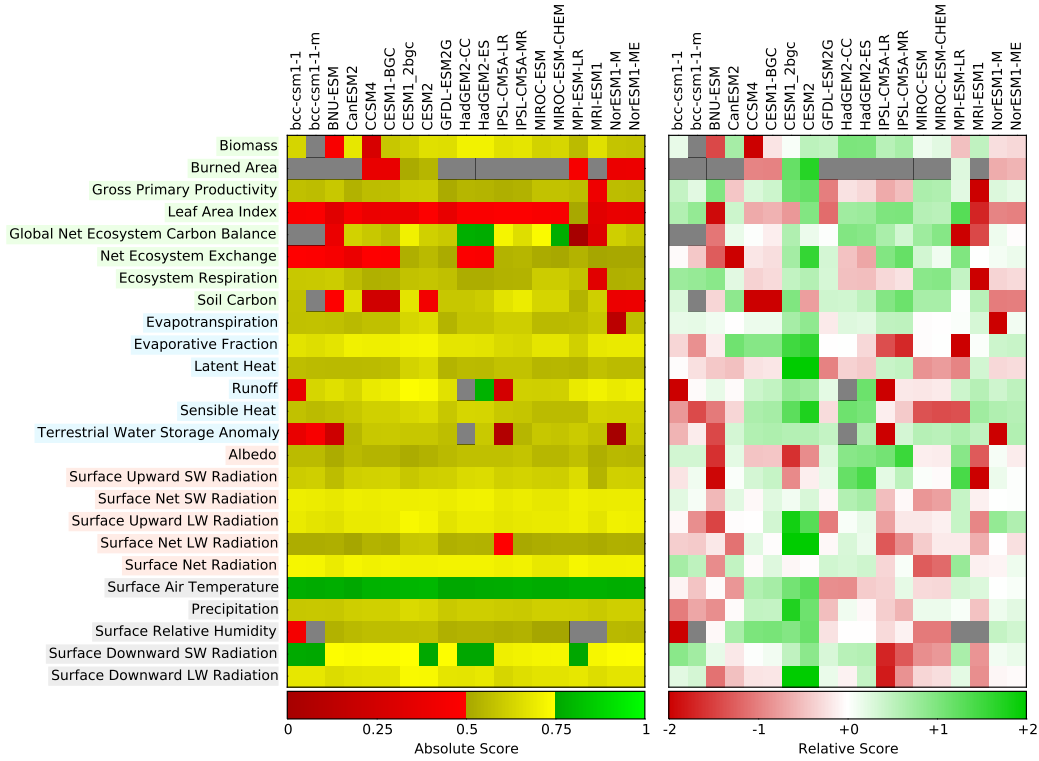


Figure 1: The ILAMB top-level graphic uses stop-light colors to show how different models or model versions (across the top) score with respect to each variable (down the left) in an absolute sense (left rectangle) and with respect to each other (right rectangle). Grey boxes reflect missing or unavailable data.

138 of this section, we discuss the analysis of a generalized variable $v(t, \mathbf{x})$ which we assume
139 represents a piecewise discontinuous function of constants in space and time. This means
140 that the temporal domain, represented by the variable t , is defined by the beginning and
141 ending of time intervals and the spatial domain, represented by the variable \mathbf{x} (bolded to
142 emphasize it is a vector quantity), represents the areas created by cell boundaries or the
143 areas associated with data sites. When necessary, we use the subscript ‘ref’ to reflect a
144 variable whose source is a reference or observational dataset, and the subscript ‘mod’ for
145 model datasets.

146 While many statistical quantities may be computed, the goal of our initial methodol-
147 ogy is to examine the mean state and variability around the mean over monthly to decadal
148 time scales and grid cell to global spatial scales. While we intend to uniformly apply this
149 analysis procedure to all variables, we also implement a mechanism to skip certain aspects
150 when deemed inappropriate. For example, if a reference dataset only contains average in-
151 formation across a span of years, the annual cycle is undefined and automatically skipped
152 in our implementation. The implementation also allows users to skip aspects of the analy-
153 sis that are deemed inappropriate even if it is possible to compute these metrics using the
154 available data. For example, the interannual variability may be poorly characterized in a
155 reference dataset even though the quantity could be computed.

156 **2.1 Preliminary Definitions**

157 Before presenting the specifics of the ILAMB methodology, we first present some
158 definitions used throughout the paper. While the following definitions are widely used in
159 the community, there are many subtle choices in their implementation that affect the inter-
160 pretation of the results. We present them here with precise meanings to emphasize where
161 a choice has been made and our reasoning for making it.

162 **2.1.1 Mean values over time**

163 When calculating mean values over the time period of the benchmark dataset, de-
164 noted by a bar superscribing the variable, we use the midpoint quadrature rule to approxi-
165 mate the integral,

$$\begin{aligned}\bar{v}(\mathbf{x}) &= \frac{1}{t_f - t_0} \int_{t_0}^{t_f} v(t, \mathbf{x}) dt \\ &\approx \frac{1}{T(\mathbf{x})} \sum_{i=1}^n v(t_i, \mathbf{x}) \Delta t_i\end{aligned}\quad (1)$$

166 where n represents the number of time intervals on which v is defined between the ini-
167 tial time, t_0 , and the final time, t_f , and Δt_i is the size of the i^{th} time interval, modified to
168 exclude time which falls outside of the integral limits,

$$\Delta t_i = \min(t_f, t_f^i) - \max(t_0, t_0^i) \quad (2)$$

169 where t_0^i and t_f^i are the initial and final time of each time interval. The average value is
170 obtained by dividing through by the amount of time in the interval, $t_f - t_0$, replaced in our
171 discrete approximation by the following function.

$$T(\mathbf{x}) = \sum_{i=1}^n \Delta t_i \quad \text{if } v(t_i, \mathbf{x}) \text{ is valid} \quad (3)$$

172 In words, Equation (3) addresses temporally discontinuous data by summing all the time
173 step interval sizes only if the corresponding variable data is marked as valid. This means
174 that if a function has some values masked or marked as invalid at some locations, we do
175 not penalize the averaged value by including this as a time at which a value is expected.
176 If an integral (or sum) is desired instead of an average, then we simply omit the division
177 by $T(\mathbf{x})$ in Equation (1).

178 **2.1.2 Mean values over space**

179 When computing spatial means over various regions of interest, denoted by a double
180 bar over a variable, we use the midpoint rule for integration to approximate the following
181 weighted spatial integral,

$$\begin{aligned}\bar{\bar{v}}(t) &= \frac{1}{\int_{\Omega} w(\mathbf{x}) d\Omega} \int_{\Omega} v(t, \mathbf{x}) w(\mathbf{x}) d\Omega \\ &\approx \frac{1}{A(\Omega)} \sum_{i=1}^{n(\Omega)} v(t, \mathbf{x}_i) w(\mathbf{x}_i) a_i\end{aligned}\quad (4)$$

182 over a region Ω , also referred to as a area-weighted mean. Here the function $w(\mathbf{x})$ is an
 183 optional generic weighting function defined over space. The summation is over $n(\Omega)$, that
 184 is the integer number of spatial cells whose centroids fall into the region of interest. A
 185 function evaluation at a location \mathbf{x}_i refers to the constant value which corresponds to that
 186 spatial cell. The value of a_i is the area of the cell, which could be some fraction of the
 187 total cell area if integrating over land in coastal regions. We then divide through by the
 188 measure, the sum of the grid areas with the weights,

$$A(\Omega) = \sum_{i=1}^{n(\Omega)} w(\mathbf{x}_i) a_i \quad \text{if } v(t, \mathbf{x}_i) \text{ is valid} \quad (5)$$

189 Note that if no weighting is required, this is a normalization by the sum of the area over
 190 which we integrate. As with the temporal mean, if an integral only is required, we simply
 191 omit the division by $A(\Omega)$. In cases where a mean over a collection of sites is needed, the
 192 spatial integral reduces to an arithmetic mean across the sites.

193 If we are spatially integrating a variable from a single source, then its spatial grid is
 194 clearly defined and Equation (4) can be directly applied to compute the quantity of inter-
 195 est. However, if the integrand involves quantities from two different sources, as in comput-
 196 ing the global bias or RMSE, then there is likely a disparity in both resolution and repre-
 197 sentation of land areas. We address resolution differences by interpolating both sources to
 198 a grid composed of the cell breaks, the location at which two neighboring cells meet, of
 199 both data sources. Consider two spatial grids whose cells are defined by the outer product
 200 of 1D vectors representing the cell breaks in spherical coordinates,

$$\mathcal{G}_1 := \theta_1 \otimes \varphi_1 \quad (6)$$

$$\mathcal{G}_2 := \theta_2 \otimes \varphi_2 \quad (7)$$

201 where θ refers to the latitude, φ to longitude, and \otimes a operator which creates a two-dimensional
 202 grid from one-dimensional vectors. We address differences in resolution by defining a
 203 composite grid which consists of the outer product of the union of these two grids' cell
 204 breaks,

$$\mathcal{G}_c := (\theta_1 \cup \theta_2) \otimes (\varphi_1 \cup \varphi_2). \quad (8)$$

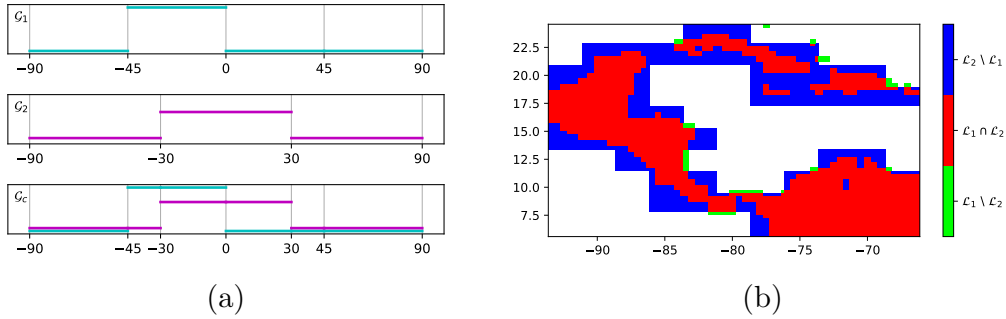


Figure 2: When comparing two spatial variables of varying resolution, we interpolate both to a common grid composed of the cell breaks of both variables over the intersection of what both variables agree is land. (a) Interpolation of sample step functions defined on grids \mathcal{G}_1 and \mathcal{G}_2 both interpolated to a composite grid \mathcal{G}_c using nearest neighbor interpolation with zero interpolation error. The vertical grid lines reflect the cell boundaries in each grid., (b) Differences in the representation of land from a reference and model dataset zoomed into Central America for emphasis. The red region represents where both sources are in agreement, the blue is land for the model but not the reference and the green is land for the reference but not the model.

205 Once constructed, quantities defined on both \mathcal{G}_1 and \mathcal{G}_2 may be interpolated to \mathcal{G}_c by
 206 nearest neighbor interpolation with zero interpolation error due to the nested nature of
 207 the grids. This can be seen visually by comparing the three plots shown in Figure 2(a). In
 208 each plot, the tick marks along the x-axis represent the cell breaks of the particular one-
 209 dimensional grid left coarse for illustration. The cyan curve represents a step function de-
 210 fined on the grid of a reference dataset \mathcal{G}_1 and the magenta curve on that of the model
 211 dataset \mathcal{G}_2 . Both are interpolated to the composed grid \mathcal{G}_c without loss of information, al-
 212 beit on a new grid containing more cells of variable size. Once on a composite grid, the
 213 quantities may be compared directly. As the ILAMB methodology has been envisioned for
 214 comparisons with model output from CMIP5, we have made an implicit assumption that
 215 each source grid, \mathcal{G}_1 and \mathcal{G}_2 , is regular and can be represented by one-dimensional vec-
 216 tors. While the implementation does provide naive interpolation for non-regular grids, the
 217 user is encouraged to employ a conservative interpolation scheme of their choosing prior
 218 to applying the ILAMB methodology.

219 In addition to resolution differences, we observe that data sources vary in the under-
 220 lying representation of the distinction between land and water. We illustrate this concept

221 in Figure 2(b) where we compare a fine scale representation of land \mathcal{L}_1 to a relatively
 222 coarse representation \mathcal{L}_2 . This is a typical situation encountered when comparing high
 223 resolution observational data to lower resolution model output. The red region represents
 224 the intersection of land areas $\mathcal{L}_1 \cap \mathcal{L}_2$, that is, where both sources report the presence of
 225 land. However, there are missed land areas from both sources, represented by the blue and
 226 green colors. As much of the disagreement over what is considered land occurs around
 227 islands in tropical regions (for example Central America and Equatorial Asia), these non-
 228 represented areas can constitute a nontrivial percentage of the total represented variable
 229 v .

230 For transparency, the ILAMB implementation is built with the capability of report-
 231 ing integrals over each of these three land areas. Unless specifically stated otherwise,
 232 when spatially integrating a quantity from a single source, we use the original grid and
 233 land areas given by that source. This is to remain as true to the original intent of the
 234 provider as we can. However, when comparing two data sources of varying resolution and
 235 land representation, we perform this integration over what both report to be land, $\mathcal{L}_1 \cap \mathcal{L}_2$
 236 (the red area in Figure 2(b)).

237 **2.1.3 Computing normalized scores from errors**

238 In the following sections 2.2 and 2.3, we detail how we compute errors and trans-
 239 form them into normalized scores on the unit interval. This approach is intended to syn-
 240 thesize model performance across a range of dimensions with respect to a given dataset.
 241 We achieve this by taking a measure of the relative error, generically represented here as
 242 ε , and passing it through the exponential function,

$$s = e^{-\alpha\varepsilon} \quad (9)$$

243 where s is a score on the interval $[0, 1]$ and α is a parameter which can be used to tune
 244 the mapping of error to score. The classic expression of relative error is prone to numer-
 245 ical instabilities for denominator values near or which cross zero. Furthermore the mag-
 246 nitude of the error can depend on the units selected. For this reason we depart from the
 247 standard definition of relative error and develop specialized expressions in Equations (13,
 248 18, 26).

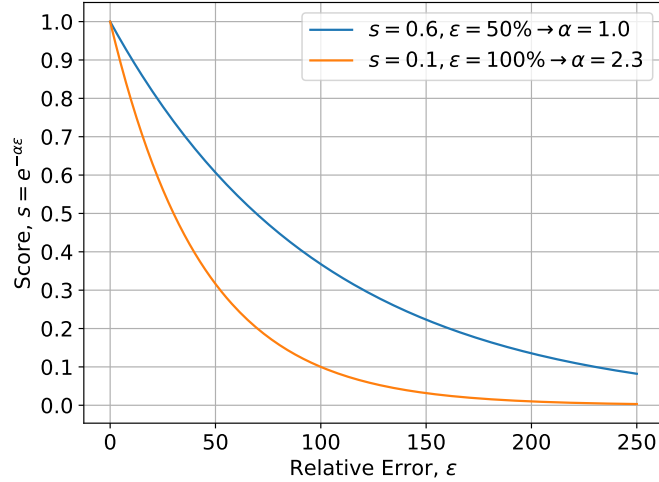


Figure 3: Mapping function of relative error ϵ to a score s on the unit interval. Two choices of α are shown: $\alpha = 1$, shown in blue, which equates a score of 0.6 to a relative error of 50%, and $\alpha = 2.3$, shown in orange, which equates a score of 0.1 to a relative error of 100%.

249 While the choice of the exponential function is arbitrary, it was chosen because
 250 it maps zero error to a score of one and smoothly reduces the score as the error grows,
 251 never reaching exactly zero. This is important as we want to improve the score when the
 252 error improves, no matter how large of error we observe. If the user wants a relative error
 253 of $\hat{\epsilon}$ to equate to a score of \hat{s} , then

$$\alpha = -\frac{\ln(\hat{s})}{\hat{\epsilon}} \quad (10)$$

254 In Figure 3 we plot this function with two choices for α , which illustrates how the relative
 255 error may be controlled. Unless stated otherwise, we use an implicit $\alpha = 1$ throughout the
 256 manuscript.

257 2.2 Mean State Analysis

258 In this section, we describe the various metrics and plots that our methodology gener-
 259 erates. While presented in terms of the abstract variable v , we also include sample plots
 260 of a comparison of the GBAF (*Jung et al., 2010*) gross primary productivity (GPP) with
 261 CLM4.5 (*Oleson et al., 2013*) for the purpose of illustration. In practice, ILAMB pro-

262 duces thousands of such plots and scalars, which are browsable in a website designed to
263 aid modelers in understanding the benchmarking results.

264 **2.2.1 Bias**

265 We find the mean value in time, $\overline{v_{\text{ref}}(\mathbf{x})}$, over the time period of the reference, as
266 well as that of the model, $\overline{v_{\text{mod}}(\mathbf{x})}$, over the same time period. These are spatial variables
267 that are included in the standard output as plots, as shown in Figure 4(a-b). We also com-
268 pute the bias,

$$\text{bias}(\mathbf{x}) = \overline{v_{\text{mod}}(\mathbf{x})} - \overline{v_{\text{ref}}(\mathbf{x})} \quad (11)$$

269 as well as its mean over a given region, $\overline{\overline{\text{bias}(\mathbf{x})}}$. To score the bias, we need to non-dimensionalize
270 it as a relative error. We have chosen to do this by using the centralized root mean square
271 of the reference data,

$$\text{crms}(\mathbf{x}) = \sqrt{\frac{1}{t_f - t_0} \int_{t_0}^{t_f} (v_{\text{ref}}(t, \mathbf{x}) - \overline{v_{\text{ref}}(\mathbf{x})})^2 dt}, \quad (12)$$

272 which makes the relative error in bias given as,

$$\varepsilon_{\text{bias}}(\mathbf{x}) = |\text{bias}(\mathbf{x})| / \text{crms}(\mathbf{x}) \quad (13)$$

273 where the $|\cdot|$ operator represents the absolute value. The bias score as a function of space
274 is,

$$s_{\text{bias}}(\mathbf{x}) = e^{-\varepsilon_{\text{bias}}(\mathbf{x})} \quad (14)$$

275 and the scalar score

$$S_{\text{bias}} = \overline{\overline{s_{\text{bias}}(\mathbf{x})}}, \quad (15)$$

276 that is, the spatially integrated bias score. The motivation behind Equation (13) is to nor-
277 malize the bias by the variability at any given spatial location. However, this also leads
278 to the consequence that in areas where the given variable v has a small magnitude, sim-
279 ple noise can lead to large relative errors. For example, in Figure 4(d) we observe a poor

280 score in the dry regions of Australia where GPP is small. Given the small contribution,
 281 it is undesirable that these errors induce a large negative contribution to the overall score.
 282 To address this issue, we introduce the concept of *mass weighting*. That is, when perform-
 283 ing the spatial integral to obtain a scalar score (Equation (15)), we weight the integral by
 284 the period mean value of the reference variable using Equation (4) with $w = \overline{v_{\text{ref}}}$. In some
 285 instances the variable is truly a mass, but other times a flux or rate. The main motivation
 286 is to weight in areas where the variable is active. So while in our conceptual example,
 287 there is large relative error in GPP over deserts, these values will not negatively contribute
 288 to the overall score as the value of GPP is low in this area.

289 We apply mass weighting when the variable v represents a mass or flux of carbon
 290 or water as in GPP or precipitation. For variables representing energy states or quantities,
 291 such as temperature and radiation, we omit the weighting and perform a spatial integral
 292 only. We report plots of the bias and its score as well as the scalar integrated mean val-
 293 ues.

294 2.2.2 Root mean squared error

295 For reference datasets with seasonal and interannual variability, we compute the root
 296 mean squared error over the time period of the reference dataset,

$$rmse(\mathbf{x}) = \sqrt{\frac{1}{t_f - t_0} \int_{t_0}^{t_f} (v_{\text{mod}}(t, \mathbf{x}) - v_{\text{ref}}(t, \mathbf{x}))^2 dt} \quad (16)$$

297 and include plots and the scalar $\overline{rmse(\mathbf{x})}$ in the standard output (Figure 5(a)). To score the
 298 root mean square error, we normalize the centralized root mean square error,

$$crmse(\mathbf{x}) = \sqrt{\frac{1}{t_f - t_0} \int_{t_0}^{t_f} ((v_{\text{mod}}(t, \mathbf{x}) - \overline{v_{\text{mod}}(\mathbf{x})}) - (v_{\text{ref}}(t, \mathbf{x}) - \overline{v_{\text{ref}}(\mathbf{x})}))^2 dt} \quad (17)$$

299 by the centralized root mean square of the reference dataset, Equation (12). This leads to
 300 a relative error of

$$\varepsilon_{rmse}(\mathbf{x}) = crmse(\mathbf{x})/crms(\mathbf{x}) \quad (18)$$

301 and a spatial RMSE score

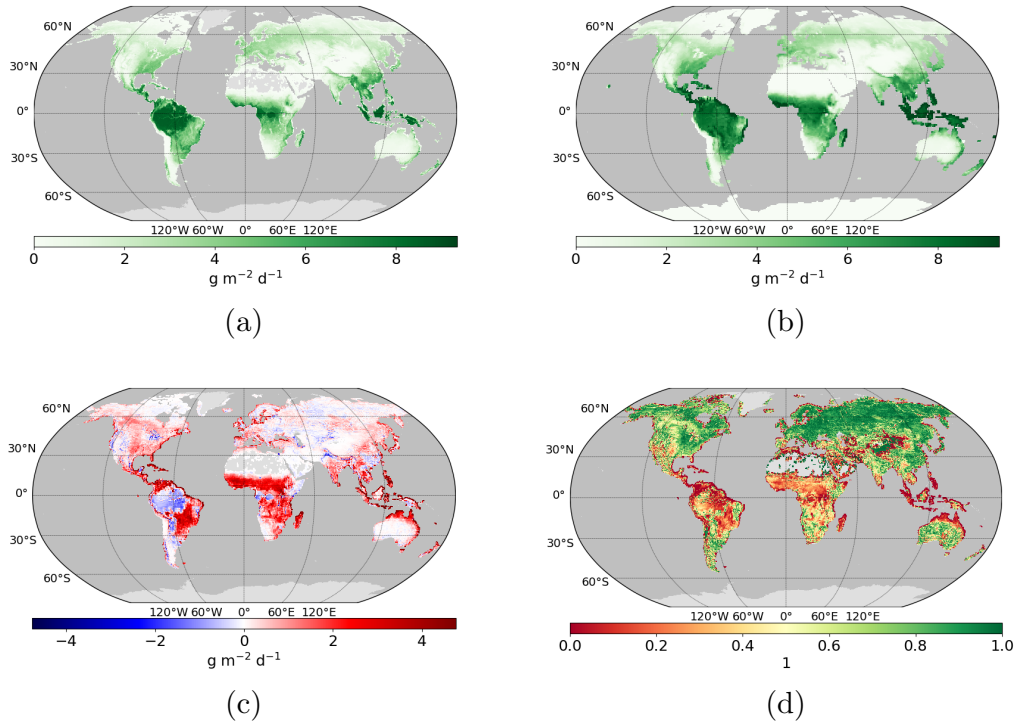


Figure 4: Comparisons of GPP between the reference (GBAF) and the model (CLM4.5) dataset. Each period mean is plotted over the original grid of the dataset. We highlight here that the reference (a) is not defined over Antarctica, Greenland, and part of the Sahara desert whereas the model (b) is defined over all land areas. Yet when the bias (c) and its score (d) is reported, the area represented is what both the reference and model agree on as land. (a) Reference period mean, $\overline{v_{\text{ref}}}(\mathbf{x})$, (b) Model period mean, $\overline{v_{\text{mod}}}(\mathbf{x})$, (c) Bias, $\text{bias}(\mathbf{x})$, (d) Bias Score, $s_{\text{bias}}(\mathbf{x})$

$$s_{rmse}(\mathbf{x}) = e^{-\varepsilon_{rmse}(\mathbf{x})}. \quad (19)$$

302 The scalar score is obtained by

$$S_{rmse} = \overline{s_{rmse}}(\mathbf{x}), \quad (20)$$

303 where we again employ mass weighting when necessary. We score the centralized root
 304 mean squared error to decouple the bias score from the RMSE score. Computing the
 305 RMSE score by normalizing the RMSE would lead to a double counting of errors. That
 306 is, a large error in bias also leads to a large error in RMSE. By scoring the centralized
 307 RMSE, we remove the bias from the RMSE, allowing the RMSE score to focus on an or-
 308 thogonal aspect of model performance.

309 2.2.3 Phase Shift

310 We evaluate the phase shift of the annual cycle of many datasets that have monthly
 311 variability by comparing the timing of the maximum of the annual cycle of the variable,
 312 $c(v)$ at each spatial cell across the time period of the reference dataset. We then approxi-
 313 mate the phase shift of the reference and model datasets by subtracting these two values,

$$\theta(\mathbf{x}) = \arg \max_t (c_{\text{mod}}(t, \mathbf{x})) - \arg \max_t (c_{\text{ref}}(t, \mathbf{x})) \quad (21)$$

314 expressed in days. As the units for phase shift are consistent across all variables, no nor-
 315 malization is needed and we can remap the shift to the unit interval by

$$s_{\text{phase}}(\mathbf{x}) = \frac{1}{2} \left(1 + \cos \left(\frac{2\pi\theta(\mathbf{x})}{365} \right) \right) \quad (22)$$

316 and then spatially integrate the score over the appropriate region to find the scalar score,

$$S_{\text{phase}} = \overline{s_{\text{phase}}}(\mathbf{x}), \quad (23)$$

317 where again mass weighting is employed when appropriate. We include plots of the phase
 318 shift and its score in the standard output and represent them here in Figure 5(c-d). In ad-
 319 dition to plots which show the time averaged variables as a map, we include line plots of

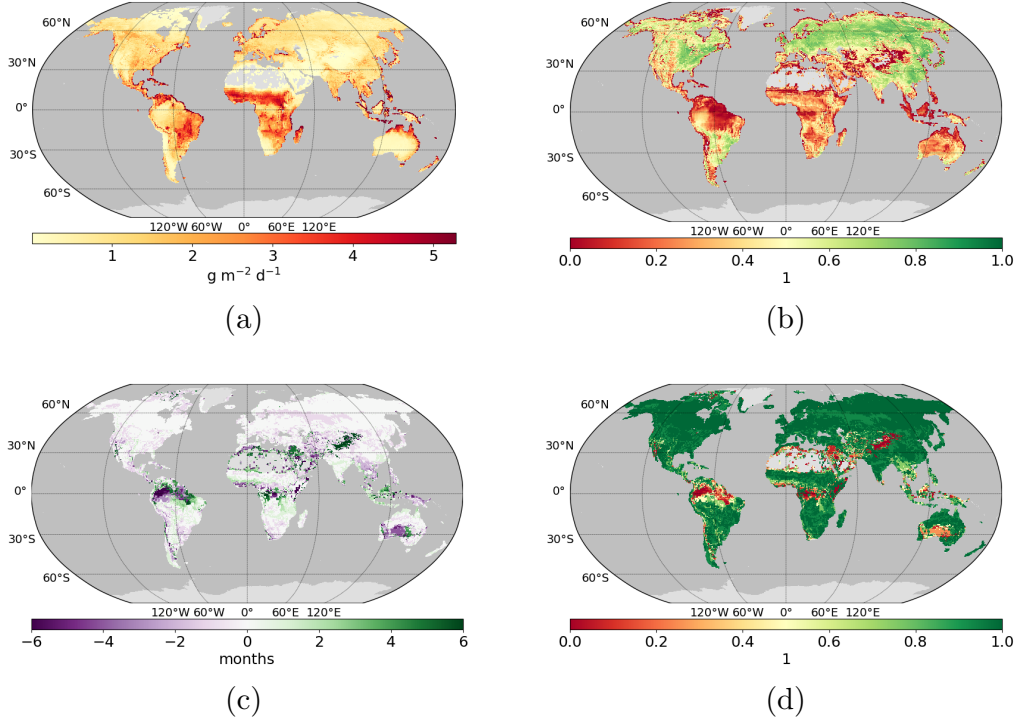


Figure 5: Comparisons of the RMSE and phase of GPP between the reference (GBAF) and the model (CLM4.5) dataset. (a) RMSE, $rmse(\mathbf{x})$, (b) RMSE score, $s_{rmse}(\mathbf{x})$, (c) Phase shift, $\theta(\mathbf{x})$, (d) Phase shift score, $s_{cycle}(\mathbf{x})$

320 the mean annual cycle and the spatially averaged variables, $\overline{v_{ref}}(t)$ and $\overline{v_{mod}}(t)$ shown in
 321 Figure 6.

322 2.2.4 Interannual Variability

323 A score for the interannual variability is computed by removing the annual cycle
 324 from both the reference and the model,

$$iav_{ref}(\mathbf{x}) = \sqrt{\frac{1}{t_f - t_0} \int_{t_0}^{t_f} (v_{ref}(t, \mathbf{x}) - c_{ref}(t, \mathbf{x}))^2 dt} \quad (24)$$

$$iav_{mod}(\mathbf{x}) = \sqrt{\frac{1}{t_f - t_0} \int_{t_0}^{t_f} (v_{mod}(t, \mathbf{x}) - c_{mod}(t, \mathbf{x}))^2 dt} \quad (25)$$

$$\varepsilon_{iav}(\mathbf{x}) = (iav_{mod}(\mathbf{x}) - iav_{ref}(\mathbf{x})) / iav_{ref}(\mathbf{x}) \quad (26)$$

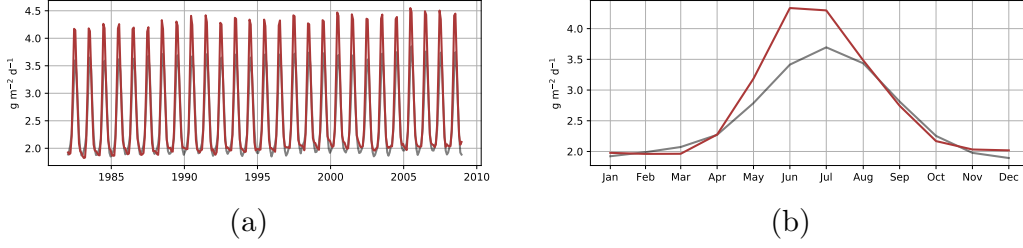


Figure 6: Spatial means of GPP of the reference (GBAF) shown in grey and the model (CLM4.5) in maroon. (a) Spatially integrated mean, $\overline{\overline{v_{\text{ref}}(t)}}$ and $\overline{\overline{v_{\text{mod}}(t)}}$, (b) Mean annual cycle, $\overline{\overline{v_{\text{ref}}(t)}}$ and $\overline{\overline{v_{\text{mod}}(t)}}$

326 and then computing a score as a function of space,

$$s_{iav}(\mathbf{x}) = e^{-\varepsilon_{iav}(\mathbf{x})}. \quad (27)$$

327 The scalar score is then obtained by

$$S_{iav} = \overline{\overline{s_{iav}(\mathbf{x})}}, \quad (28)$$

328 where mass weighting is used when necessary. We include plots of the variability and
 329 the score in the standard output and show them here in Figure 7. Note that while here
 330 we have shown the interannual variability of the GBAF product for illustration, in the de-
 331 fault ILAMB configuration, the interannual variability is currently omitted for the GBAF
 332 products because its representativeness is considered to be poor (see Figure 10 of (Kumar
 333 *et al.*, 2016)).

334 2.2.5 Spatial Distribution

335 We score the spatial distribution of the time averaged variable by generating a Tay-
 336 lor (Taylor, 2001) diagram. We do this by computing the normalized standard deviation,

$$\sigma = \frac{\text{stdev}(\overline{\overline{v_{\text{mod}}(\mathbf{x})}})}{\text{stdev}(\overline{\overline{v_{\text{ref}}(\mathbf{x})}})} \quad (29)$$

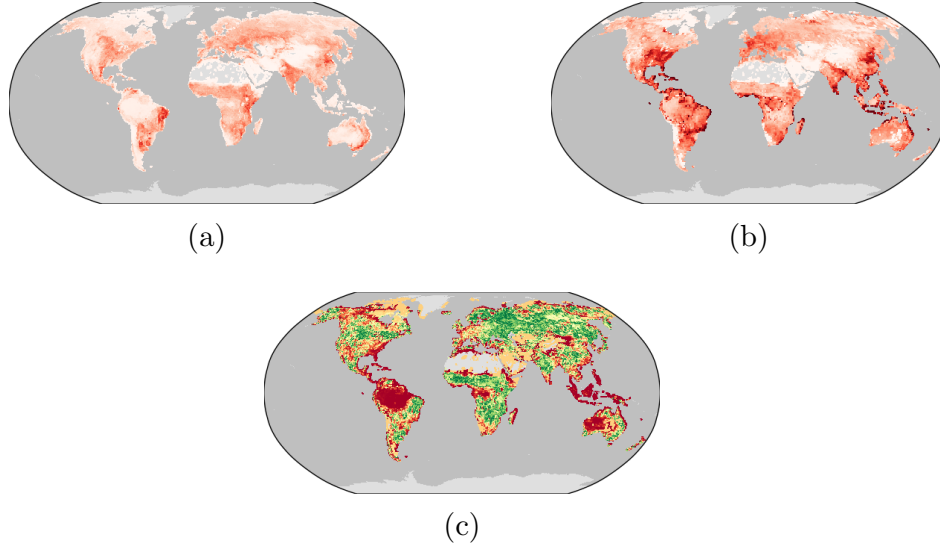


Figure 7: Comparisons of the interannual variability of GPP between the reference (GBAF) and the model (CLM4.5) dataset. (a) Reference interannual variability, $iav_{ref}(\mathbf{x})$, (b) Model interannual variability, $iav_{mod}(\mathbf{x})$, (c) Interannual variability score, $s_{iav}(\mathbf{x})$

337 and the spatial correlation R of the period mean values $\bar{v}_{ref}(\mathbf{x})$ and $\bar{v}_{mod}(\mathbf{x})$, and then as-
 338 signing a score by the following relationship

$$S_{dist} = \frac{2(1 + R)}{(\sigma + \frac{1}{\sigma})^2}, \quad (30)$$

339 where the main idea is that we penalize the score when R and σ deviate from a value of
 340 1. We include the Taylor plot in the standard output and represent it here in Figure 8.

341 2.2.6 Overall Score

342 The overall score for a given variable and data product is a composite of the suite of
 343 metrics defined above. We use a weighted sum,

$$S_{overall} = \frac{S_{bias} + 2S_{rmse} + S_{phase} + S_{iav} + S_{dist}}{1 + 2 + 1 + 1 + 1}, \quad (31)$$

344 where the RMSE score is doubly weighted to emphasize its importance.

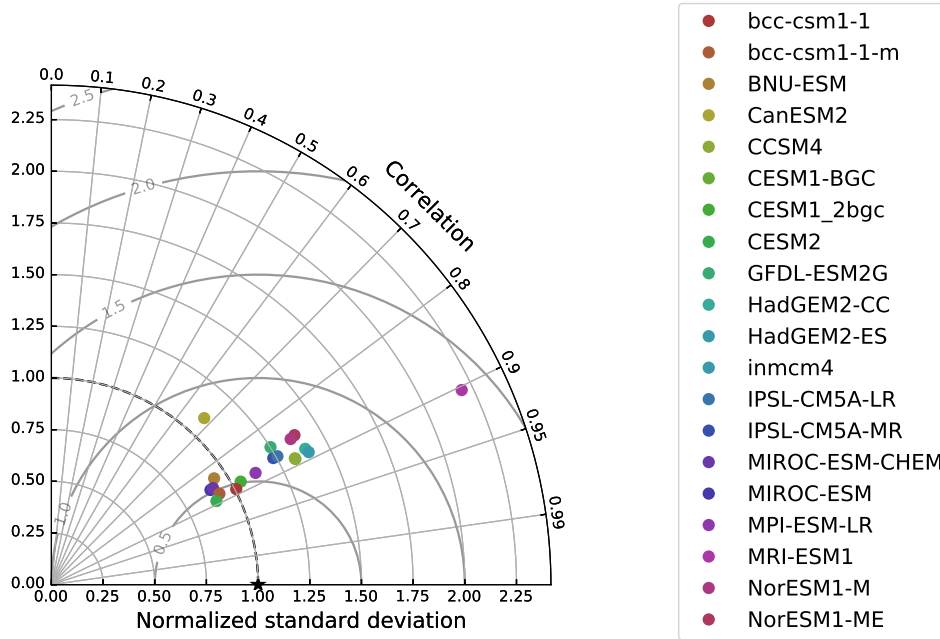


Figure 8: Taylor diagram comparing the spatial distribution of GPP of the reference (GBAF) shown as a black star to the CMIP5 models shown in colors.

2.3 Relationship Analysis

As models are frequently calibrated using the mean state scalar measures described in Section 2.2, a higher score does not necessarily reflect a more process-oriented model. In order to assess the representation of mechanistic processes in models, we also evaluate variable-to-variable relationships. For example, we look at how well models represent the relationship that GPP has with precipitation, evapotranspiration, and temperature. For the purposes of this section, we represent a generic dependent variable as v , as before, and score its relationship with an independent variable u . We then quantify the variable-to-variable relationship of the time period mean, $\bar{u}(\mathbf{x})$ on $\bar{v}(\mathbf{x})$, derived from the combination of reference datasets to the relationship diagnosed in models. We use the mean values over the reference time period to establish relationships as they represent a logical starting point. In the future, we plan to extend the relationship analysis to include seasonal and interannual variability.

2.3.1 Functional Response

We estimate a functional response by a 1D histogram, binned in terms of the independent variable $\bar{u}(\mathbf{x})$ with a number of bins, initially set to $n_{\text{bins}} = 25$. Then in each bin, we compute the mean value of the corresponding dependent variable, $\bar{v}(\mathbf{x})$ to approximate the functional dependence of u on v . We represent this binning with the operator \mathcal{F} that operates on the dependent and independent variables. We use it to compute functions from both the reference and model datasets.

$$f_{\text{ref}}(u) = \mathcal{F}(\bar{v}_{\text{ref}}(\mathbf{x}), \bar{u}_{\text{ref}}(\mathbf{x})) \quad (32)$$

$$f_{\text{mod}}(u) = \mathcal{F}(\bar{v}_{\text{mod}}(\mathbf{x}), \bar{u}_{\text{mod}}(\mathbf{x})), \quad (33)$$

where both curves are plotted in Figure 9(a) for the case of GPP compared to surface air temperature. These response curves are then scored by computing a relative error based on the RMSE,

$$\varepsilon_{\text{func}}^u = \sqrt{\frac{\int (f_{\text{ref}}(u) - f_{\text{mod}}(u))^2 du}{\int f_{\text{ref}}(u)^2 du}}, \quad (34)$$

where the integrals are approximated by the midpoint rule over the bins of the independent variable $\bar{u}(\mathbf{x})$. Then we use Equation (9) to map this relative error to a score by,

$$S_{\text{func}}^u = e^{-\varepsilon_{\text{func}}^u}. \quad (35)$$

The superscript u reinforces that this score represents functional performance with respect to a given independent variable u . The ILAMB implementation allows for any number of independent variables to be studied. In terms of our sample, ILAMB scores the functional relationship of GPP with respect to each independent variable separately (precipitation, evapotranspiration, temperature, *etc.*) and then computes the mean of these scores for the overall relationship score.

2.3.2 Hellinger Distance

In addition to the one-dimensional histograms, we also build normalized two-dimensional histograms ($n_{\text{bins}} = 25$ in both dimensions) from the time averaged data $\bar{v}(\mathbf{x})$ and $\bar{u}(\mathbf{x})$, represented here by the operator \mathcal{D} . We represent these distributions by,

$$d_{\text{ref}}(u) = \mathcal{D}(\overline{v_{\text{ref}}}(\mathbf{x}), \overline{u_{\text{ref}}}(\mathbf{x})), \quad (36)$$

$$d_{\text{mod}}(u) = \mathcal{D}(\overline{v_{\text{mod}}}(\mathbf{x}), \overline{u_{\text{mod}}}(\mathbf{x})), \quad (37)$$

380 as depicted in Figure 9(b–c). If we represent individual elements from these distributions
 381 $d_{\text{ref}}(u) = (p_1, \dots, p_{n_{\text{bins}}^2})$ and $d_{\text{mod}}(u) = (q_1, \dots, q_{n_{\text{bins}}^2})$, we can compute the so-called Hellinger
 382 distance (Law *et al.*, 2015)

$$S_{\text{dist}}^u = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^{n_{\text{bins}}^2} (\sqrt{p_i} - \sqrt{q_i})^2}$$

383 as a measure of how similar two distributions are to each other. While there are other
 384 choices, such as the Kullback-Leibler divergence, which are more commonly employed
 385 (Dirmeyer *et al.*, 2014), the Hellinger distance comes with the added benefit of being al-
 386 ready normalized [0, 1] and thus, further normalization is not necessary to use this directly
 387 as a score.

388 However, we only report the Hellinger distance as a scalar and do not include it
 389 in the scoring of the relationships. This is due to the fact that a bias in an independent
 390 variable can cause a density shift in the 2D distribution that would cause the score to un-
 391 reasonably decrease. In terms of our example, a bias in precipitation (e.g. arising from
 392 a coupled model) could result in a poor relationship score with GPP, even if there is no
 393 underlying deficiency in the land-model simulated precipitation versus GPP relationship.

394 **3 Datasets**

395 In this section we explain how we utilize the methodology presented in Section 2 to
 396 evaluate model performance with respect to a collection of datasets (Tables 2–5) assem-
 397 bled by the ILAMB community. Errors in measurements, lack of measured or reported
 398 uncertainties, and inconsistencies in measurement methodology or instrumentation leading
 399 to ambiguous confidence in derived or synthesized data products all represent challenges
 400 in using observational data for benchmarking. In addition, the spatial and temporal cover-
 401 age of different data products can vary substantially.

402 To account for the lack of quantitative uncertainties and scale mismatches between
 403 observations and models, and to bring a quantitative objectivity to model–data compari-

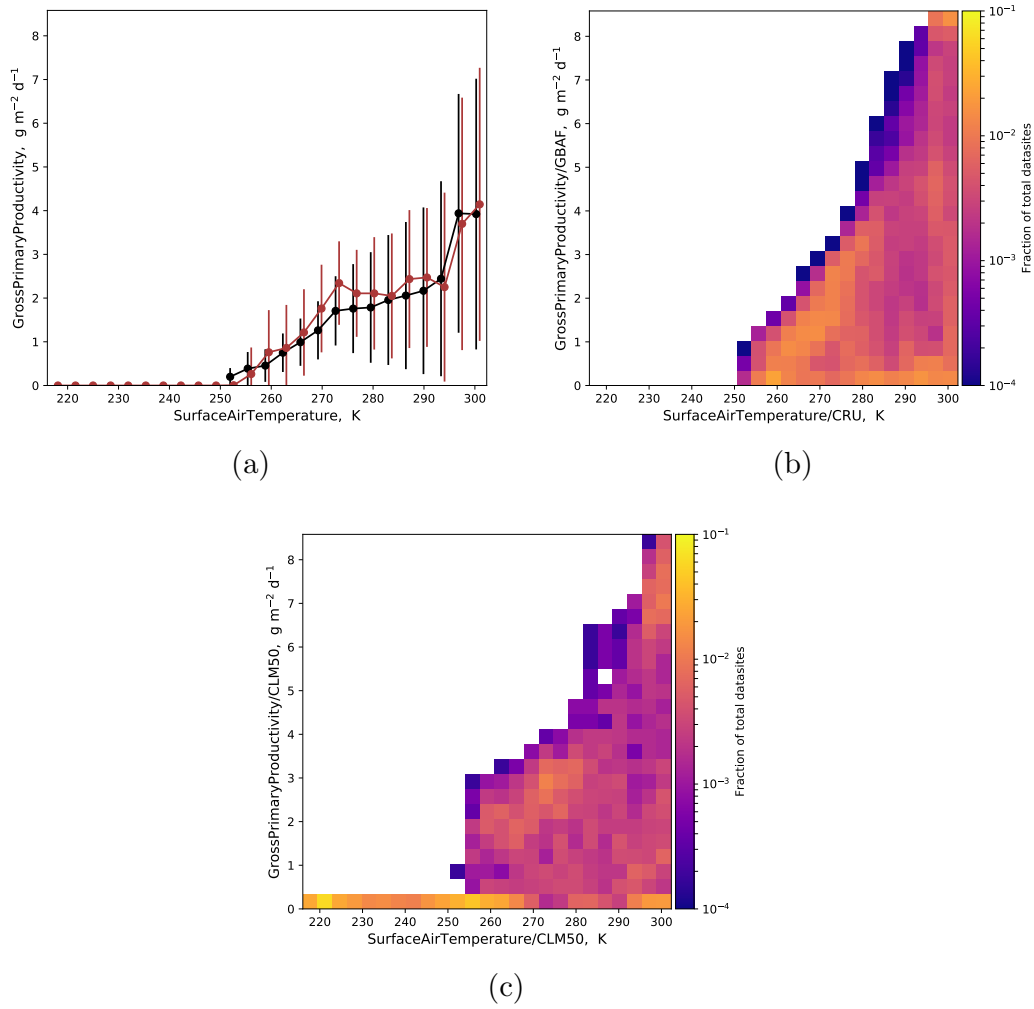


Figure 9: Variable-to-variable relationship plots which are a part of the standard output from the ILAMB methodology. (a) Functional responses, the reference $f_{ref}(u)$ in black, and the model $f_{mod}(u)$ in maroon. Data points reflect the mean for each independent value and the error bars reflect the standard deviation range., (b) Reference distribution, $d_{ref}(u)$, (c) Model distribution, $d_{mod}(u)$

404 son, we developed a three-element rubric for weighting datasets as represented in Table 1.
 405 The first weight is based on a qualitative estimate of the certainty we have in a particu-
 406 lar dataset. This weight encompasses both our certainty in the process used to obtain the
 407 observational information as well as the presence of quantitative uncertainty in the mea-
 408 surements themselves. A second weight for each dataset reflects its spatial and temporal
 409 coverage. The datasets employed in ILAMB are diverse and include site-level data, re-
 410 analysis data products, and remotely-sensed data. As our aim is to provide insight in land
 411 model performance on global and decadal scales, we give more weight to global products
 412 which are time series that extend for several years. The weights are combined multiplica-
 413 tively to assign a total weight for each dataset. Then we normalize the weight by the sum
 414 of the weights of all the datasets for a given variable. For example, from Table 2 we see
 415 that there are two datasets used to benchmark GPP: Fluxnet and GBAF. For the Fluxnet
 416 product, we assign a certainty weight of 3 because while the collection is discussed in
 417 the published literature, there is no quantitative uncertainty provided. We assign a scale
 418 weight of 3 because the collection of sites covers multiple years of a substantial region of
 419 the globe yet has sparse coverage over important regions such as the tropics. The GBAF
 420 product is assigned a certainty weight of 3 for the same reason and a scale weight of 5
 421 as it provides global coverage spanning multiple years. Then the total weight for the GPP
 422 variable which the GBAF dataset carries is

$$w_{\text{GBAF}}^{\text{GPP}} = \frac{3 \cdot 5}{3 \cdot 3 + 3 \cdot 5} \approx 63\%.$$

423 We use these weights to blend the overall score (Equation (31)) from each dataset for each
 424 variable. In this way ILAMB remains flexible to adding datasets as they are developed,
 425 allowing more weight to be given to those that the community believes are more credible
 426 and that are more comparable in scale to models.

427 A third weight reflects how useful the measured variable is in the focus of a model
 428 intercomparison project. Here, as an example, we show weighting for an analysis of model
 429 performance in representing the carbon cycle. We use these weights to blend the overall
 430 scores from each variable into a complete score across all variables for a given model.
 431 This allows ILAMB to include comparisons that are important for a complete understand-
 432 ing of the carbon cycle without necessarily allowing them to heavily influence the overall
 433 score. For example, the radiation and energy cycle datasets in Table 4 are all weighted

Table 1: The ILAMB rubric used to assign relative weights of a dataset. A score for each dataset is assigned in each of three areas. These scores are then combined multiplicatively and used to determine relative importance for a dataset with respect to a given variable.

Score	Certainty	Scale	Process
1	No given uncertainty, significant methodological issues affecting quality	Site level observations with limited space/time coverage	Observations that have limited influence on the targeted Earth system dynamics
2	No given uncertainty, some methodological issues affecting quality	Partial regional coverage, up to 1 year	Observations have direct influence on the targeted Earth system dynamics
3	No given uncertainty, methodology has some peer review	Regional coverage, at least 1 year	Observations useful to constrain processes that contribute to the targeted Earth system dynamics
4	Qualitative uncertainty, methodology accepted	Important regional coverage, at least 1 year	Observations well-suited to constrain important processes
5	Well-defined and relatively low uncertainty	Global scale spanning multiple years	Observations well-suited for discriminating critical processes among models

434 comparatively low because, while they help one understand the carbon cycle, they are not
435 as influential in the overall behavior.

436 We emphasize that this rubric is particular to our overarching goal of understanding
437 the carbon cycle on global and decadal scales. However, the implementation is flexible
438 and allows for an arbitrary weighting scheme to be developed that suits the needs of the
439 user, community, or model intercomparison project that it serves.

440 The references and weights for each dataset that we have selected may be found in
 441 Tables 2–5. Each table represents a different aspect of the model: the ecosystem and car-
 442 bon cycle in Table 2, the hydrological cycle in Table 3, the radiation and energy cycle in
 443 Table 4, and the forcings in Table 5. For the majority of these datasets, we make a direct
 444 comparison of the observed quantity to model outputs, or algebraic combinations of model
 445 outputs using the methodology described in section 2. However, there are a few special
 446 cases which require specific handling which we describe in the next section.

447 3.1 Special cases

448 In general, a consistent methodology is applied to compare model output with each
 449 dataset. This consistency across variables and datasets is a strength of the ILAMB method-
 450 ology. However, this is not always possible, and here we enumerate a few exceptions and
 451 how they are handled.

452 3.1.1 Evaporative Fraction

453 To test the partitioning of surface energy, we compare the evaporative fraction de-
 454 rived from the GBAF (*Jung et al.*, 2010) data product to that of the models. The evapora-
 455 tive fraction is an algebraic expression in terms of the latent heat $L_e(t, \mathbf{x})$ and the sensible
 456 heat $S_h(t, \mathbf{x})$, given as

$$457 \quad ef(t, \mathbf{x}) = \frac{L_e(t, \mathbf{x})}{L_e(t, \mathbf{x}) + S_h(t, \mathbf{x})}. \quad (38)$$

458 The expression can cause nonsensical results because in winter, the sensible heat flux can
 459 be negative, leading to a change of sign in the evaporative fraction. The expression can
 460 also lead to large evaporative fraction values since the magnitudes of both the latent and
 461 sensible heat can become small. For this reason, we apply a mask to ef , L_e , and S_h only
 462 considering values for which $S_h > 0$, $L_e > 0$, and $S_h + L_e > \phi$, where $\phi = 20 [Wm^{-2}]$ is a
 463 surface energy threshold.

464 Equation (38) is used to study how models partition the surface energy throughout
 465 the relevant season. Thus we use that expression when computing the RMSE or seasonal
 466 cycle. However, when comparing period mean values and the bias, Equation (38) leads to
 467 a combination of averaging methods. For this reason, when computing the mean evapora-

Table 2: References and weighting of datasets used to measure the ecosystem and carbon cycle. Weights are chosen using the rubric in Table 1 and reflect a focus on understanding the carbon cycle.

Variable/Dataset	Certainty	Scale	Process
Biomass			5
Tropical (<i>Saatchi et al., 2011</i>)	4	4	
NBCD2000 (<i>Kellndorfer et al., 2013</i>)	4	2	
USForest (<i>Blackard et al., 2008</i>)	4	2	
BurnedArea			4
GFED4S (<i>Giglio et al., 2010</i>)	4	5	
GrossPrimaryProductivity			5
Fluxnet (<i>Lasslop et al., 2010</i>)	3	3	
GBAF (<i>Jung et al., 2010</i>)	3	5	
LeafAreaIndex			3
AVHRR (<i>Myneni et al., 1997</i>)	3	5	
MODIS (<i>De Kauwe et al., 2011</i>)	3	5	
GlobalNetEcosystemCarbonBalance			5
GCP (<i>Le Quéré et al., 2016</i>)	4	5	
Hoffman (<i>Hoffman et al., 2014</i>)	4	5	
NetEcosystemExchange			5
Fluxnet (<i>Lasslop et al., 2010</i>)	3	3	
GBAF (<i>Jung et al., 2010</i>)	2	2	
EcosystemRespiration			4
Fluxnet (<i>Lasslop et al., 2010</i>)	2	3	
GBAF (<i>Jung et al., 2010</i>)	2	2	
SoilCarbon			5
HWSD (<i>Todd-Brown et al., 2013</i>)	3	5	
NCSCDV22 (<i>Hugelius et al., 2013</i>)	3	4	

Table 3: References and weighting of datasets used to measure the hydrology cycle. Weights are chosen using the rubric in Table 1 and reflect a focus on understanding the carbon cycle.

Variable/Dataset	Certainty	Scale	Process
Evapotranspiration			5
GLEAM (<i>Miralles et al., 2011</i>)	3	5	
MODIS (<i>De Kauwe et al., 2011</i>)	3	5	
EvaporativeFraction			5
GBAF (<i>Jung et al., 2010</i>)	3	3	
LatentHeat			5
Fluxnet (<i>Lasslop et al., 2010</i>)	3	1	
GBAF (<i>Jung et al., 2010</i>)	3	3	
Runoff			5
Dai (<i>Dai and Trenberth, 2002</i>)	3	5	
SensibleHeat			2
Fluxnet (<i>Lasslop et al., 2010</i>)	3	3	
GBAF (<i>Jung et al., 2010</i>)	3	5	
TerrestrialWaterStorageAnomaly			5
GRACE (<i>Swenson and Wahr, 2006</i>)	5	5	

Table 4: References and weighting of datasets used to measure the radiation and energy cycle. Weights are chosen using the rubric in Table 1 and reflect a focus on understanding the carbon cycle.

Variable/Dataset	Certainty	Scale	Process
Albedo			1
CERES (<i>Kato et al., 2013</i>)	4	5	
GEWEX.SRB (<i>Stackhouse Jr. et al., 2011</i>)	4	5	
MODIS (<i>De Kauwe et al., 2011</i>)	4	5	
SurfaceUpwardSWRadiation			1
CERES (<i>Kato et al., 2013</i>)	4	4	
GEWEX.SRB (<i>Stackhouse Jr. et al., 2011</i>)	4	5	
WRMC.BSRN (<i>König-Langlo et al., 2013</i>)	4	3	
SurfaceNetSWRadiation			1
CERES (<i>Kato et al., 2013</i>)	4	5	
GEWEX.SRB (<i>Stackhouse Jr. et al., 2011</i>)	4	5	
WRMC.BSRN (<i>König-Langlo et al., 2013</i>)	4	3	
SurfaceUpwardLWRadiation			1
CERES (<i>Kato et al., 2013</i>)	4	5	
GEWEX.SRB (<i>Stackhouse Jr. et al., 2011</i>)	4	5	
WRMC.BSRN (<i>König-Langlo et al., 2013</i>)	4	3	
SurfaceNetLWRadiation			1
CERES (<i>Kato et al., 2013</i>)	4	5	
GEWEX.SRB (<i>Stackhouse Jr. et al., 2011</i>)	4	5	
WRMC.BSRN (<i>König-Langlo et al., 2013</i>)	4	3	
SurfaceNetRadiation			2
CERES (<i>Kato et al., 2013</i>)	4	5	
Fluxnet (<i>Lasslop et al., 2010</i>)	4	3	
GEWEX.SRB (<i>Stackhouse Jr. et al., 2011</i>)	4	5	
WRMC.BSRN (<i>König-Langlo et al., 2013</i>)	4	3	

Table 5: References and weighting of datasets used to measure the forcings. Weights are chosen using the rubric in Table 1 and reflect a focus on understanding the carbon cycle.

Variable/Dataset	Certainty	Scale	Process
SurfaceAirTemperature			2
CRU (<i>Harris et al.</i> , 2014)	5	5	
Fluxnet (<i>Lasslop et al.</i> , 2010)	3	3	
Precipitation			2
CMAP (<i>Xie and Arkin</i> , 1997)	4	5	
Fluxnet (<i>Lasslop et al.</i> , 2010)	3	3	
GPCC (<i>Schneider et al.</i> , 2014)	4	5	
GPCP2 (<i>Adler et al.</i> , 2012)	4	5	
SurfaceRelativeHumidity			3
ERA (<i>Dee et al.</i> , 2011)	2	5	
SurfaceDownwardSWRadiation			2
CERES (<i>Kato et al.</i> , 2013)	4	5	
Fluxnet (<i>Lasslop et al.</i> , 2010)	4	3	
GEWEX.SRB (<i>Stackhouse Jr. et al.</i> , 2011)	4	5	
WRMC.BSRN (<i>König-Langlo et al.</i> , 2013)	4	3	
SurfaceDownwardLWRadiation			1
CERES (<i>Kato et al.</i> , 2013)	4	5	
GEWEX.SRB (<i>Stackhouse Jr. et al.</i> , 2011)	4	5	
WRMC.BSRN (<i>König-Langlo et al.</i> , 2013)	4	3	

468 tive fraction over time and the bias, we use a ratio of means in place of the mean of the
469 ratio,

$$470 \quad \overline{ef}(\mathbf{x}) = \frac{\overline{L_e}(\mathbf{x})}{\overline{L_e}(\mathbf{x}) + \overline{S_h}(\mathbf{x})}. \quad (39)$$

471 Beyond this change, the evaporative fraction is evaluated using the methodology defined in
472 Section 2.

473 **3.1.2 Albedo**

474 We compare the albedo derived from observational data products (*Kato et al.*, 2013;
475 *Stackhouse Jr. et al.*, 2011; *König-Langlo et al.*, 2013) to that of models using the follow-
476 ing expression,

$$477 \quad al(t, \mathbf{x}) = \frac{R_u(t, \mathbf{x})}{R_d(t, \mathbf{x})}. \quad (40)$$

478 where R_u and R_d is the upward and downward shortwave radiation, respectively. As with
479 the evaporative fraction in Section 3.1.1, the albedo expression can become numerically
480 unstable when R_d approaches 0. Thus we again apply a mask, ignoring regions where no
481 significant incoming radiation is observed, $R_d < \delta$. Equation (40) is used when comparing
482 the RMSE and seasonal cycle. When the period mean and bias are computed, we compute
483 the period mean average albedo based on the ratio of averages,

$$484 \quad \overline{al}(\mathbf{x}) = \frac{\overline{R_u}(\mathbf{x})}{\overline{R_d}(\mathbf{x})}. \quad (41)$$

485 **3.1.3 Global Net Ecosystem Carbon Balance**

486 The observational datasets for the global net ecosystem carbon balance (*Le Quéré*
487 *et al.*, 2016; *Hoffman et al.*, 2014) represent global totals, yet models return this value as
488 fluxes defined over space. To create a model quantity commensurate with the observa-
489 tional data, ILAMB must integrate over the globe using Equation (4). As the observa-
490 tional dataset is a time series, much of our scoring methodology does not apply. For this
491 discussion we will represent the global rate of carbon as nbp [$PgC yr^{-1}$]. We compute the
492 accumulation of nbp

$$anbp(t) = \int_{t_0}^t nbp(t) dt \quad (42)$$

493 and score the difference in accumulated total at the end of the time period. The precise
494 method differs slightly in each observational dataset.

495 The Global Carbon Project (GCP) dataset is derived by taking the land sink (un-
496 certainty of $\pm 0.8 [PgC yr^{-1}]$) and subtracting the emissions from land-use change (uncer-
497 tainty of $\pm 0.5 [PgC yr^{-1}]$). This means that the total uncertainty of the accumulated nbp
498 at the end of 2010 is $\sqrt{0.5^2 + 0.8^2} \cdot (2010 - 1959) = 48.1 [PgC]$. We use this uncertainty
499 to normalize the difference in accumulation at the end of the time period as a measure of
500 relative error,

$$\epsilon_{GCP} = \left| \frac{anbp_{mod}(2010) - anbp_{ref}(2010)}{48.1} \right| \quad (43)$$

501 and then again Equation (9) to compute a score of the difference

$$S_{GCP}^{diff} = e^{-\alpha_{nbp} \epsilon_{GCP}}, \quad (44)$$

502 where $\alpha_{nbp} = 0.287$ and is chosen such that if a model falls within the certainty bounds
503 of the accumulated amount through 2010, the corresponding score is at minimum 0.75.
504 We see this as an important first step in incorporating uncertainty into the comparison
505 methodology. We use the uncertainty to tune the scoring methodology, giving a good
506 score to models that fall inside this uncertainty bound. We also compare the global rates
507 of carbon across the time period in the form of a Taylor score of the time series, S_{GCP}^{dist}
508 Equation (30) where the correlation and standard deviation are taken across the temporal
509 dimension. Then the overall score is

$$S_{GCP}^{nbp} = \frac{1}{2} \left(S_{GCP}^{diff} + S_{GCP}^{dist} \right) \quad (45)$$

511 In the *Hoffman et al. (2014)* dataset, we only score the accumulated amount at the
512 end of the observed period. We omit providing a Taylor scoring of the rates because there
513 appears to be some smoothing of the rate data inherent in the process of producing this
514 dataset. However, this dataset explicitly provides a lower and upper bound on uncertainty
515 as a function of time throughout the dataset. So we determine the integrated uncertainty
516 at the end of 2010 by accumulating the upper (52.4 [PgC]) and lower (-32.1 [PgC]) limit

517 of uncertainty, computing the difference, and then halving the value resulting in an uncer-
 518 tainty of 42.3 [PgC]. We then use the same approach to score the difference,

$$\varepsilon_{\text{Hoffman}} = \left| \frac{a_{\text{mod}}(2010) - a_{\text{ref}}(2010)}{42.3} \right| \quad (46)$$

$$S_{\text{Hoffman}}^{\text{nbp}} = e^{-\alpha_{\text{nbp}} \varepsilon_{\text{Hoffman}}} \quad (47)$$

519 **3.1.4 Runoff**

520 We use the *Dai and Trenberth* (2002) river discharge dataset to assess model perfor-
 521 mance of runoff for the world's 50 largest river basins. First, we compute the mean annual
 522 runoff from the model over the time period of the observational dataset. Then we take the
 523 river discharge data and distribute it over the area of the river basins and compare this to
 524 the mean runoff over the same basin. This simple approach was taken to allow us to com-
 525 pare runoff across models even if they do not have a river routing model.

526 We include plots of the mean runoff of the reference and model over river basins
 527 and the bias, represented in Figure 10. We also include regional mean runoff plots for
 528 each of the river basins included, but only show that of the Amazon river basin in Fig-
 529 ure 10(d). The model performance is then scored using the bias (Section 2.2.1), the inter-
 530 annual variability (Section 2.2.4), and the spatial distribution (Section 2.2.5) metrics.

531 **3.1.5 Terrestrial Water Storage Anomaly**

532 We use the Gravity Recovery and Climate Experiment (GRACE) (*Swenson and*
 533 *Wahr*, 2006) dataset to assess the terrestrial water storage anomaly (twsa) in models. How-
 534 ever, there are a few challenges in producing a fair comparison. The first of those is that
 535 models report only the storage and so the anomaly must be computed. The more seri-
 536 ous challenge is that the resolution of this data is quite coarse (300–400 [km]) and thus,
 537 pointwise comparisons are not appropriate (*Swenson*, 2013). Instead we compare mean
 538 anomaly values over 30 of the world's largest river basins. In this way the comparison is
 539 more fair as it is over large areas and automatically omits dry areas which are not of inter-
 540 est.

541 We include plots of the magnitude of the mean anomaly of the reference and model
 542 over river basins and the RMSE, represented in Figure 11. We also include regional mean
 543 anomaly plots for each of the river basins, but show only that of the Amazon river basin

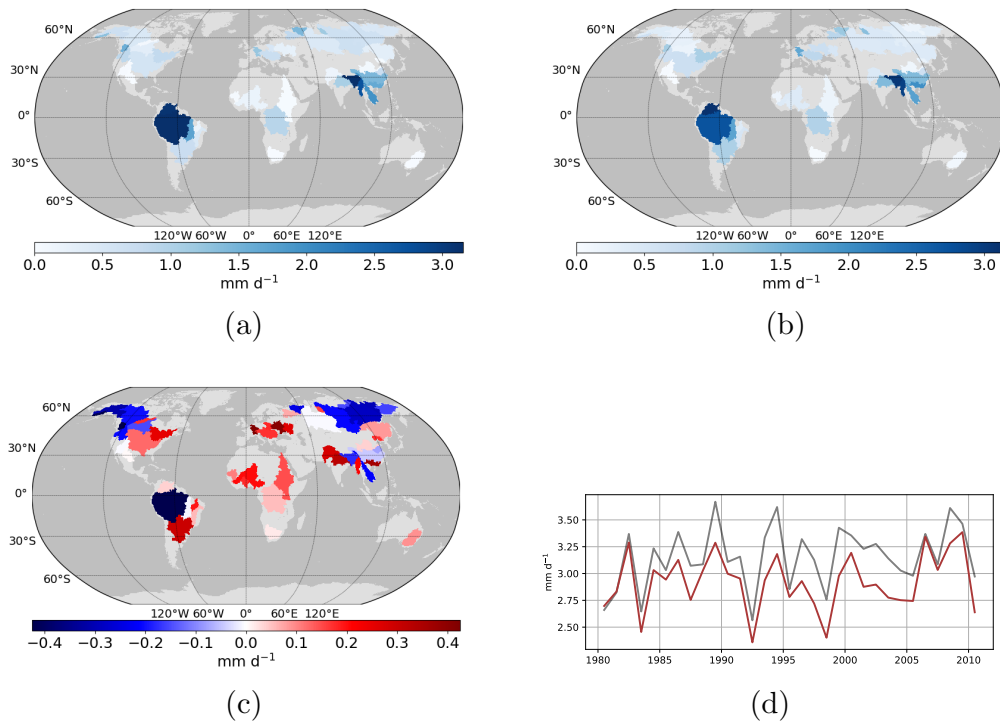


Figure 10: Comparisons of runoff between the reference (*Dai and Trenberth, 2002*) and the model (CLM4.5) dataset. (a) Reference mean runoff, (b) Model mean runoff, (c) Mean runoff bias, (d) Annual mean runoff for the Amazon river basin where the reference is shown in grey and the model in maroon.

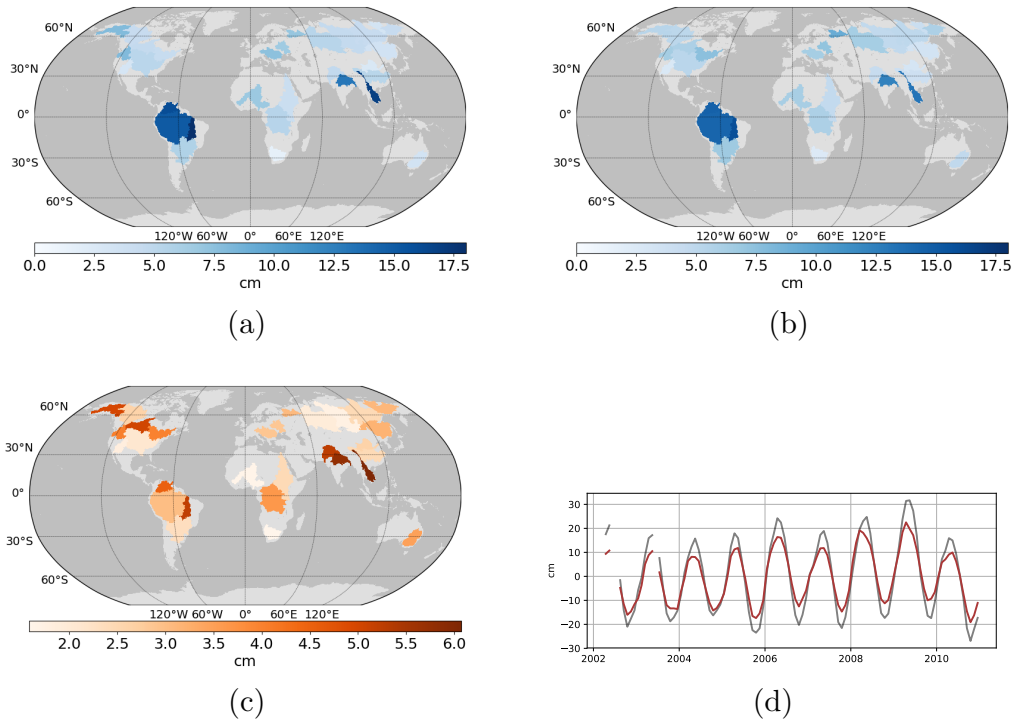


Figure 11: Comparisons of the terrestrial water storage anomaly between the reference (GRACE) and the model (CLM4.5) dataset. (a) Reference mean anomaly magnitude, (b) Model mean anomaly magnitude, (c) Mean anomaly RMSE, (d) Annual mean anomaly for the Amazon river basin where the reference is shown in grey and the model in maroon.

544 in Figure 11(d). The model performance is then scored using the RMSE (Section 2.2.2)
 545 and the interannual variability (Section 2.2.4) metrics.

546 **4 Software**

547 We have implemented the methodology described in Sections 2 and 3 into a soft-
 548 ware package that is freely available to the community. We previously developed a proto-
 549 type implementation (*Mu et al., 2015*) based on the NCAR Command Language (NCL).
 550 We then moved the algorithm into an open-source, openly-developed python package (*Col-
 551 lier et al., 2016*) in an effort to produce a product to which the community can more easily
 552 make contributions. The referenced digital object identifier (DOI) will lead to the software
 553 repository, where the source code and documentation can be found. The documentation
 554 includes the public interface as well as tutorials that span topics such as installation, basic
 555 usage, adding models or benchmark datasets, and formatting benchmark datasets.

556 The ILAMB package is designed to ingest datasets which follow the Climate and
557 Forecast (CF) convention (*Eaton et al., 2017*). The CF website explains that the “conven-
558 tions define metadata that provide a definitive description of what the data in each variable
559 represents, and the spatial and temporal properties of the data. This enables users of data
560 from different sources to decide which quantities are comparable, and facilitates building
561 applications with powerful extraction, regridding, and display capabilities.” We have built
562 the ILAMB package to embody this philosophy, making it directly useful to those who ad-
563 here to this standard. While model intercomparison efforts, such as CMIP5, have encour-
564 aged the use of these conventions among modelers, the observational community has not
565 yet widely put them into practice. Much of the work in adding datasets to the collection is
566 in encoding them to follow this convention.

567 For the purpose of communicating how the ILAMB package works, consider the
568 configure file shown in Figure 12, which defines a set of observational datasets that will
569 be used to confront models. The ‘h1’ bracket is a heading used to categorize variables,
570 represented by the ‘h2’ heading. This comparison involves the surface upward shortwave
571 radiation and the albedo, both of which are variables belonging to the radiation and en-
572 ergy cycle. Inside each ‘h2’ heading, we specify the variable name that will be compared
573 (‘rsus’ is the netCDF variable name for surface upward shortwave radiation). However,
574 we provide a mechanism for variable synonyms in this case by specifying alternate vari-
575 able names. If the ILAMB system cannot find the main variable, it will try to find any
576 alternates that the user specifies. This allows the software to encourage the use of stan-
577 dard variable names, but accounts for modeling groups wanting to use ILAMB without
578 pre-processing. Also note the ‘derived’ keyword in the albedo section. While the compo-
579 nents of albedo are part of standard model output, the albedo is not. The ILAMB package
580 allows for users to specify algebraic relationships in the configure file process. This makes
581 the process automatic and transparent to those who may read this configure file.

582 The ILAMB package will ingest this configure file and try to build commensu-
583 rate quantities from model outputs. While observational datasets come in different forms
584 (globally gridded remote sensing products, tower data collections, etc.), the ILAMB sys-
585 tem reads the spatial and temporal information found in the file and uses it to trim, sub-
586 sample, and/or coarsen the model data as appropriate.

```
[h1: Radiation and Energy Cycle]

[h2: Surface Upward SW Radiation]
variable = "rsus"
alternate_vars = "FSNS"

[CERES]
source = "DATA/rsus/CERES/rsus_0.5x0.5.nc"

[h2: Albedo]
variable = "albedo"
derived = "rsus/rsds"

[CERES]
source = "DATA/albedo/CERES/albedo_0.5x0.5.nc"
```

Figure 12: Sample ILAMB configure file defining comparisons to the surface upward shortwave radiation and albedo variables from the CERES (*Kato et al., 2013*) product.

5 Discussion

The ILAMB framework is designed to be both powerful and flexible. While we have made choices in the default configuration, described above, focused on global analysis for decadal to centennial scale ESMs, ILAMB allows the user to customize selection of variables, weighting of datasets, and spatial subsetting that make it useful for assessing results from mesoscale weather forecasting or other models. We envision developing a library of sample configuration files, targeting various well-known models and model applications.

As much of the usefulness of ILAMB depends on the quality of the underlying observational data, we recommend that data providers include explicit representations of the underlying spatial grids including the areas over which quantities have been averaged. Observational datasets frequently report mean values in a cell taken over an area which may include land but also portions of lakes, rivers, and oceans. This leads to ambiguity with regard to the contribution of land cover types to the measurement itself and subsequently adds to the uncertainty when comparing values to model output.

5.1 Interpreting the Overall Score

The thrust of this paper is to detail a methodology for computing a single overall score that captures a model's skill in reproducing patterns found in the observed record. However, we do not view the absolute value of the score as particularly meaningful beyond the precise definition described in this paper. In general, no model can achieve a perfect score for any given variable for several reasons.

First, there is measurement error and uncertainty in the observational data that makes a perfect score unlikely or undesirable against even a single dataset. This is what motivates some in the community (*Abramowitz, 2005; Best et al., 2015*) to pose that benchmarking requires an expectation of performance which is admittedly lacking in our approach. Second, despite that every attempt is made to employ multiple independent datasets of high quality for confrontation with models, these datasets are inconsistent with each other, making a perfect score across all datasets impossible. We do this as comparisons with multiple observational and synthesized datasets for a single variable offer the user more information about the robustness of model predictions within the limits of observational uncertainty at varying spatial and temporal scales. Third, a lower score with respect to a given variable is not necessarily a sign of a poor model. It may rather highlight the

618 need for better measurement campaigns or improved metrics (i.e., sometimes we learn that
619 our measurements are incomplete or do not acknowledge important uncertainties, or our
620 metrics are inappropriate for a given dataset).

621 The overall score is meant to aid the scientist in discovering when meaningful changes
622 have occurred in the model or across models. The holistic nature of the ILAMB suite of
623 datasets and metrics helps provide a synthesis of model performance that directs the atten-
624 tion of the user to relevant aspects. While we present Figure 1 as the main result of the
625 ILAMB methodology, it is intended to merely indicate variables of particular interest for
626 further consideration. ILAMB output is presented as a hierarchy of interactive webpages
627 that employ JavaScript features to present information to users in a logical and intuitive
628 fashion. From the graphical overview, the user can select individual variables and datasets
629 from the “Results Table” tab to be led to pages which detail the contributing factors to the
630 model’s overall score. On this new page, pre-defined spatial regions can be individually
631 selected, causing the tabular data and diagnostics to be updated automatically to reflect
632 information relevant only to that region. Although all the tabular information, scores, and
633 graphical diagnostics are pre-computed and generated when ILAMB is run, the web-based
634 interface is designed to facilitate discovery and understanding of model results. The over-
635 all score does not replace the scientist, it guides her/him to the relevant plots and diagnos-
636 tics.

637 **5.2 How is ILAMB Used?**

638 The ILAMB package is particularly useful for verification, i.e., during model devel-
639 opment to confirm that new model code improves performance in a targeted area without
640 degrading performance in another area, and for validation, i.e., when comparing perfor-
641 mance of one model or model version to that of other models or model versions.

642 In developing and applying the ILAMB package, we have incorporated a wide va-
643 riety of representative observational datasets (see Tables 2, 3, 4, and 5), and we have fa-
644 vored data that have the most open data policies. In many cases, these data have been av-
645 eraged or remapped to be more directly comparable with model output. As this collection
646 of datasets grows, maintaining and distributing the latest versions will be challenging and
647 require community collaboration. For tracking the evolving performance of models over
648 the long term, it may be necessary to maintain access to older versions of data as well

649 as the latest version since corrections to observational datasets can significantly impact
650 model performance scores. Various technologies could fill this role, and the Observations
651 for Climate Model Intercomparisons (obs4MIPs; [https://www.earthsystemcog.org/
652 projects/obs4mips/](https://www.earthsystemcog.org/projects/obs4mips/)) activity shows promise as a potential solution to this challenge.
653 The preferred solution would ideally support versioning and allow for long-lived versions
654 associated with ILAMB releases. In the interim, we have implemented a simple scheme
655 for sharing summarized and remapped datasets through a webserver.

656 The ILAMB package is currently being used by individual model developers and
657 international modeling centers. ILAMB offers developers a quick and easy method for
658 checking the impacts of new model development before committing code changes. For
659 modeling centers, ILAMB provides a systematic assessment of historical simulation ex-
660 periments and enables tracking of performance of model revisions. ILAMB will also be
661 useful for model intercomparison projects (MIPs) as a starting point for evaluating model
662 variability and uncertainty. As a part of such MIPs, investigators may wish to develop
663 custom metrics or incorporate datasets specific to their purposes. ILAMB could be exe-
664 cuted automatically as model results are uploaded to a system like the Earth System Grid
665 Federation (ESGF; <https://esgf.llnl.gov/>) to give users a “first look” at variation
666 in results and to determine if output should be downloaded for a particular study. ILAMB
667 diagnostics can also be useful for parameter sensitivity studies or for optimization experi-
668 ments in combination with an automated modeling framework like the Predictive Ecosys-
669 tem Analyzer (PEcAn; <http://pecanproject.org/>; *LeBauer et al.*, 2013; *Dietze et al.*,
670 2014). For the assessments community, the results of a multi-model ILAMB evaluation
671 could be useful for understanding which model results would be appropriate for use in
672 studying impacts and which models may poorly capture processes relevant to the impacts
673 under consideration.

674 **5.3 Future Work**

675 Development of the ILAMB package is ongoing, and the terrestrial modeling and
676 observational communities are being engaged to identify *in situ* and remote sensing datasets,
677 to define additional evaluation metrics, and to use the package for a wide variety of MIPs
678 (*Hoffman et al.*, 2017). While most effort has been invested in global- and regional-scale
679 model evaluation, new work is focused on improved benchmarking for site-level time se-
680 ries, spatial transects, and seasonal and diurnal variability. Future development will in-

681 clude incorporation of experiment-specific model evaluation metrics derived from prior
682 studies, including Free-Air CO₂ Enrichment (FACE) (Zaehle *et al.*, 2014; Walker *et al.*,
683 2014, 2015), nutrient addition, rainfall exclusion, and warming experiments (Bouskill *et al.*,
684 2014; Zhu *et al.*, 2016). Partner activities, like NASA's Permafrost Benchmarking System
685 project and the Arctic-Boreal Vulnerability Experiment (ABOVE), are integrating addi-
686 tional datasets and building metrics for specific regions, study areas, or processes of inter-
687 est. We are applying the ILAMB methodology and code base to develop a marine biogeo-
688 chemical model benchmarking tool, called the International Ocean Model Benchmarking
689 (IOMB) package.

690 Based on previous prototypes and community discussion, we developed the ILAMB
691 model benchmarking package for evaluating the fidelity of land carbon cycle models.
692 The package generates graphical diagnostics and computes a comprehensive set of statis-
693 tics through model–data comparisons, and scores model performance for a wide variety
694 of variables for a suite of observational datasets. Rigorously defined model evaluation
695 metrics and strategies for handling multiple resolutions and land masks are documented
696 above. The ILAMB package is open source and is becoming widely adopted by modeling
697 centers and for informing model intercomparison studies. We are actively seeking commu-
698 nity involvement in adding more evaluation metrics and new observational datasets.

699 **Acknowledgments**

700 This manuscript has been authored by UT-Battelle, LLC under Contract No. DE-
701 AC05-00OR22725 with the U.S. Department of Energy. The United States Government
702 retains and the publisher, by accepting the article for publication, acknowledges that the
703 United States Government retains a non-exclusive, paid-up, irrevocable, world-wide license
704 to publish or reproduce the published form of this manuscript, or allow others to do so,
705 for United States Government purposes. The Department of Energy will provide public
706 access to these results of federally sponsored research in accordance with the DOE Public
707 Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

708 This research was supported through the Reducing Uncertainties in Biogeochemical
709 Interactions through Synthesis and Computation Scientific Focus Area (RUBISCO SFA),
710 which is sponsored by the Regional and Global Climate Modeling (RGCM) Program in
711 the Climate and Environmental Sciences Division (CESD) of the Office of Biological and

712 Environmental Research (BER) in the U.S. Department of Energy Office of Science. Oak
713 Ridge National Laboratory (ORNL) is managed by UT-Battelle, LLC for the U.S. Depart-
714 ment of Energy under Contract No. DE-AC05-00OR22725. The National Center for At-
715 mospheric Research (NCAR) is managed by the University Corporation for Atmospheric
716 Research (UCAR) on behalf of the National Science Foundation (NSF). Lawrence Berke-
717 ley National Laboratory (LBNL) is managed and operated by the Regents of the Univer-
718 sity of California under Contract No. DE-AC02-05CH11231.

719 References

- 720 Abramowitz, G. (2005), Towards a benchmark for land surface models, *Geophys. Res.*
721 *Let.*, *32*(22), L22,702, doi:10.1029/2005GL024419.
- 722 Abramowitz, G. (2012), Towards a public, standardized, diagnostic benchmarking system
723 for land surface models, *Geosci. Model Dev.*, *5*(3), 819–827, doi:10.5194/gmd-5-819-
724 2012.
- 725 Adler, R. F., G. Gu, and G. J. Huffman (2012), Estimating climatological bias errors for
726 the Global Precipitation Climatology Project (GPCP), *J. Appl. Meteor. Climatol.*, *51*(1),
727 84–99, doi:10.1175/JAMC-D-11-052.1.
- 728 Anav, A., P. Friedlingstein, M. Kidston, L. Bopp, P. Ciais, P. Cox, C. Jones, M. Jung,
729 R. Myneni, and Z. Zhu (2013), Evaluating the land and ocean components of the global
730 carbon cycle in the CMIP5 Earth system models, *J. Clim.*, *26*(18), 6801–6843, doi:
731 10.1175/JCLI-D-12-00417.1.
- 732 Arora, V. K., G. J. Boer, P. Friedlingstein, M. Eby, C. D. Jones, J. R. Christian, G. Bo-
733 nan, L. Bopp, V. Brovkin, P. Cadule, T. Hajima, T. Ilyina, K. Lindsay, J. F. Tjiputra,
734 and T. Wu (2013), Carbon-concentration and carbon-climate feedbacks in CMIP5 Earth
735 system models, *J. Clim.*, *26*(15), 5289–5314, doi:10.1175/JCLI-D-12-00494.1.
- 736 Best, M. J., G. Abramowitz, H. R. Johnson, A. J. Pitman, G. Balsamo, A. Boone,
737 M. Cuntz, B. Decharme, P. A. Dirmeyer, J. Dong, M. Ek, Z. Guo, V. Haverd, B. J. J.
738 van den Hurk, G. S. Nearing, B. Pak, C. Peters-Lidard, J. A. Santanello Jr., L. Stevens,
739 and N. Vuichard (2015), The plumbing of land surface models: Benchmarking model
740 performance, *J. Hydrometeor.*, *16*(3), 1425–1442, doi:10.1175/JHM-D-14-0158.1.
- 741 Blackard, J. A., M. V. Finco, E. H. Helmer, G. R. Holden, M. L. Hoppus, D. M. Ja-
742 cobs, A. J. Lister, G. G. Moisen, M. D. Nelson, R. Riemann, B. Ruefenacht, D. Sala-
743 janu, D. L. Weyermann, K. C. Winterberger, T. J. Brandeis, R. L. Czaplewski, R. E.

- 744 McRoberts, P. L. Patterson, and R. P. Tymcio (2008), Mapping U.S. forest biomass us-
745 ing nationwide forest inventory data and moderate resolution information, *Remote Sens.*
746 *Environ.*, *112*(4), 1658–1677, doi:10.1016/j.rse.2007.08.021, Remote Sensing Data As-
747 similation Special Issue.
- 748 Blyth, E., D. B. Clark, R. Ellis, C. Huntingford, S. Los, M. Pryor, M. Best, and S. Sitch
749 (2011), A comprehensive set of benchmark tests for a land surface model of simulta-
750 neous fluxes of water and carbon at both the global and seasonal scale, *Geosci. Model*
751 *Dev.*, *4*(2), 255–269, doi:10.5194/gmd-4-255-2011.
- 752 Bouskill, N. J., W. J. Riley, and J. Tang (2014), Meta-analysis of high-latitude nitrogen-
753 addition and warming studies implies ecological mechanisms overlooked by land mod-
754 els, *Biogeosci.*, *11*(23), 6969–6983, doi:10.5194/bg-11-6969-2014.
- 755 Cadule, P., P. Friedlingstein, L. Bopp, S. Sitch, C. D. Jones, P. Ciais, S. L. Piao, and
756 P. Peylin (2010), Benchmarking coupled climate-carbon models against long-term
757 atmospheric CO₂ measurements, *Global Biogeochem. Cycles*, *24*(2), GB2016, doi:
758 10.1029/2009GB003556.
- 759 Ciais, P., C. Sabine, G. Bala, L. Bopp, V. Brovkin, J. Canadell, A. Chhabra, R. DeFries,
760 J. Galloway, M. Heimann, C. Jones, C. Le Quéré, R. B. Myneni, S. Piao, and P. Thorn-
761 ton (2013), Carbon and other biogeochemical cycles, in *Climate Change 2013: The*
762 *Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report*
763 *of the Intergovernmental Panel on Climate Change*, edited by T. F. Stocker, D. Qin, G.-
764 K. Plattner, M. Tignor, S. K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex, and P. M.
765 Midgley, pp. 465–570, Cambridge University Press, Cambridge, United Kingdom and
766 New York, NY, USA.
- 767 Collier, N., F. Hoffman, M. Mu, J. T. Randerson, and W. J. Riley (2016), In-
768 ternational Land Model Benchmarking (ILAMB) package v2, online, doi:
769 10.18139/ILAMB.v002.00/1251621.
- 770 Cox, P. M., R. A. Betts, C. D. Jones, S. A. Spall, and I. J. Totterdell (2000), Acceleration
771 of global warming due to carbon-cycle feedbacks in a coupled climate model, *Nature*,
772 *408*(6809), 184–187, doi:10.1038/35041539.
- 773 Dai, A., and K. E. Trenberth (2002), Estimates of freshwater discharge from continents:
774 Latitudinal and seasonal variations, *J. Hydrometeor.*, *3*(6), 660–687, doi:10.1175/1525-
775 7541(2002)003<0660:EOFDFC>2.0.CO;2.

- 776 Dalmonech, D., and S. Zaehle (2013), Towards a more objective evaluation of modelled
777 land-carbon trends using atmospheric CO₂ and satellite-based vegetation activity obser-
778 vations, *Biogeosci.*, *10*(6), 4189–4210, doi:10.5194/bg-10-4189-2013.
- 779 De Kauwe, M. G., M. I. Disney, T. Quaife, P. Lewis, and M. Williams (2011), An assess-
780 ment of the MODIS Collection 5 leaf area index product for a region of mixed conifer-
781 ous forest, *Remote Sens. Environ.*, *115*(2), 767–780.
- 782 Dee, D. P., S. M. Uppala, A. J. Simmons, P. Berrisford, P. Poli, S. Kobayashi, U. Andrae,
783 M. A. Balmaseda, G. Balsamo, P. Bauer, P. Bechtold, A. C. M. Beljaars, L. van de
784 Berg, J. Bidlot, N. Bormann, C. Delsol, R. Dragani, M. Fuentes, A. J. Geer, L. Haim-
785 berger, S. B. Healy, H. Hersbach, E. V. HÅşlm, L. Isaksen, P. KÄëllberg, M. Köh-
786 ler, M. Matricardi, A. P. McNally, B. M. Monge-Sanz, J.-J. Morcrette, B.-K. Park,
787 C. Peubey, P. de Rosnay, C. Tavolato, J.-N. ThÄlpaut, and F. Vitart (2011), The ERA-
788 Interim reanalysis: Configuration and performance of the data assimilation system, *Q. J.*
789 *Roy. Meteor. Soc.*, *137*(656), 553–597, doi:10.1002/qj.828.
- 790 Denman, K. L., G. Brasseur, A. Chidthaisong, P. Ciaais, P. M. Cox, R. E. Dickinson,
791 D. Hauglustaine, C. Heinze, E. Holland, D. Jacob, U. Lohmann, S. Ramachandran,
792 P. L. d. Dias, S. C. Wofsy, and X. Zhang (2007), Couplings between changes in the
793 climate system and biogeochemistry, in *Climate Change 2007: The Physical Science Ba-*
794 *sis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovern-*
795 *mental Panel on Climate Change*, edited by S. Solomon, D. Qin, M. Manning, Z. Chen,
796 M. Marquis, K. B. Averyt, M. Tignor, and H. L. Miller, pp. 499–587, Cambridge Uni-
797 versity Press, Cambridge, United Kingdom and New York, NY, USA.
- 798 Dietze, M. C., S. P. Serbin, C. Davidson, A. R. Desai, X. Feng, R. Kelly, R. Kooper,
799 D. LeBauer, J. Mantooth, K. McHenry, and D. Wang (2014), A quantitative assessment
800 of a terrestrial biosphere model’s data needs across North American biomes, *J. Geo-*
801 *phys. Res. Biogeosci.*, *119*(3), 286–300, doi:10.1002/2013JG002392.
- 802 Dirmeyer, P. A., J. Wei, M. G. Bosilovich, and D. M. Mocko (2014), Comparing evapo-
803 rative sources of terrestrial precipitation and their extremes in merra using relative en-
804 tropy, *J. Hydrometeor.*, *15*(1), 102–116, doi:10.1175/JHM-D-13-053.1.
- 805 Eaton, B., J. Gregory, B. Drach, K. Taylor, S. Hankin, J. Blower, J. Caron, R. Signell,
806 P. Bentley, G. Rappa, H. Höck, A. Pamment, M. Jukes, M. Raspaud, and R. Horne
807 (2017), Netcdf climate and forecast (cf) metadata conventions, <http://cfconventions.org/>.

- 808 Eyring, V., M. Righi, A. Lauer, M. Evaldsson, S. Wenzel, C. Jones, A. Anav, O. An-
809 dews, I. Cionni, E. L. Davin, C. Deser, C. Ehbrecht, P. Friedlingstein, P. Gleckler,
810 K.-D. Gottschaldt, S. Hagemann, M. Juckes, S. Kindermann, J. Krasting, D. Kunert,
811 R. Levine, A. Loew, J. Mäkelä, G. Martin, E. Mason, A. S. Phillips, S. Read, C. Rio,
812 R. Roehrig, D. Senftleben, A. Sterl, L. H. van Ulft, J. Walton, S. Wang, and K. D.
813 Williams (2016), ESMValTool (v1.0) – A community diagnostic and performance met-
814 rics tool for routine evaluation of Earth system models in CMIP, *Geosci. Model Dev.*,
815 *9*(5), 1747–1802, doi:10.5194/gmd-9-1747-2016.
- 816 Friedlingstein, P., L. Bopp, P. Ciais, J.-L. Dufresne, L. Fairhead, H. LeTreut, P. Monfray,
817 and J. Orr (2001), Positive feedback between future climate change and the carbon cy-
818 cle, *Geophys. Res. Lett.*, *28*(8), 1543–1546, doi:10.1029/2000GL012015.
- 819 Friedlingstein, P., P. M. Cox, R. A. Betts, L. Bopp, W. von Bloh, V. Brovkin, S. C. Doney,
820 M. Eby, I. Fung, B. Govindasamy, J. John, C. D. Jones, F. Joos, T. Kato, M. Kawamiya,
821 W. Knorr, K. Lindsay, H. D. Matthews, T. Raddatz, P. Rayner, C. Reick, E. Roeckner,
822 K.-G. Schnitzler, R. Schnur, K. Strassmann, S. Thompson, A. J. Weaver, C. Yoshikawa,
823 and N. Zeng (2006), Climate–carbon cycle feedback analysis: Results from the C⁴MIP
824 model intercomparison, *J. Clim.*, *19*(14), 3373–3353, doi:10.1175/JCLI3800.1.
- 825 Friedlingstein, P., M. Meinshausen, V. K. Arora, C. D. Jones, A. Anav, S. K. Liddicoat,
826 and R. Knutti (2014), Uncertainties in CMIP5 climate projections due to carbon cycle
827 feedbacks, *J. Clim.*, *27*(2), 511–526, doi:10.1175/JCLI-D-12-00579.1.
- 828 Fung, I. Y., S. C. Doney, K. Lindsay, and J. John (2005), Evolution of carbon
829 sinks in a changing climate, *Proc. Nat. Acad. Sci.*, *102*(32), 11,201–11,206, doi:
830 10.1073/pnas.0504949102.
- 831 Ghimire, B., W. J. Riley, C. D. Koven, M. Mu, and J. T. Randerson (2016), Representing
832 leaf and root physiological traits in CLM improves global carbon and nitrogen cycling
833 predictions, *J. Adv. Model. Earth Syst.*, *8*(2), 598–613, doi:10.1002/2015MS000538.
- 834 Giglio, L., J. T. Randerson, G. R. van der Werf, P. S. Kasibhatla, G. J. Collatz, D. C.
835 Morton, and R. S. DeFries (2010), Assessing variability and long-term trends in burned
836 area by merging multiple satellite fire products, *Biogeosci.*, *7*(3), 1171–1186, doi:
837 10.5194/bg-7-1171-2010.
- 838 Gleckler, P. J., K. E. Taylor, and C. Doutriaux (2008), Performance metrics for climate
839 models, *J. Geophys. Res.*, *113*(D6), D06,104, doi:10.1029/2007JD008972.

- 840 Gleckler, P. J., C. Doutriaux, P. J. Durack, K. E. Taylor, Y. Zhang, D. N. Williams, E. Ma-
841 son, and J. Servonnat (2016), A more powerful reality test for climate models, *Eos*
842 *Trans. AGU*, 97, doi:10.1029/2016EO051663.
- 843 Gregory, J. M., C. D. Jones, P. Cadule, and P. Friedlingstein (2009), Quantifying carbon
844 cycle feedbacks, *J. Clim.*, 22(19), 5232–5250, doi:10.1175/2009JCLI2949.1.
- 845 Harris, I., P. Jones, T. Osborn, and D. Lister (2014), Updated high-resolution grids of
846 monthly climatic observations – the cru ts3.10 dataset, *Int. J. Climatol.*, 34(3), 623–
847 642, doi:10.1002/joc.3711.
- 848 Hoffman, F. M., J. T. Randerson, V. K. Arora, Q. Bao, P. Cadule, D. Ji, C. D. Jones,
849 M. Kawamiya, S. Khattiwala, K. Lindsay, A. Obata, E. Shevliakova, K. D. Six, J. F.
850 Tjiputra, E. M. Volodin, and T. Wu (2014), Causes and implications of persistent at-
851 mospheric carbon dioxide biases in Earth System Models, *J. Geophys. Res. Biogeosci.*,
852 119(2), 141–162, doi:10.1002/2013JG002381.
- 853 Hoffman, F. M., C. D. Koven, G. Keppel-Aleks, D. M. Lawrence, W. J. Riley, J. T. Ran-
854 derson, A. Ahlström, G. Abramowitz, D. D. Baldocchi, M. J. Best, B. Bond-Lamberty,
855 M. G. De Kauwe, A. S. Denning, A. R. Desai, V. Eyring, J. B. Fisher, R. A. Fisher,
856 P. J. Gleckler, M. Huang, G. Hugelius, A. K. Jain, N. Y. Kiang, H. Kim, R. D. Koster,
857 S. V. Kumar, H. Li, Y. Luo, J. Mao, N. G. McDowell, U. Mishra, P. R. Moorcroft,
858 G. S. H. Pau, D. M. Ricciuto, K. Schaefer, C. R. Schwalm, S. P. Serbin, E. Shevliakova,
859 A. G. Slater, J. Tang, M. Williams, J. Xia, C. Xu, R. Joseph, and D. Koch (2017), In-
860 ternational Land Model Benchmarking (ILAMB) 2016 workshop report, *Tech. Rep.*
861 *DOE/SC-0186*, U.S. Department of Energy, Office of Science, Germantown, Maryland,
862 USA, doi:10.2172/1330803.
- 863 Hugelius, G., J. G. Bockheim, P. Camill, B. Elberling, G. Grosse, J. W. Harden, K. John-
864 son, T. Jorgenson, C. D. Koven, P. Kuhry, G. Michaelson, U. Mishra, J. Palmtag,
865 C.-L. Ping, J. O'Donnell, L. Schirrmeister, E. A. G. Schuur, Y. Sheng, L. C. Smith,
866 J. Strauss, and Z. Yu (2013), A new data set for estimating organic carbon storage to
867 3 m depth in soils of the northern circumpolar permafrost region, *Earth Syst. Sci. Data*,
868 5(2), 393–402, doi:10.5194/essd-5-393-2013.
- 869 Jung, M., M. Reichstein, P. Ciais, S. I. Seneviratne, J. Sheffield, M. L. Goulden, G. Bo-
870 nan, A. Cescatti, J. Chen, R. de Jeu, A. J. Dolman, W. Eugster, D. Gerten, D. Gianelle,
871 N. Gobron, J. Heinke, J. Kimball, B. E. Law, L. Montagnani, Q. Mu, B. Mueller,
872 K. Oleson, D. Papale, A. D. Richardson, O. Roupsard, S. Running, E. Tomelleri,

- 873 N. Viovy, U. Weber, C. Williams, E. Wood, S. Zaehle, and K. Zhang (2010), Recent
874 decline in the global land evapotranspiration trend due to limited moisture supply, *Nature*,
875 *467*, 951–954, doi:10.1038/nature09396.
- 876 Kato, S., N. G. Loeb, F. G. Rose, D. R. Doelling, D. A. Rutan, T. E. Caldwell, L. Yu,
877 and R. A. Weller (2013), Surface irradiances consistent with CERES-derived top-
878 of-atmosphere shortwave and longwave irradiances, *J. Clim.*, *26*(9), 2719–2740, doi:
879 10.1175/JCLI-D-12-00436.1.
- 880 Kelley, D. I., I. C. Prentice, S. P. Harrison, H. Wang, M. Simard, J. B. Fisher, and K. O.
881 Willis (2013), A comprehensive benchmarking system for evaluating global vegetation
882 models, *Biogeosci.*, *10*(5), 3313–3340, doi:10.5194/bg-10-3313-2013.
- 883 Kelldorfer, J., W. Walker, K. Kirsch, G. Fiske, J. Bishop, L. Lapoint, M. Hoppus, and
884 J. Westfall (2013), NACP aboveground biomass and carbon baseline data, V.2 (NBCD
885 2000), U.S.A., 2000, doi:10.3334/ornl daac/1161.
- 886 König-Langlo, G., R. Sieger, H. Schmithüsen, A. Bücker, F. Richter, and E. Dutton
887 (2013), The Baseline Surface Radiation Network and its World Radiation Monitoring
888 Centre at the Alfred Wegener Institute, *Tech. Rep. WCRP-2*, Alfred Wegener Institute,
889 doi:10013/epic.42596.d001.
- 890 Kumar, J., F. M. Hoffman, W. W. Hargrove, and N. Collier (2016), Understanding the rep-
891 resentativeness of FLUXNET for upscaling carbon flux from eddy covariance measure-
892 ments, *Earth System Science Data Discussions*, *2016*, 1–25, doi:10.5194/essd-2016-36.
- 893 Kumar, S. V., C. D. Peters-Lidard, J. Santanello, K. Harrison, Y. Liu, and M. Shaw
894 (2012), Land surface Verification Toolkit (LVT) — A generalized framework for land
895 surface model evaluation, *Geosci. Model Dev.*, *5*(3), 869–886, doi:10.5194/gmd-5-869-
896 2012.
- 897 Lasslop, G., M. Reichstein, D. Papale, A. D. Richardson, A. Arneth, A. Barr, P. Stoy, and
898 G. Wohlfahrt (2010), Separation of net ecosystem exchange into assimilation and respi-
899 ration using a light response curve approach: critical issues and global evaluation, *Glob.*
900 *Change Biol.*, *16*(1), 187–208, doi:10.1111/j.1365-2486.2009.02041.x.
- 901 Law, K., A. Stuart, and K. Zygalkis (2015), *Data Assimilation: A Mathematical Intro-*
902 *duction*, *Texts in Applied Mathematics*, vol. 63, 1 ed., 242 pp., Springer International
903 Publishing, doi:10.1007/978-3-319-20325-6.
- 904 Le Quéré, C., R. M. Andrew, J. G. Canadell, S. Sitch, J. I. Korsbakken, G. P. Peters,
905 A. C. Manning, T. A. Boden, P. P. Tans, R. A. Houghton, R. F. Keeling, S. Alin,

- 906 O. D. Andrews, P. Anthoni, L. Barbero, L. Bopp, F. Chevallier, L. P. Chini, P. Ciais,
907 K. Currie, C. Delire, S. C. Doney, P. Friedlingstein, T. Gkritzalis, I. Harris, J. Hauck,
908 V. Haverd, M. Hoppema, K. Klein Goldewijk, A. K. Jain, E. Kato, A. Körtzinger,
909 P. Landschützer, N. Lefèvre, A. Lenton, S. Lienert, D. Lombardozi, J. R. Melton,
910 N. Metz, F. Millero, P. M. S. Monteiro, D. R. Munro, J. E. M. S. Nabel, S.-I. Nakaoka,
911 K. O'Brien, A. Olsen, A. M. Omar, T. Ono, D. Pierrot, B. Poulter, C. Rödenbeck,
912 J. Salisbury, U. Schuster, J. Schwinger, R. Séférian, I. Skjelvan, B. D. Stocker, A. J.
913 Sutton, T. Takahashi, H. Tian, B. Tilbrook, I. T. van der Laan-Luijkx, G. R. van der
914 Werf, N. Viovy, A. P. Walker, A. J. Wiltshire, and S. Zaehle (2016), Global carbon bud-
915 get 2016, *Earth Syst. Sci. Data*, 8(2), 605–649, doi:10.5194/essd-8-605-2016.
- 916 LeBauer, D. S., D. Wang, K. T. Richter, C. C. Davidson, and M. C. Dietze (2013), Fa-
917 cilitating feedbacks between field measurements and ecosystem models, *Ecol. Monogr.*,
918 83(2), 133–154, doi:10.1890/12-0137.1.
- 919 Luo, Y. Q., J. T. Randerson, G. Abramowitz, C. Bacour, E. Blyth, N. Carvalhais, P. Ciais,
920 D. Dalmonech, J. B. Fisher, R. Fisher, P. Friedlingstein, K. Hibbard, F. Hoffman,
921 D. Huntzinger, C. D. Jones, C. Koven, D. Lawrence, D. J. Li, M. Mahecha, S. L.
922 Niu, R. Norby, S. L. Piao, X. Qi, P. Peylin, I. C. Prentice, W. Riley, M. Reichstein,
923 C. Schwalm, Y. P. Wang, J. Y. Xia, S. Zaehle, and X. H. Zhou (2012), A framework for
924 benchmarking land models, *Biogeosci.*, 9(10), 3857–3874, doi:10.5194/bg-9-3857-2012.
- 925 Mahowald, N. M., J. T. Randerson, K. Lindsay, E. Muñoz, S. C. Doney, P. Lawrence,
926 S. Schlunegger, D. S. Ward, D. Lawrence, and F. M. Hoffman (2017), Interactions be-
927 tween land use change and carbon cycle feedbacks, *Global Biogeochem. Cycles*, 31(1),
928 96–113, doi:10.1002/2016GB005374.
- 929 Meehl, G. A., C. Covey, T. Delworth, M. Latif, B. McAvaney, J. F. B. Mitchell, R. J.
930 Stouffer, and K. E. Taylor (2007), The WCRP CMIP3 multimodel dataset: A new
931 era in climate change research, *Bull. Am. Meteorol. Soc.*, 88(9), 1383–1394, doi:
932 10.1175/BAMS-88-9-1383.
- 933 Miralles, D., T. Holmes, R. D. Jeu, J. Gash, A. Meesters, and A. Dolman (2011), Global
934 land-surface evaporation estimated from satellite-based observations, *Hydrol. Earth Syst.*
935 *Sci.*, 15, 453–469.
- 936 Moore, J. K., W. Fu, F. Primeau, G. L. Britten, K. Lindsay, M. Long, S. C. Doney,
937 N. Mahowald, F. M. Hoffman, and J. T. Randerson (2018), Sustained climate warm-
938 ing drives declining marine biological productivity, *Science*, 359(6380), 1139–1143,

939 doi:10.1126/science.aao6379.

940 Mu, M., J. T. Randerson, W. J. Riley, C. D. Koven, G. Keppel-Aleks, D. M. Lawrence, ,
941 and F. M. Hoffman (2015), International Land Model Benchmarking (ILAMB) package
942 v1, online, doi:10.18139/ILAMB.v001.00/1251597.

943 Myneni, R. B., R. R. Nemani, and S. Running (1997), Algorithm for the estimation of
944 global land cover, LAI and FPAR based on radiative transfer models, *IEEE Trans.*
945 *Geosci. Remote Sens.*, *35*, 1380–1392.

946 Oleson, K. W., D. M. Lawrence, G. B. Bonan, B. Drewniak, M. Huang, C. D. Koven,
947 S. Levis, F. Li, W. J. Riley, Z. M. Subin, S. C. Swenson, P. E. Thornton, A. Bozbiyik,
948 R. Fisher, C. L. Heald, E. Kluzek, J.-F. Lamarque, P. J. Lawrence, L. R. Leung, W. Lip-
949 scomb, S. Muszala, D. M. Ricciuto, W. Sacks, Y. Sun, J. Tang, and Z.-L. Yang (2013),
950 Technical description of version 4.5 of the Community Land Model (CLM), *Techni-*
951 *cal Note NCAR/TN-503+STR*, National Center for Atmospheric Research, Boulder, Col-
952 orado, USA.

953 Piao, S., S. Sitch, P. Ciais, P. Friedlingstein, P. Peylin, X. Wang, A. Ahlström, A. Anav,
954 J. G. Canadell, N. Cong, C. Huntingford, M. Jung, S. Levis, P. E. Levy, J. Li, X. Lin,
955 M. R. Lomas, M. Lu, Y. Luo, Y. Ma, R. B. Myneni, B. Poulter, Z. Sun, T. Wang,
956 N. Viovy, S. Zaehle, and N. Zeng (2013), Evaluation of terrestrial carbon cycle models
957 for their response to climate variability and to CO₂ trends, *Glob. Change Biol.*, *19*(7),
958 2117–2132, doi:10.1111/gcb.12187.

959 Randerson, J. T., F. M. Hoffman, P. E. Thornton, N. M. Mahowald, K. Lindsay, Y.-H.
960 Lee, C. D. Nevison, S. C. Doney, G. Bonan, R. Stöckli, C. Covey, S. W. Running,
961 and I. Y. Fung (2009), Systematic assessment of terrestrial biogeochemistry in cou-
962 pled climate–carbon models, *Glob. Change Biol.*, *15*(9), 2462–2484, doi:10.1111/j.1365-
963 2486.2009.01912.x.

964 Randerson, J. T., K. Lindsay, E. Munoz, W. Fu, J. K. Moore, F. M. Hoffman, N. M. Ma-
965 howald, and S. C. Doney (2015), Multicentury changes in ocean and land contribu-
966 tions to the climate–carbon feedback, *Global Biogeochem. Cycles*, *29*(6), 744–759, doi:
967 10.1002/2014GB005079.

968 Reichler, T., and J. Kim (2008), How well do coupled models simulate today’s climate?,
969 *Bull. Am. Meteorol. Soc.*, *89*(3), 303–311, doi:10.1175/BAMS-89-3-303.

970 Saatchi, S. S., N. L. Harris, S. Brown, M. Lefsky, E. T. A. Mitchard, W. Salas,
971 B. R. Zutta, W. Buermann, S. L. Lewis, S. Hagen, S. Petrova, L. White, M. Sil-

- 972 man, and A. Morel (2011), Benchmark map of forest carbon stocks in tropical
973 regions across three continents, *Proc. Nat. Acad. Sci.*, *108*(24), 9899–9904, doi:
974 10.1073/pnas.1019576108.
- 975 Schneider, U., A. Becker, P. Finger, A. Meyer-Christoffer, M. Ziese, and B. Rudolf (2014),
976 GPCP’s new land surface precipitation climatology based on quality-controlled in situ
977 data and its role in quantifying the global water cycle, *Theor. Appl. Climatol.*, *115*(1),
978 15–40, doi:10.1007/s00704-013-0860-x.
- 979 Stackhouse Jr., P. W., S. K. Gupta, S. J. Cox, J. C. Mikovitz, T. Zhang, and L. M. Hinkel-
980 man (2011), The NASA/GEWEX surface radiation budget release 3.0: 24.5-year
981 dataset, *GEWEX News*, *21*(1), 10–12.
- 982 Swenson, S. (2013), GRACE: Gravity recovery and climate experiment: Surface mass, to-
983 tal water storage, and derived variables, in *The Climate Data Guide*, edited by National
984 Center for Atmospheric Research Staff, National Center for Atmospheric Research.
- 985 Swenson, S., and J. Wahr (2006), Post-processing removal of correlated errors in GRACE
986 data, *Geophys. Res. Lett.*, *33*(8), L08,402, doi:10.1029/2005GL025285.
- 987 Taylor, K. E. (2001), Summarizing multiple aspects of model performance in a single dia-
988 gram, *J. Geophys. Res. Atmos.*, *106*(D7), 7183–7192, doi:10.1029/2000JD900719.
- 989 Taylor, K. E., R. J. Stouffer, and G. A. Meehl (2012), An overview of CMIP5 and the ex-
990 periment design, *Bull. Am. Meteorol. Soc.*, *93*(4), 485–498, doi:10.1175/BAMS-D-11-
991 00094.1.
- 992 Todd-Brown, K. E. O., J. T. Randerson, W. M. Post, F. M. Hoffman, C. Tarnocai, E. A. G.
993 Schuur, and S. D. Allison (2013), Causes of variation in soil carbon simulations from
994 CMIP5 Earth system models and comparison with observations, *Biogeosci.*, *10*(3),
995 1717–1736, doi:10.5194/bg-10-1717-2013.
- 996 Walker, A. P., P. J. Hanson, M. G. De Kauwe, B. E. Medlyn, S. Zaehle, S. Asao, M. Di-
997 etze, T. Hickler, C. Huntingford, C. M. Iversen, A. Jain, M. Lomas, Y. Luo, H. Mc-
998 Carthy, W. J. Parton, I. C. Prentice, P. E. Thornton, S. Wang, Y.-P. Wang, D. Warlind,
999 E. Weng, J. M. Warren, F. I. Woodward, R. Oren, and R. J. Norby (2014), Comprehen-
1000 sive ecosystem model–data synthesis using multiple data sets at two temperate forest
1001 free-air CO₂ enrichment experiments: Model performance at ambient CO₂ concentra-
1002 tion, *J. Geophys. Res. Biogeosci.*, *119*(5), 2169–8961, doi:10.1002/2013JG002553.
- 1003 Walker, A. P., S. Zaehle, B. E. Medlyn, M. G. De Kauwe, S. Asao, T. Hickler, W. Par-
1004 ton, D. M. Ricciuto, Y.-P. Wang, D. Wärlind, and R. J. Norby (2015), Predicting

- 1005 long-term carbon sequestration in response to CO₂ enrichment: How and why do
1006 current ecosystem models differ?, *Global Biogeochem. Cycles*, 29(4), 476–495, doi:
1007 10.1002/2014GB004995.
- 1008 Xie, P., and P. A. Arkin (1997), Global precipitation: A 17-year monthly anal-
1009 ysis based on gauge observations, satellite estimates, and numerical model
1010 outputs, *Bull. Am. Meteorol. Soc.*, 78(11), 2539–2558, doi:10.1175/1520-
1011 0477(1997)078<2539:GPAYMA>2.0.CO;2.
- 1012 Zaehle, S., B. E. Medlyn, M. G. De Kauwe, A. P. Walker, M. C. Dietze, T. Hickler,
1013 Y. Luo, Y.-P. Wang, B. El-Masri, P. Thornton, A. Jain, S. Wang, D. Warlind, E. Weng,
1014 W. Parton, C. M. Iversen, A. Gallet-Budynek, H. McCarthy, A. Finzi, P. J. Hanson,
1015 I. C. Prentice, R. Oren, and R. J. Norby (2014), Evaluation of 11 terrestrial carbon–
1016 nitrogen cycle models against observations from two temperate free-air CO₂ enrichment
1017 studies, *New Phytol.*, 202(3), 803–822, doi:10.1111/nph.12697.
- 1018 Zhu, Q., W. J. Riley, J. Tang, and C. D. Koven (2016), Multiple soil nutrient competition
1019 between plants, microbes, and mineral surfaces: Model development, parameterization,
1020 and example applications in several tropical forests, *Biogeosci.*, 13(1), 341–363, doi:
1021 10.5194/bg-13-341-2016.

1022 References

- 1023 Abramowitz, G. (2005), Towards a benchmark for land surface models, *Geophys. Res.*
1024 *Let.*, 32(22), L22,702, doi:10.1029/2005GL024419.
- 1025 Abramowitz, G. (2012), Towards a public, standardized, diagnostic benchmarking system
1026 for land surface models, *Geosci. Model Dev.*, 5(3), 819–827, doi:10.5194/gmd-5-819-
1027 2012.
- 1028 Adler, R. F., G. Gu, and G. J. Huffman (2012), Estimating climatological bias errors for
1029 the Global Precipitation Climatology Project (GPCP), *J. Appl. Meteor. Climatol.*, 51(1),
1030 84–99, doi:10.1175/JAMC-D-11-052.1.
- 1031 Anav, A., P. Friedlingstein, M. Kidston, L. Bopp, P. Ciais, P. Cox, C. Jones, M. Jung,
1032 R. Myneni, and Z. Zhu (2013), Evaluating the land and ocean components of the global
1033 carbon cycle in the CMIP5 Earth system models, *J. Clim.*, 26(18), 6801–6843, doi:
1034 10.1175/JCLI-D-12-00417.1.
- 1035 Arora, V. K., G. J. Boer, P. Friedlingstein, M. Eby, C. D. Jones, J. R. Christian, G. Bo-
1036 nan, L. Bopp, V. Brovkin, P. Cadule, T. Hajima, T. Ilyina, K. Lindsay, J. F. Tjiputra,

- 1037 and T. Wu (2013), Carbon-concentration and carbon-climate feedbacks in CMIP5 Earth
1038 system models, *J. Clim.*, *26*(15), 5289–5314, doi:10.1175/JCLI-D-12-00494.1.
- 1039 Best, M. J., G. Abramowitz, H. R. Johnson, A. J. Pitman, G. Balsamo, A. Boone,
1040 M. Cuntz, B. Decharme, P. A. Dirmeyer, J. Dong, M. Ek, Z. Guo, V. Haverd, B. J. J.
1041 van den Hurk, G. S. Nearing, B. Pak, C. Peters-Lidard, J. A. Santanello Jr., L. Stevens,
1042 and N. Vuichard (2015), The plumbing of land surface models: Benchmarking model
1043 performance, *J. Hydrometeor.*, *16*(3), 1425–1442, doi:10.1175/JHM-D-14-0158.1.
- 1044 Blackard, J. A., M. V. Finco, E. H. Helmer, G. R. Holden, M. L. Hoppus, D. M. Ja-
1045 cobs, A. J. Lister, G. G. Moisen, M. D. Nelson, R. Riemann, B. Ruefenacht, D. Sala-
1046 janu, D. L. Weyermann, K. C. Winterberger, T. J. Brandeis, R. L. Czaplewski, R. E.
1047 McRoberts, P. L. Patterson, and R. P. Tymcio (2008), Mapping U.S. forest biomass us-
1048 ing nationwide forest inventory data and moderate resolution information, *Remote Sens.*
1049 *Environ.*, *112*(4), 1658–1677, doi:10.1016/j.rse.2007.08.021, Remote Sensing Data As-
1050 similation Special Issue.
- 1051 Blyth, E., D. B. Clark, R. Ellis, C. Huntingford, S. Los, M. Pryor, M. Best, and S. Sitch
1052 (2011), A comprehensive set of benchmark tests for a land surface model of simulta-
1053 neous fluxes of water and carbon at both the global and seasonal scale, *Geosci. Model*
1054 *Dev.*, *4*(2), 255–269, doi:10.5194/gmd-4-255-2011.
- 1055 Bouskill, N. J., W. J. Riley, and J. Tang (2014), Meta-analysis of high-latitude nitrogen-
1056 addition and warming studies implies ecological mechanisms overlooked by land mod-
1057 els, *Biogeosci.*, *11*(23), 6969–6983, doi:10.5194/bg-11-6969-2014.
- 1058 Cadule, P., P. Friedlingstein, L. Bopp, S. Sitch, C. D. Jones, P. Ciais, S. L. Piao, and
1059 P. Peylin (2010), Benchmarking coupled climate-carbon models against long-term
1060 atmospheric CO₂ measurements, *Global Biogeochem. Cycles*, *24*(2), GB2016, doi:
1061 10.1029/2009GB003556.
- 1062 Ciais, P., C. Sabine, G. Bala, L. Bopp, V. Brovkin, J. Canadell, A. Chhabra, R. DeFries,
1063 J. Galloway, M. Heimann, C. Jones, C. Le Quéré, R. B. Myneni, S. Piao, and P. Thorn-
1064 ton (2013), Carbon and other biogeochemical cycles, in *Climate Change 2013: The*
1065 *Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report*
1066 *of the Intergovernmental Panel on Climate Change*, edited by T. F. Stocker, D. Qin, G.-
1067 K. Plattner, M. Tignor, S. K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex, and P. M.
1068 Midgley, pp. 465–570, Cambridge University Press, Cambridge, United Kingdom and
1069 New York, NY, USA.

- 1070 Collier, N., F. Hoffman, M. Mu, J. T. Randerson, and W. J. Riley (2016), In-
1071 ternational Land Model Benchmarking (ILAMB) package v2, online, doi:
1072 10.18139/ILAMB.v002.00/1251621.
- 1073 Cox, P. M., R. A. Betts, C. D. Jones, S. A. Spall, and I. J. Totterdell (2000), Acceleration
1074 of global warming due to carbon–cycle feedbacks in a coupled climate model, *Nature*,
1075 408(6809), 184–187, doi:10.1038/35041539.
- 1076 Dai, A., and K. E. Trenberth (2002), Estimates of freshwater discharge from continents:
1077 Latitudinal and seasonal variations, *J. Hydrometeor.*, 3(6), 660–687, doi:10.1175/1525-
1078 7541(2002)003<0660:EOFDFC>2.0.CO;2.
- 1079 Dalmonech, D., and S. Zaehle (2013), Towards a more objective evaluation of modelled
1080 land-carbon trends using atmospheric CO₂ and satellite-based vegetation activity obser-
1081 vations, *Biogeosci.*, 10(6), 4189–4210, doi:10.5194/bg-10-4189-2013.
- 1082 De Kauwe, M. G., M. I. Disney, T. Quaife, P. Lewis, and M. Williams (2011), An assess-
1083 ment of the MODIS Collection 5 leaf area index product for a region of mixed conifer-
1084 ous forest, *Remote Sens. Environ.*, 115(2), 767–780.
- 1085 Dee, D. P., S. M. Uppala, A. J. Simmons, P. Berrisford, P. Poli, S. Kobayashi, U. Andrae,
1086 M. A. Balmaseda, G. Balsamo, P. Bauer, P. Bechtold, A. C. M. Beljaars, L. van de
1087 Berg, J. Bidlot, N. Bormann, C. Delsol, R. Dragani, M. Fuentes, A. J. Geer, L. Haim-
1088 berger, S. B. Healy, H. Hersbach, E. V. HÅşlm, L. Isaksen, P. KÄällberg, M. Köh-
1089 ler, M. Matricardi, A. P. McNally, B. M. Monge-Sanz, J.-J. Morcrette, B.-K. Park,
1090 C. Peubey, P. de Rosnay, C. Tavolato, J.-N. ThÄlpaut, and F. Vitart (2011), The ERA-
1091 Interim reanalysis: Configuration and performance of the data assimilation system, *Q. J.*
1092 *Roy. Meteor. Soc.*, 137(656), 553–597, doi:10.1002/qj.828.
- 1093 Denman, K. L., G. Brasseur, A. Chidthaisong, P. Ciais, P. M. Cox, R. E. Dickinson,
1094 D. Hauglustaine, C. Heinze, E. Holland, D. Jacob, U. Lohmann, S. Ramachandran,
1095 P. L. d. Dias, S. C. Wofsy, and X. Zhang (2007), Couplings between changes in the
1096 climate system and biogeochemistry, in *Climate Change 2007: The Physical Science Ba-*
1097 *sis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovern-*
1098 *mental Panel on Climate Change*, edited by S. Solomon, D. Qin, M. Manning, Z. Chen,
1099 M. Marquis, K. B. Averyt, M. Tignor, and H. L. Miller, pp. 499–587, Cambridge Uni-
1100 versity Press, Cambridge, United Kingdom and New York, NY, USA.
- 1101 Dietze, M. C., S. P. Serbin, C. Davidson, A. R. Desai, X. Feng, R. Kelly, R. Kooper,
1102 D. LeBauer, J. Mantooh, K. McHenry, and D. Wang (2014), A quantitative assessment

- 1103 of a terrestrial biosphere model's data needs across North American biomes, *J. Geo-*
1104 *phys. Res. Biogeosci.*, *119*(3), 286–300, doi:10.1002/2013JG002392.
- 1105 Dirmeyer, P. A., J. Wei, M. G. Bosilovich, and D. M. Mocko (2014), Comparing evapo-
1106 rative sources of terrestrial precipitation and their extremes in merra using relative en-
1107 tropy, *J. Hydrometeor.*, *15*(1), 102–116, doi:10.1175/JHM-D-13-053.1.
- 1108 Eaton, B., J. Gregory, B. Drach, K. Taylor, S. Hankin, J. Blower, J. Caron, R. Signell,
1109 P. Bentley, G. Rappa, H. Höck, A. Pamment, M. Jukes, M. Raspaud, and R. Horne
1110 (2017), Netcdf climate and forecast (cf) metadata conventions, <http://cfconventions.org/>.
- 1111 Eyring, V., M. Righi, A. Lauer, M. Evaldsson, S. Wenzel, C. Jones, A. Anav, O. An-
1112 drews, I. Cionni, E. L. Davin, C. Deser, C. Ehbrecht, P. Friedlingstein, P. Gleckler,
1113 K.-D. Gottschaldt, S. Hagemann, M. Jukes, S. Kindermann, J. Krasting, D. Kunert,
1114 R. Levine, A. Loew, J. Mäkelä, G. Martin, E. Mason, A. S. Phillips, S. Read, C. Rio,
1115 R. Roehrig, D. Senftleben, A. Sterl, L. H. van Ulft, J. Walton, S. Wang, and K. D.
1116 Williams (2016), ESMValTool (v1.0) – A community diagnostic and performance met-
1117 rics tool for routine evaluation of Earth system models in CMIP, *Geosci. Model Dev.*,
1118 *9*(5), 1747–1802, doi:10.5194/gmd-9-1747-2016.
- 1119 Friedlingstein, P., L. Bopp, P. Ciais, J.-L. Dufresne, L. Fairhead, H. LeTreut, P. Monfray,
1120 and J. Orr (2001), Positive feedback between future climate change and the carbon cy-
1121 cle, *Geophys. Res. Lett.*, *28*(8), 1543–1546, doi:10.1029/2000GL012015.
- 1122 Friedlingstein, P., P. M. Cox, R. A. Betts, L. Bopp, W. von Bloh, V. Brovkin, S. C. Doney,
1123 M. Eby, I. Fung, B. Govindasamy, J. John, C. D. Jones, F. Joos, T. Kato, M. Kawamiya,
1124 W. Knorr, K. Lindsay, H. D. Matthews, T. Raddatz, P. Rayner, C. Reick, E. Roeckner,
1125 K.-G. Schnitzler, R. Schnur, K. Strassmann, S. Thompson, A. J. Weaver, C. Yoshikawa,
1126 and N. Zeng (2006), Climate–carbon cycle feedback analysis: Results from the C⁴MIP
1127 model intercomparison, *J. Clim.*, *19*(14), 3373–3353, doi:10.1175/JCLI3800.1.
- 1128 Friedlingstein, P., M. Meinshausen, V. K. Arora, C. D. Jones, A. Anav, S. K. Liddicoat,
1129 and R. Knutti (2014), Uncertainties in CMIP5 climate projections due to carbon cycle
1130 feedbacks, *J. Clim.*, *27*(2), 511–526, doi:10.1175/JCLI-D-12-00579.1.
- 1131 Fung, I. Y., S. C. Doney, K. Lindsay, and J. John (2005), Evolution of carbon
1132 sinks in a changing climate, *Proc. Nat. Acad. Sci.*, *102*(32), 11,201–11,206, doi:
1133 10.1073/pnas.0504949102.
- 1134 Ghimire, B., W. J. Riley, C. D. Koven, M. Mu, and J. T. Randerson (2016), Representing
1135 leaf and root physiological traits in CLM improves global carbon and nitrogen cycling

- 1136 predictions, *J. Adv. Model. Earth Syst.*, 8(2), 598–613, doi:10.1002/2015MS000538.
- 1137 Giglio, L., J. T. Randerson, G. R. van der Werf, P. S. Kasibhatla, G. J. Collatz, D. C.
1138 Morton, and R. S. DeFries (2010), Assessing variability and long-term trends in burned
1139 area by merging multiple satellite fire products, *Biogeosci.*, 7(3), 1171–1186, doi:
1140 10.5194/bg-7-1171-2010.
- 1141 Gleckler, P. J., K. E. Taylor, and C. Doutriaux (2008), Performance metrics for climate
1142 models, *J. Geophys. Res.*, 113(D6), D06,104, doi:10.1029/2007JD008972.
- 1143 Gleckler, P. J., C. Doutriaux, P. J. Durack, K. E. Taylor, Y. Zhang, D. N. Williams, E. Ma-
1144 son, and J. Servonnat (2016), A more powerful reality test for climate models, *Eos*
1145 *Trans. AGU*, 97, doi:10.1029/2016EO051663.
- 1146 Gregory, J. M., C. D. Jones, P. Cadule, and P. Friedlingstein (2009), Quantifying carbon
1147 cycle feedbacks, *J. Clim.*, 22(19), 5232–5250, doi:10.1175/2009JCLI2949.1.
- 1148 Harris, I., P. Jones, T. Osborn, and D. Lister (2014), Updated high-resolution grids of
1149 monthly climatic observations – the cru ts3.10 dataset, *Int. J. Climatol.*, 34(3), 623–
1150 642, doi:10.1002/joc.3711.
- 1151 Hoffman, F. M., J. T. Randerson, V. K. Arora, Q. Bao, P. Cadule, D. Ji, C. D. Jones,
1152 M. Kawamiya, S. Khatiwala, K. Lindsay, A. Obata, E. Shevliakova, K. D. Six, J. F.
1153 Tjiputra, E. M. Volodin, and T. Wu (2014), Causes and implications of persistent at-
1154 mospheric carbon dioxide biases in Earth System Models, *J. Geophys. Res. Biogeosci.*,
1155 119(2), 141–162, doi:10.1002/2013JG002381.
- 1156 Hoffman, F. M., C. D. Koven, G. Keppel-Aleks, D. M. Lawrence, W. J. Riley, J. T. Ran-
1157 derson, A. Ahlström, G. Abramowitz, D. D. Baldocchi, M. J. Best, B. Bond-Lamberty,
1158 M. G. De Kauwe, A. S. Denning, A. R. Desai, V. Eyring, J. B. Fisher, R. A. Fisher,
1159 P. J. Gleckler, M. Huang, G. Hugelius, A. K. Jain, N. Y. Kiang, H. Kim, R. D. Koster,
1160 S. V. Kumar, H. Li, Y. Luo, J. Mao, N. G. McDowell, U. Mishra, P. R. Moorcroft,
1161 G. S. H. Pau, D. M. Ricciuto, K. Schaefer, C. R. Schwalm, S. P. Serbin, E. Shevliakova,
1162 A. G. Slater, J. Tang, M. Williams, J. Xia, C. Xu, R. Joseph, and D. Koch (2017), In-
1163 ternational Land Model Benchmarking (ILAMB) 2016 workshop report, *Tech. Rep.*
1164 *DOE/SC-0186*, U.S. Department of Energy, Office of Science, Germantown, Maryland,
1165 USA, doi:10.2172/1330803.
- 1166 Hugelius, G., J. G. Bockheim, P. Camill, B. Elberling, G. Grosse, J. W. Harden, K. John-
1167 son, T. Jorgenson, C. D. Koven, P. Kuhry, G. Michaelson, U. Mishra, J. Palmtag,
1168 C.-L. Ping, J. O’Donnell, L. Schirrmeister, E. A. G. Schuur, Y. Sheng, L. C. Smith,

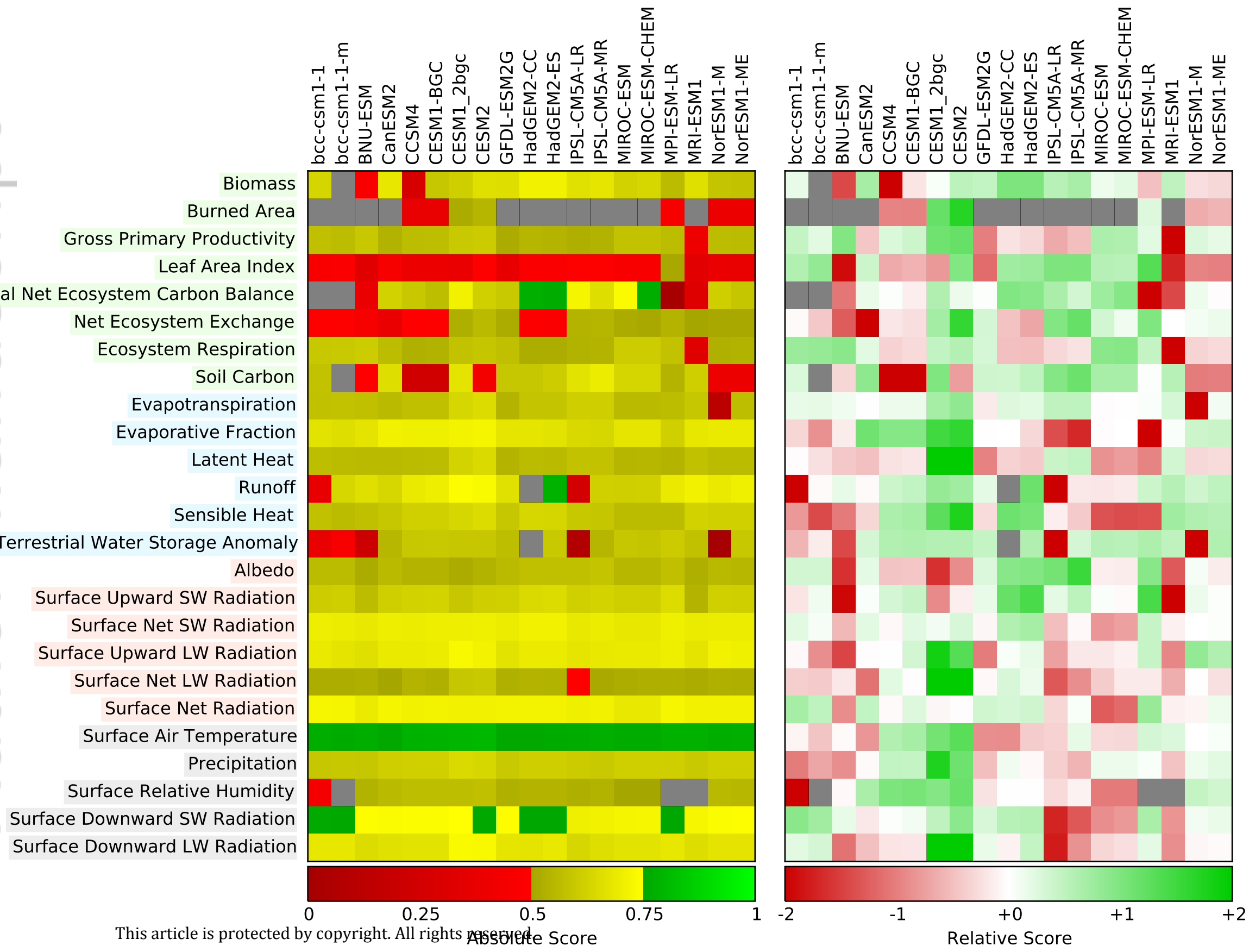
- 1169 J. Strauss, and Z. Yu (2013), A new data set for estimating organic carbon storage to
1170 3 m depth in soils of the northern circumpolar permafrost region, *Earth Syst. Sci. Data*,
1171 5(2), 393–402, doi:10.5194/essd-5-393-2013.
- 1172 Jung, M., M. Reichstein, P. Ciais, S. I. Seneviratne, J. Sheffield, M. L. Goulden, G. Bo-
1173 nan, A. Cescatti, J. Chen, R. de Jeu, A. J. Dolman, W. Eugster, D. Gerten, D. Gianelle,
1174 N. Gobron, J. Heinke, J. Kimball, B. E. Law, L. Montagnani, Q. Mu, B. Mueller,
1175 K. Oleson, D. Papale, A. D. Richardson, O. Roupsard, S. Running, E. Tomelleri,
1176 N. Viovy, U. Weber, C. Williams, E. Wood, S. Zaehle, and K. Zhang (2010), Recent
1177 decline in the global land evapotranspiration trend due to limited moisture supply, *Na-
1178 ture*, 467, 951–954, doi:10.1038/nature09396.
- 1179 Kato, S., N. G. Loeb, F. G. Rose, D. R. Doelling, D. A. Rutan, T. E. Caldwell, L. Yu,
1180 and R. A. Weller (2013), Surface irradiances consistent with CERES-derived top-
1181 of-atmosphere shortwave and longwave irradiances, *J. Clim.*, 26(9), 2719–2740, doi:
1182 10.1175/JCLI-D-12-00436.1.
- 1183 Kelley, D. I., I. C. Prentice, S. P. Harrison, H. Wang, M. Simard, J. B. Fisher, and K. O.
1184 Willis (2013), A comprehensive benchmarking system for evaluating global vegetation
1185 models, *Biogeosci.*, 10(5), 3313–3340, doi:10.5194/bg-10-3313-2013.
- 1186 Kellendorfer, J., W. Walker, K. Kirsch, G. Fiske, J. Bishop, L. Lapoint, M. Hoppus, and
1187 J. Westfall (2013), NACP aboveground biomass and carbon baseline data, V.2 (NBCD
1188 2000), U.S.A., 2000, doi:10.3334/ornl daac/1161.
- 1189 König-Langlo, G., R. Sieger, H. Schmithüsen, A. Bücker, F. Richter, and E. Dutton
1190 (2013), The Baseline Surface Radiation Network and its World Radiation Monitoring
1191 Centre at the Alfred Wegener Institute, *Tech. Rep. WCRP-2*, Alfred Wegener Institute,
1192 doi:10013/epic.42596.d001.
- 1193 Kumar, J., F. M. Hoffman, W. W. Hargrove, and N. Collier (2016), Understanding the rep-
1194 resentativeness of FLUXNET for upscaling carbon flux from eddy covariance measure-
1195 ments, *Earth System Science Data Discussions*, 2016, 1–25, doi:10.5194/essd-2016-36.
- 1196 Kumar, S. V., C. D. Peters-Lidard, J. Santanello, K. Harrison, Y. Liu, and M. Shaw
1197 (2012), Land surface Verification Toolkit (LVT) — A generalized framework for land
1198 surface model evaluation, *Geosci. Model Dev.*, 5(3), 869–886, doi:10.5194/gmd-5-869-
1199 2012.
- 1200 Lasslop, G., M. Reichstein, D. Papale, A. D. Richardson, A. Arneth, A. Barr, P. Stoy, and
1201 G. Wohlfahrt (2010), Separation of net ecosystem exchange into assimilation and respi-

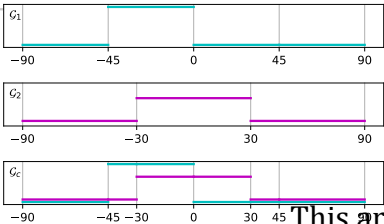
- 1202 ration using a light response curve approach: critical issues and global evaluation, *Glob.*
1203 *Change Biol.*, 16(1), 187–208, doi:10.1111/j.1365-2486.2009.02041.x.
- 1204 Law, K., A. Stuart, and K. Zygalkis (2015), *Data Assimilation: A Mathematical Intro-*
1205 *duction, Texts in Applied Mathematics*, vol. 63, 1 ed., 242 pp., Springer International
1206 Publishing, doi:10.1007/978-3-319-20325-6.
- 1207 Le Quéré, C., R. M. Andrew, J. G. Canadell, S. Sitch, J. I. Korsbakken, G. P. Peters,
1208 A. C. Manning, T. A. Boden, P. P. Tans, R. A. Houghton, R. F. Keeling, S. Alin,
1209 O. D. Andrews, P. Anthoni, L. Barbero, L. Bopp, F. Chevallier, L. P. Chini, P. Ciais,
1210 K. Currie, C. Delire, S. C. Doney, P. Friedlingstein, T. Gkritzalis, I. Harris, J. Hauck,
1211 V. Haverd, M. Hoppema, K. Klein Goldewijk, A. K. Jain, E. Kato, A. Körtzinger,
1212 P. Landschützer, N. Lefèvre, A. Lenton, S. Lienert, D. Lombardozzi, J. R. Melton,
1213 N. Metzl, F. Millero, P. M. S. Monteiro, D. R. Munro, J. E. M. S. Nabel, S.-I. Nakaoka,
1214 K. O'Brien, A. Olsen, A. M. Omar, T. Ono, D. Pierrot, B. Poulter, C. Rödenbeck,
1215 J. Salisbury, U. Schuster, J. Schwinger, R. Séférian, I. Skjelvan, B. D. Stocker, A. J.
1216 Sutton, T. Takahashi, H. Tian, B. Tilbrook, I. T. van der Laan-Luijkx, G. R. van der
1217 Werf, N. Viovy, A. P. Walker, A. J. Wiltshire, and S. Zaehle (2016), Global carbon bud-
1218 get 2016, *Earth Syst. Sci. Data*, 8(2), 605–649, doi:10.5194/essd-8-605-2016.
- 1219 LeBauer, D. S., D. Wang, K. T. Richter, C. C. Davidson, and M. C. Dietze (2013), Fa-
1220 cilitating feedbacks between field measurements and ecosystem models, *Ecol. Monogr.*,
1221 83(2), 133–154, doi:10.1890/12-0137.1.
- 1222 Luo, Y. Q., J. T. Randerson, G. Abramowitz, C. Bacour, E. Blyth, N. Carvalhais, P. Ciais,
1223 D. Dalmonech, J. B. Fisher, R. Fisher, P. Friedlingstein, K. Hibbard, F. Hoffman,
1224 D. Huntzinger, C. D. Jones, C. Koven, D. Lawrence, D. J. Li, M. Mahecha, S. L.
1225 Niu, R. Norby, S. L. Piao, X. Qi, P. Peylin, I. C. Prentice, W. Riley, M. Reichstein,
1226 C. Schwalm, Y. P. Wang, J. Y. Xia, S. Zaehle, and X. H. Zhou (2012), A framework for
1227 benchmarking land models, *Biogeosci.*, 9(10), 3857–3874, doi:10.5194/bg-9-3857-2012.
- 1228 Mahowald, N. M., J. T. Randerson, K. Lindsay, E. Muñoz, S. C. Doney, P. Lawrence,
1229 S. Schlunegger, D. S. Ward, D. Lawrence, and F. M. Hoffman (2017), Interactions be-
1230 tween land use change and carbon cycle feedbacks, *Global Biogeochem. Cycles*, 31(1),
1231 96–113, doi:10.1002/2016GB005374.
- 1232 Meehl, G. A., C. Covey, T. Delworth, M. Latif, B. McAvaney, J. F. B. Mitchell, R. J.
1233 Stouffer, and K. E. Taylor (2007), The WCRP CMIP3 multimodel dataset: A new
1234 era in climate change research, *Bull. Am. Meteorol. Soc.*, 88(9), 1383–1394, doi:

- 1235 10.1175/BAMS-88-9-1383.
- 1236 Miralles, D., T. Holmes, R. D. Jeu, J. Gash, A. Meesters, and A. Dolman (2011), Global
1237 land-surface evaporation estimated from satellite-based observations, *Hydrol. Earth Syst.*
1238 *Sci.*, *15*, 453–469.
- 1239 Moore, J. K., W. Fu, F. Primeau, G. L. Britten, K. Lindsay, M. Long, S. C. Doney,
1240 N. Mahowald, F. M. Hoffman, and J. T. Randerson (2018), Sustained climate warm-
1241 ing drives declining marine biological productivity, *Science*, *359*(6380), 1139–1143,
1242 doi:10.1126/science.aao6379.
- 1243 Mu, M., J. T. Randerson, W. J. Riley, C. D. Koven, G. Keppel-Aleks, D. M. Lawrence, ,
1244 and F. M. Hoffman (2015), International Land Model Benchmarking (ILAMB) package
1245 v1, online, doi:10.18139/ILAMB.v001.00/1251597.
- 1246 Myneni, R. B., R. R. Nemani, and S. Running (1997), Algorithm for the estimation of
1247 global land cover, LAI and FPAR based on radiative transfer models, *IEEE Trans.*
1248 *Geosci. Remote Sens.*, *35*, 1380–1392.
- 1249 Oleson, K. W., D. M. Lawrence, G. B. Bonan, B. Drewniak, M. Huang, C. D. Koven,
1250 S. Levis, F. Li, W. J. Riley, Z. M. Subin, S. C. Swenson, P. E. Thornton, A. Bozbiyik,
1251 R. Fisher, C. L. Heald, E. Kluzek, J.-F. Lamarque, P. J. Lawrence, L. R. Leung, W. Lip-
1252 scomb, S. Muszala, D. M. Ricciuto, W. Sacks, Y. Sun, J. Tang, and Z.-L. Yang (2013),
1253 Technical description of version 4.5 of the Community Land Model (CLM), *Techni-*
1254 *cal Note NCAR/TN-503+STR*, National Center for Atmospheric Research, Boulder, Col-
1255 orado, USA.
- 1256 Piao, S., S. Sitch, P. Ciais, P. Friedlingstein, P. Peylin, X. Wang, A. Ahlström, A. Anav,
1257 J. G. Canadell, N. Cong, C. Huntingford, M. Jung, S. Levis, P. E. Levy, J. Li, X. Lin,
1258 M. R. Lomas, M. Lu, Y. Luo, Y. Ma, R. B. Myneni, B. Poulter, Z. Sun, T. Wang,
1259 N. Viovy, S. Zaehle, and N. Zeng (2013), Evaluation of terrestrial carbon cycle models
1260 for their response to climate variability and to CO₂ trends, *Glob. Change Biol.*, *19*(7),
1261 2117–2132, doi:10.1111/gcb.12187.
- 1262 Randerson, J. T., F. M. Hoffman, P. E. Thornton, N. M. Mahowald, K. Lindsay, Y.-H.
1263 Lee, C. D. Nevison, S. C. Doney, G. Bonan, R. Stöckli, C. Covey, S. W. Running,
1264 and I. Y. Fung (2009), Systematic assessment of terrestrial biogeochemistry in cou-
1265 pled climate–carbon models, *Glob. Change Biol.*, *15*(9), 2462–2484, doi:10.1111/j.1365-
1266 2486.2009.01912.x.

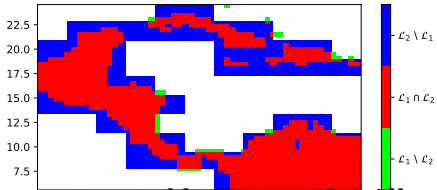
- 1267 Randerson, J. T., K. Lindsay, E. Munoz, W. Fu, J. K. Moore, F. M. Hoffman, N. M. Ma-
1268 howald, and S. C. Doney (2015), Multicentury changes in ocean and land contribu-
1269 tions to the climate–carbon feedback, *Global Biogeochem. Cycles*, 29(6), 744–759, doi:
1270 10.1002/2014GB005079.
- 1271 Reichler, T., and J. Kim (2008), How well do coupled models simulate today’s climate?,
1272 *Bull. Am. Meteorol. Soc.*, 89(3), 303–311, doi:10.1175/BAMS-89-3-303.
- 1273 Saatchi, S. S., N. L. Harris, S. Brown, M. Lefsky, E. T. A. Mitchard, W. Salas,
1274 B. R. Zutta, W. Buermann, S. L. Lewis, S. Hagen, S. Petrova, L. White, M. Sil-
1275 man, and A. Morel (2011), Benchmark map of forest carbon stocks in tropical
1276 regions across three continents, *Proc. Nat. Acad. Sci.*, 108(24), 9899–9904, doi:
1277 10.1073/pnas.1019576108.
- 1278 Schneider, U., A. Becker, P. Finger, A. Meyer-Christoffer, M. Ziese, and B. Rudolf (2014),
1279 GPCP’s new land surface precipitation climatology based on quality-controlled in situ
1280 data and its role in quantifying the global water cycle, *Theor. Appl. Climatol.*, 115(1),
1281 15–40, doi:10.1007/s00704-013-0860-x.
- 1282 Stackhouse Jr., P. W., S. K. Gupta, S. J. Cox, J. C. Mikovitz, T. Zhang, and L. M. Hinkel-
1283 man (2011), The NASA/GEWEX surface radiation budget release 3.0: 24.5-year
1284 dataset, *GEWEX News*, 21(1), 10–12.
- 1285 Swenson, S. (2013), GRACE: Gravity recovery and climate experiment: Surface mass, to-
1286 tal water storage, and derived variables, in *The Climate Data Guide*, edited by National
1287 Center for Atmospheric Research Staff, National Center for Atmospheric Research.
- 1288 Swenson, S., and J. Wahr (2006), Post-processing removal of correlated errors in GRACE
1289 data, *Geophys. Res. Lett.*, 33(8), L08,402, doi:10.1029/2005GL025285.
- 1290 Taylor, K. E. (2001), Summarizing multiple aspects of model performance in a single dia-
1291 gram, *J. Geophys. Res. Atmos.*, 106(D7), 7183–7192, doi:10.1029/2000JD900719.
- 1292 Taylor, K. E., R. J. Stouffer, and G. A. Meehl (2012), An overview of CMIP5 and the ex-
1293 periment design, *Bull. Am. Meteorol. Soc.*, 93(4), 485–498, doi:10.1175/BAMS-D-11-
1294 00094.1.
- 1295 Todd-Brown, K. E. O., J. T. Randerson, W. M. Post, F. M. Hoffman, C. Tarnocai, E. A. G.
1296 Schuur, and S. D. Allison (2013), Causes of variation in soil carbon simulations from
1297 CMIP5 Earth system models and comparison with observations, *Biogeosci.*, 10(3),
1298 1717–1736, doi:10.5194/bg-10-1717-2013.

- 1299 Walker, A. P., P. J. Hanson, M. G. De Kauwe, B. E. Medlyn, S. Zaehle, S. Asao, M. Di-
1300 etze, T. Hickler, C. Huntingford, C. M. Iversen, A. Jain, M. Lomas, Y. Luo, H. Mc-
1301 Carthy, W. J. Parton, I. C. Prentice, P. E. Thornton, S. Wang, Y.-P. Wang, D. Warlind,
1302 E. Weng, J. M. Warren, F. I. Woodward, R. Oren, and R. J. Norby (2014), Comprehen-
1303 sive ecosystem model–data synthesis using multiple data sets at two temperate forest
1304 free-air CO₂ enrichment experiments: Model performance at ambient CO₂ concentra-
1305 tion, *J. Geophys. Res. Biogeosci.*, *119*(5), 2169–8961, doi:10.1002/2013JG002553.
- 1306 Walker, A. P., S. Zaehle, B. E. Medlyn, M. G. De Kauwe, S. Asao, T. Hickler, W. Par-
1307 ton, D. M. Ricciuto, Y.-P. Wang, D. Wärlind, and R. J. Norby (2015), Predicting
1308 long-term carbon sequestration in response to CO₂ enrichment: How and why do
1309 current ecosystem models differ?, *Global Biogeochem. Cycles*, *29*(4), 476–495, doi:
1310 10.1002/2014GB004995.
- 1311 Xie, P., and P. A. Arkin (1997), Global precipitation: A 17-year monthly anal-
1312 ysis based on gauge observations, satellite estimates, and numerical model
1313 outputs, *Bull. Am. Meteorol. Soc.*, *78*(11), 2539–2558, doi:10.1175/1520-
1314 0477(1997)078<2539:GPAYMA>2.0.CO;2.
- 1315 Zaehle, S., B. E. Medlyn, M. G. De Kauwe, A. P. Walker, M. C. Dietze, T. Hickler,
1316 Y. Luo, Y.-P. Wang, B. El-Masri, P. Thornton, A. Jain, S. Wang, D. Warlind, E. Weng,
1317 W. Parton, C. M. Iversen, A. Gallet-Budynek, H. McCarthy, A. Finzi, P. J. Hanson,
1318 I. C. Prentice, R. Oren, and R. J. Norby (2014), Evaluation of 11 terrestrial carbon–
1319 nitrogen cycle models against observations from two temperate free-air CO₂ enrichment
1320 studies, *New Phytol.*, *202*(3), 803–822, doi:10.1111/nph.12697.
- 1321 Zhu, Q., W. J. Riley, J. Tang, and C. D. Koven (2016), Multiple soil nutrient competition
1322 between plants, microbes, and mineral surfaces: Model development, parameterization,
1323 and example applications in several tropical forests, *Biogeosci.*, *13*(1), 341–363, doi:
1324 10.5194/bg-13-341-2016.



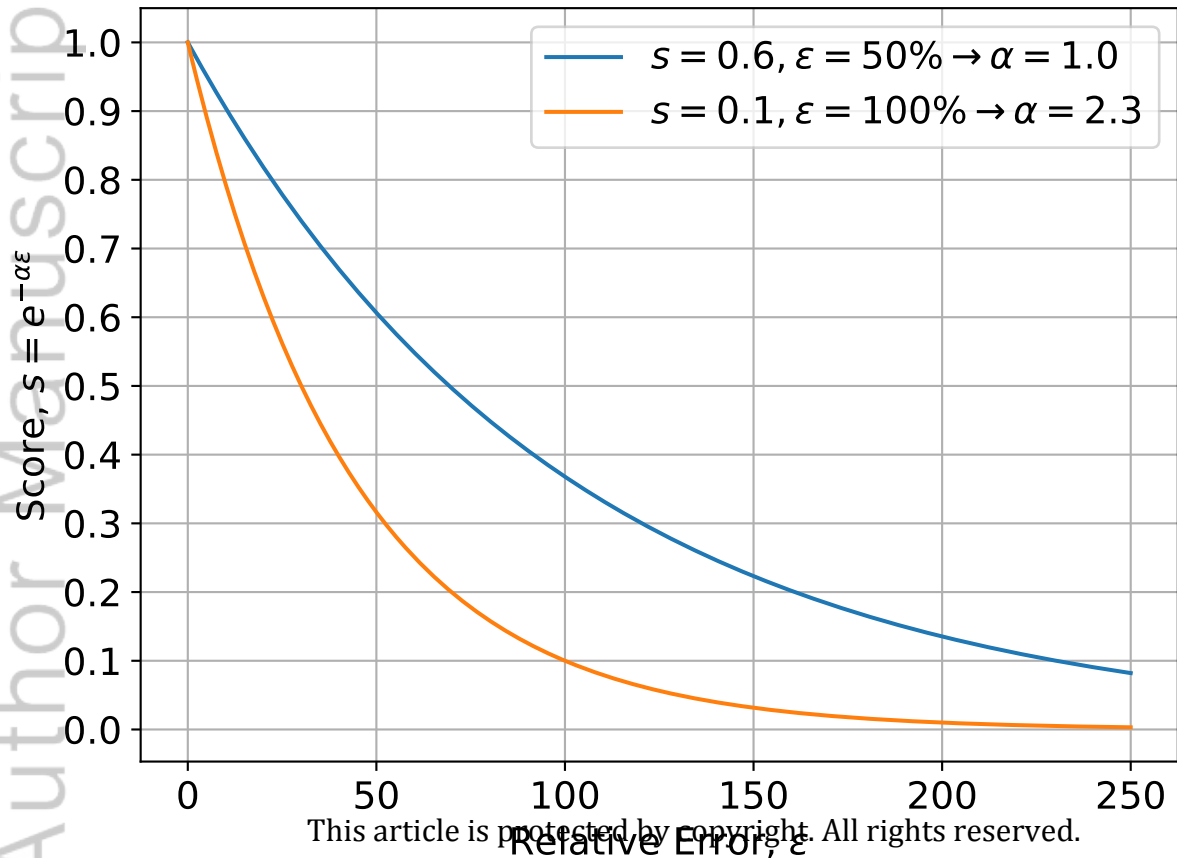


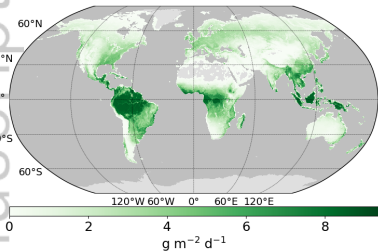
(a)



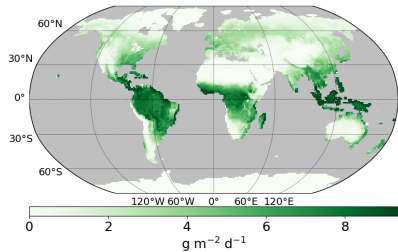
(b)

This article is protected by copyright. All rights reserved.

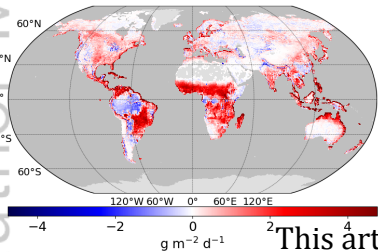




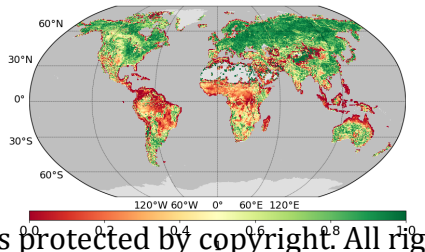
(a)



(b)

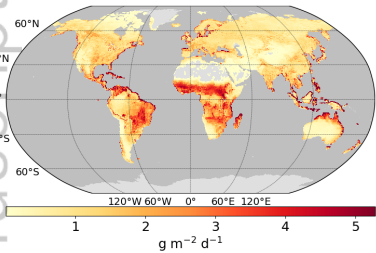


(c)

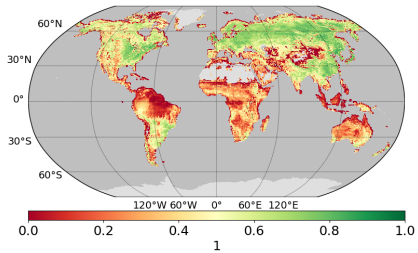


(d)

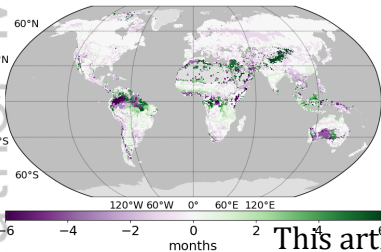
This article is protected by copyright. All rights reserved.



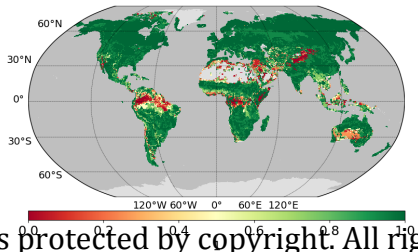
(a)



(b)

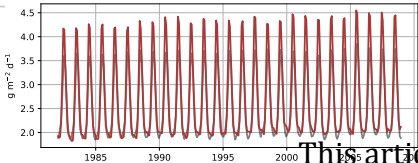


(c)

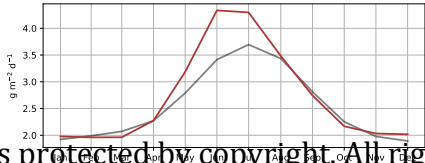


(d)

This article is protected by copyright. All rights reserved.

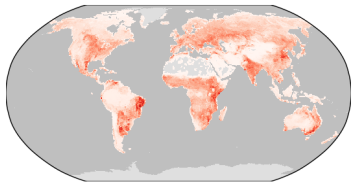


(a)

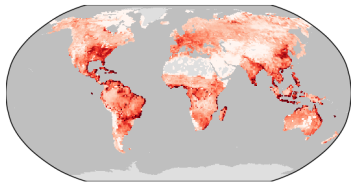


(b)

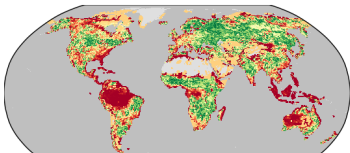
This article is protected by copyright. All rights reserved.



(a)

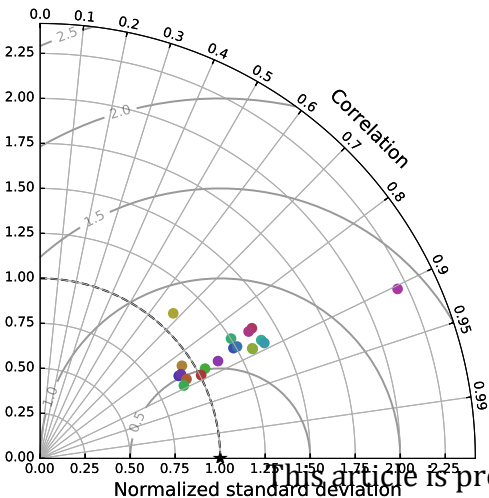


(b)



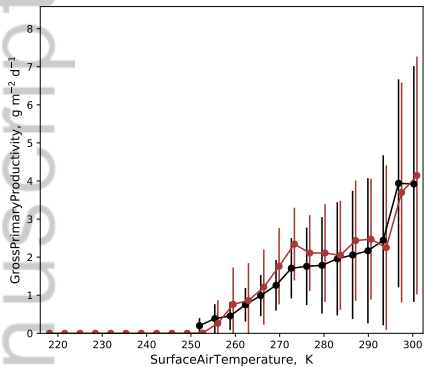
(c)

This article is protected by copyright. All rights reserved.

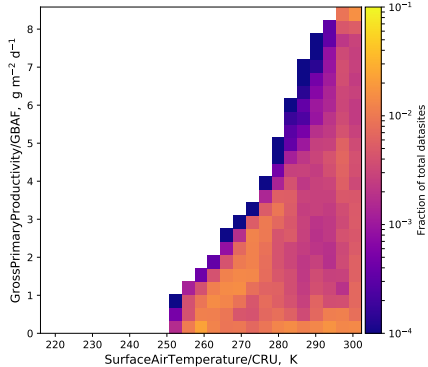


This article is protected by copyright. All rights reserved.

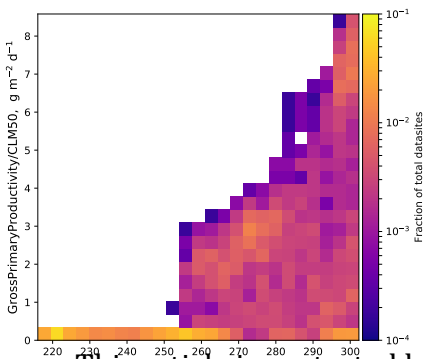
- bcc-csm1-1
- bcc-csm1-1-m
- BNU-ESM
- CanESM2
- CCSM4
- CESM1-BGC
- CESM1_2bgc
- CESM2
- GFDL-ESM2G
- HadGEM2-CC
- HadGEM2-ES
- inmcm4
- IPSL-CM5A-LR
- IPSL-CM5A-MR
- MIROC-ESM-CHEM
- MIROC-ESM
- MPI-ESM-LR
- MRI-ESM1
- NorESM1-M
- NorESM1-ME



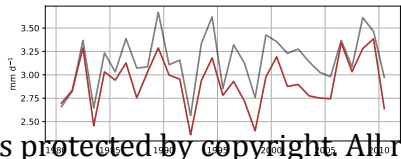
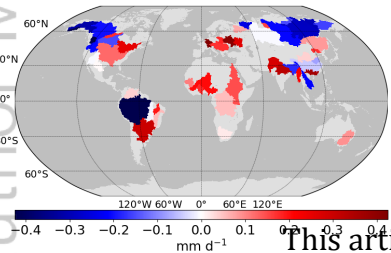
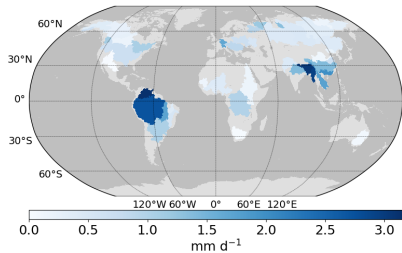
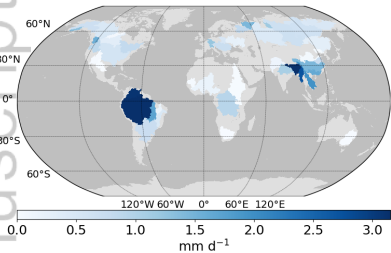
(a)

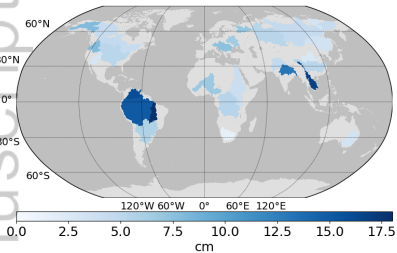


(b)

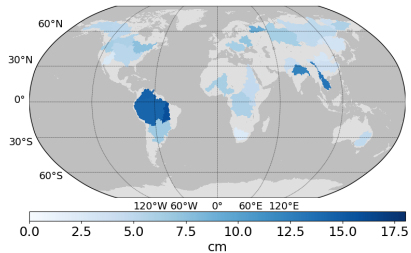


(c)

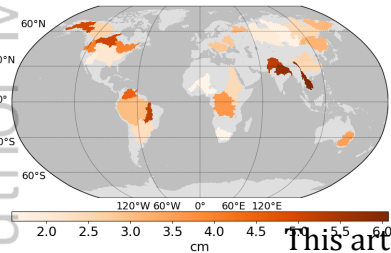




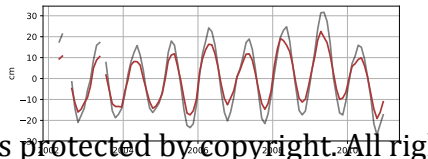
(a)



(b)



(c)



(d)

This article is protected by copyright. All rights reserved.