

**A model-data intercomparison of CO₂ exchange across North America: Results
from the North American Carbon Program Site Synthesis**

Manuscript for submission to JGR-Biogeosciences

Original version: October 7, 2009

Original submission: November 24, 2009

Revised submission: June 25, 2010

Second revision: July 23, 2010

Corresponding author:

Christopher R. Schwalm¹

¹ Graduate School of Geography, Clark University, Worcester, MA 01610, USA

cschwalm@clarku.edu, Tel: 508-793-7711

Contributing authors:

Christopher A. Williams², Kevin Schaefer³, Ryan Anderson⁴, M. Altaf Arain⁵, Ian Baker⁶, Alan Barr⁷, T. Andrew Black⁸, Guangsheng Chen⁹, Jing Ming Chen¹⁰, Philippe Ciais¹¹, Kenneth J. Davis¹², Ankur Desai¹³, Michael Dietze¹⁴, Danilo Dragoni¹⁵, Marc L. Fischer¹⁶, Lawrence B. Flanagan¹⁷, Robert Grant¹⁸, Lianhong Gu¹⁹, David Hollinger²⁰, R. César Izaurralde²¹, Chris Kucharik²², Peter Lafleur²³, Beverly E. Law²⁴, Longhui Li²⁵, Zhengpeng Li²⁶, Shuguang Liu²⁷, Erandathie Lokupitiya²⁸, Yiqi Luo²⁹, Siyan Ma³⁰, Hank Margolis³¹, Roser Matamala³², Harry McCaughey³³, Russell K. Monson³⁴, Walter C. Oechel³⁵, Changhui Peng³⁶, Benjamin Poulter³⁷, David T. Price³⁸, Dan M. Riciutto³⁹, William Riley⁴⁰, Alok Kumar Sahoo⁴¹, Michael Sprintsin⁴², Jianfeng Sun⁴³, Hanqin Tian⁴⁴, Christina Tonitto⁴⁵, Hans Verbeeck⁴⁶, Shashi B. Verma⁴⁷

² Graduate School of Geography, Clark University, Worcester, MA 01610, USA

cwilliams@clarku.edu, Tel: 508-793-7323

³ National Snow and Ice Data Center, University of Colorado at Boulder, Boulder, CO 80309, USA

kevin.schaefer@nsidc.org, Tel: 303-492-8869

⁴ Numerical Terradynamic Simulation Group, University of Montana, Missoula, MT 59812, USA

ryan.anderson@ntsg.umt.edu, Tel: (406) 243-6263

⁵ School of Geography and Earth Sciences, McMaster University, Hamilton, ON L8S 4K1, Canada

arainm@mcmaster.ca, Tel: 905-525-9140

⁶ Atmospheric Science Department, Colorado State University, Fort Collins, CO 80523,
USA

baker@atmos.colostate.edu, Tel: 970-491-4948

⁷ Climate Research Division, Atmospheric Science and Technology Directorate,
Saskatoon, SK S7N 3H5, Canada

alan.barr@ec.gc.ca, Tel: 306-975-4324

⁸ Faculty of Land and Food Systems, University of British Columbia, Vancouver, BC
V6T 1Z4, Canada

andrew.black@ubc.ca, Tel: 604-822-2730

⁹ School of Forestry and Wildlife Sciences, Auburn University, Auburn, AL 36849, USA
chengul@auburn.edu, Tel: 334-844-8057

¹⁰ Department of Geography and Program in Planning, University of Toronto, Toronto,
ON M5S 3G3, Canada

chenj@geog.utoronto.ca, Tel: 416-978-7085

¹¹ Laboratoire des Sciences du Climat et de l'Environnement, CE Orme des Merisiers, Gif
sur Yvette, 91191 France

philippe.ciais@cea.fr, Tel: +33 1 6908 9506

¹² Department of Meteorology, Pennsylvania State University, University Park, PA
16802, USA

davis@meteo.psu.edu, Tel: 814-863-8601

¹³ Center for Climatic Research, University of Wisconsin - Madison, Madison, WI 53706,
USA

desai@aos.wisc.edu, Tel: 608-265-9201

- ¹⁴ Department of Plant Biology, University of Illinois - Urbana Champaign, Urbana, IL
61801, USA
mdietze@life.uiuc.edu, Tel: 217-265-8020
- ¹⁵ Department of Geography, Indiana University, Bloomington, IN 47405, USA
ddragoni@indiana.edu, Tel: 812-855-5557
- ¹⁶ Atmospheric Science Department, Lawrence Berkeley National Laboratory, Berkeley,
CA 94720, USA
mlfischer@lbl.gov, Tel: 510-486-5539
- ¹⁷ Department of Biological Sciences, University of Lethbridge, Lethbridge, AB T1K
3M4, Canada
larry.flanagan@uleth.ca, Tel: 403-380-1858.
- ¹⁸ Department of Renewable Resources, University of Alberta, Edmonton, AB T6G 2E3,
Canada
robert.grant@afhe.ualberta.ca, Tel: 780-492-6609
- ¹⁹ Environmental Sciences Division, Oak Ridge National Laboratory, Oak Ridge, TN
37831, USA
lianhong-gu@ornl.gov, Tel: 865-241-5925
- ²⁰ Northern Research Station, USDA Forest Service, NH 03824, USA
davidh@hypatia.unh.edu, Tel: 603-868-7673
- ²¹ Pacific Northwest National Laboratory and University of Maryland, College Park, MD
20740, USA
cesar.izaurralde@pnl.gov, Tel: 301-314-6751

- ²² Department of Agronomy & Nelson Institute Center for Sustainability and the Global Environment, University of Wisconsin - Madison, Madison, WI 53706, USA
kucharik@wisc.edu, Tel: 608-263-1859
- ²³ Department of Geography, Trent University, Peterborough, ON K9J 7B8, Canada
plafleur@trentu.ca, Tel: 705-748-1011
- ²⁴ College of Forestry, Oregon State University, Corvallis, OR 97331, USA
bev.law@oregonstate.edu, Tel: 541-737-6111
- ²⁵ Laboratoire des Sciences du Climat et de l'Environnement, CE Orme des Merisiers, Gif sur Yvette, 91191 France
longhui.li@lsce.ipsl.fr, Tel: +33 1 6908 4102
- ²⁶ ASRC Research and Technology Solutions, Sioux Falls, SD 57198, USA
zli@usgs.gov, Tel: 605-594-6864
- ²⁷ Earth Resources Observation and Science, Sioux Falls, SD 57198, USA
sliu@usgs.gov, Tel: 605-594-6168
- ²⁸ Department of Atmospheric Science, Colorado State University, Fort Collins, CO 80523, USA
erandi@atmos.colostate.edu, Tel: 970-491-8915
- ²⁹ Department of Botany and Microbiology, University of Oklahoma, Norman, OK 73019, USA
yluo@ou.edu, Tel: 405-325-1651
- ³⁰ Department of Environmental Science, Policy and Management & Berkeley Atmospheric Science Center, University of California, Berkeley, Berkeley, CA 94720, USA

sma@berkeley.edu, Tel: 510-642-2421

- ³¹ Centre d'études de la forêt, Faculté de foresterie, de géographie et de géomatique,
Université Laval, Québec, QC G1V 0A6, Canada
hank.margolis@sbf.ulaval.ca, Tel: 418-656-7120
- ³² Argonne National Laboratory, Biosciences Division, Argonne, IL 60439, USA
matamala@anl.gov, Tel: 630-252-9270
- ³³ Department of Geography, Queen's University Kingston, ON K7L 3N6, Canada
mccaughe@post.queensu.ca, Tel: 613-533-6035
- ³⁴ Department of Ecology and Evolutionary Biology, University of Colorado at Boulder,
Boulder, CO 80309, USA
russell.monson@colorado.edu, Tel: 303-492-6319
- ³⁵ Department of Biology, San Diego State University, San Diego, CA 92182, USA
oechel@sunstroke.sdsu.edu, Tel: 619-594-4818
- ³⁶ Department of Biology Sciences, University of Quebec at Montreal, Montreal, QC
H3C 3P8, Canada
peng.changhui@uqam.ca, Tel: 514-987-3000
- ³⁷ Swiss Federal Research Institute WSL, Birmensdorf, CH-8903, Switzerland
benjamin.poulter@wsl.ch, Tel: +41 44 7392 215
- ³⁸ Northern Forestry Centre, Canadian Forest Service, Edmonton, AB T6H 3S5, Canada
dprice@nrcan.gc.ca, Tel: 780-435-7249
- ³⁹ Environmental Sciences Division, Oak Ridge National Laboratory, Oak Ridge, TN
37831, USA
ricciutodm@ornl.gov, Tel: 865-574-7067

- ⁴⁰ Climate and Carbon Sciences, Earth Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA
wjriley@lbl.gov, Tel: 510-486-5036
- ⁴¹ Department of Civil and Environmental Engineering, Princeton University, Princeton NJ 08544, USA
sahoo@Princeton.edu, Tel: 609-258-6383
- ⁴² Department of Geography and Program in Planning, University of Toronto, Toronto, ON M5S 3G3, Canada
misprin@gmail.com, Tel: 416-946-7715
- ⁴³ Department of Biology Sciences, University of Quebec at Montreal, Montreal, QC H3C 3P8, Canada
jianfeng_sun@yahoo.ca, Tel: 514-987-3000
- ⁴⁴ School of Forestry and Wildlife Sciences, Auburn University, Auburn, AL 36849, USA
tianhan@auburn.edu, Tel: 334-844-1059
- ⁴⁵ Department of Ecology and Evolutionary Biology, Cornell University, Ithaca, NY 14853, USA
ctonitto@cornell.edu, Tel: 607-254-4257
- ⁴⁶ Laboratory of Plant Ecology, Ghent University, 9000 Ghent, Belgium
hans.verbeeck@ugent.be, Tel: +32 9 264 61 13
- ⁴⁷ School of Natural Resources, University of Nebraska-Lincoln, Lincoln, NE 68583, USA
svermal@unl.edu, Tel: 402-472-6702

Abstract

Our current understanding of terrestrial carbon processes is represented in various models used to integrate and scale measurements of CO₂ exchange from remote sensing and other spatiotemporal data. Yet assessments are rarely conducted to determine how well models simulate carbon processes across vegetation types and environmental conditions. Using standardized data from the North American Carbon Program we compare observed and simulated monthly CO₂ exchange from 44 eddy covariance flux towers in North America and 22 terrestrial biosphere models. The analysis period spans ~220 site-years, 10 biomes, and includes two large scale drought events, providing a natural experiment to evaluate model skill as a function of drought and seasonality. We evaluate models' ability to simulate the seasonal cycle of CO₂ exchange using multiple model skill metrics and analyze links between model characteristics, site history, and model skill. Overall model performance was poor; the difference between observations and simulations was ~10-times observational uncertainty, with forested ecosystems better predicted than non-forested. Model-data agreement was highest in summer and in temperate evergreen forests. In contrast, model performance declined in spring and fall, especially in ecosystems with large deciduous components, and in dry periods during the growing season. Models used across multiple biomes and sites, the mean model ensemble, and a model using assimilated parameter values showed high consistency with observations. Models with the highest skill across all biomes all used prescribed canopy phenology, calculated *NEE* as the difference between *GPP* and ecosystem respiration, and did not use a daily time step.

Keywords: carbon modeling, ecosystem models, model validation, carbon exchange, drought, North American Carbon Program

Introduction

There is a continued need for models to improve consistency and agreement with observations [Friedlingstein et al., 2006], both overall and under more frequent extreme climatic events related to global environmental change such as drought [Trenberth et al., 2007]. Past validation studies of terrestrial biosphere models have focused only on few models and sites, typically in close proximity and primarily in forested biomes [e.g., Amthor et al., 2001; Delpierre et al., 2009; Grant et al., 2005; Hanson et al., 2004; Granier et al. 2007; Ichii et al., 2009; Ito, 2008; Siqueira et al., 2006; Zhou et al., 2008]. Furthermore, assessing model-data agreement relative to drought requires, in addition to high quality observed CO₂ exchange data, a reliable drought metric as well as a natural experiment across sites and drought conditions.

Drought is a reoccurring phenomenon in all climates [Larcher, 1995] and is characterized by a partial loss in plant function due to water limitation and heat stress. For terrestrial CO₂ exchange, drought typically reduces photosynthesis more than respiration [Baldocchi, 2008; Ciais et al., 2005; Schwalm et al., 2009], resulting in decreased net carbon uptake from the atmosphere. In the recent past drought conditions have become more prevalent globally [Dai et al., 2004] and in North America [Cook et al., 2004b]. Both incidence and severity of drought [Seager et al., 2007b] as well as heatwaves [Meehl et al., 2004] are expected to further increase in conjunction with global warming

[Houghton et al., 2001; Huntington, 2006; Sheffield & Wood, 2008; Trenberth et al., 2007].

In this study, we evaluate model performance using terrestrial CO₂ flux data and simulated fluxes collected from 1991 to 2007. This timeframe included two widespread droughts in North America: (i) the turn-of-the-century drought from 1998 to 2004 that was centered in the western interior of North America [Seager et al., 2007a] and (ii) a smaller-scale drought event in the southern continental United States from winter of 2005/2006 through October 2007 [Seager et al., 2009]. During these events Palmer Drought Severity Index values [Cook et al., 2007; Dai et al., 2004] and precipitation anomalies [Seager et al., 2007a; 2009] were highly negative over broad geographic areas. Ongoing eddy covariance measurements [Baldocchi et al., 2001], active throughout the aforementioned drought periods, provided flux data across gradients of time, space, seasonality, and drought. We use these data to examine model skill relative to site-specific drought severity, climatic season, and time. We also link model behavior to model architecture and site-specific attributes. Specifically, we address the following questions: Are current state-of-the-art terrestrial biosphere models capable of simulating CO₂ exchange subject to gradients in dryness and seasonality? Are these models able to reproduce the seasonal variation of observed CO₂ exchange across sites? Are certain characteristics of model structure coincident with better model-data agreement? Which biomes are simulated poorly/well?

Methods

Observed and simulated CO₂ exchange

Modeled and observed net ecosystem exchange (*NEE*, net carbon balance including soils where positive values indicate outgassing of CO₂ to the atmosphere) data were analyzed from 21 terrestrial biosphere models (Table 1) and 44 eddy covariance (EC) sites spanning ~220 site-years and 10 biomes in North America (Table 2). All terrestrial biosphere models analyzed simulated carbon cycling with process based formulations of varying detail for component carbon fluxes. Simulated *NEE* was based on model-specific runs using gap-filled observed weather at each site and locally observed values of soil texture according to a standard protocol [Ricciuto et al., 2009; <http://www.nacarbon.org/nacp/>], including a target *NEE* of zero integrated over the last five years of the simulation period. In addition, a mean model ensemble (hereafter: MEAN) was also analyzed. MEAN was calculated as the mean monthly value across all simulations. Furthermore, in contrast to other models, the parameter values used in the model LoTEC were optimized using data assimilation [Ricciuto et al., 2008]. LoTEC simulations were however retained when calculating MEAN as their effect on model skill was negligible due to the relatively small number of site-months simulated.

Gaps in the meteorological data record occurred at EC sites due to data quality control or instrument failure. Missing values of air temperature, humidity, shortwave radiation, and precipitation data, i.e., key model inputs, were filled using DAYMET [Thornton et al., 1997] before 2003 or the nearest available climate station in the National Climatic Data Center's Global Surface Summary of the Day (GSOD) database. Daily GSOD and DAYMET data were temporally downscaled to hourly or half-hourly values using the

phasing from observed mean diurnal cycles calculated from a 15-day moving window. The phasing used a sine wave assuming peak values at 15:00 local standard time (LST) and lowest values at 3:00 LST. In the absence of station data a 10-day running mean diurnal cycle was used [Ricciuto et al., 2009; http://nacp.ornl.gov/docs/Site_Synthesis_Protocol_v7.pdf].

EC data were produced by AmeriFlux and Fluxnet-Canada investigators and processed as a synthesis product of the North American Carbon Program (NACP) Site Level Interim Synthesis [<http://www.nacarbon.org/nacp/>]. The observed *NEE* were corrected for storage, despiked (i.e., outlying values removed), filtered to remove conditions of low turbulence (friction velocity filtered), and gap-filled to create a continuous time series [Barr et al., 2004]. The time series included estimates of random uncertainty and uncertainty due to friction velocity filtering [Barr et al., 2004; 2009]. In this analysis, *NEE* was aggregated to monthly values using only non-gap-filled data, i.e., observed values deemed spurious and subsequently infilled were not considered. Coincident modeled *NEE* values were similarly excluded. This removed the influence of gap-filling algorithms in the comparison of observed and modeled *NEE*.

Drought level was quantified using the 3-month Standard Precipitation Index (SPI, McKee et al., 1993). Monthly SPI values were taken from the U.S. Drought Monitor [<http://drought.unl.edu/DM/>] whereby each tower was matched to nearby meteorological station(s) indicative of local drought conditions given proximity, topography, and human impact. This study used three drought levels: dry required SPI < -0.8, wet corresponded

to $SPI > +0.8$, otherwise normal conditions existed. Climatic season was defined by four seasons of three months each with winter given by December, January, and February.

Model skill

Model-data mismatch was evaluated using normalized mean absolute error (*NMAE*) [Medlyn et al., 2005], the reduced χ^2 statistic (χ^2) [Taylor, 1996] as well as Taylor diagrams and skill (*S*) [Taylor, 2001]. The first metric quantifies bias, the “average distance” between observations and simulations in units of observed mean *NEE*:

$$NMAE = \sum_{ijkl} \frac{NEE_{obs} - NEE_{sim}}{n \overline{NEE}_{obs}}, \quad (1)$$

where the overbar indicates averaging across all values, n is sample size, the subscript *obs* is for observations and *sim* is for modeled estimates. The summation is for any arbitrary data group (denoted by subscripts on the summation operator only) where subscript i is for site, j is for model, k is for climatic season, l is for drought level.

The second metric used to evaluate model performance was the reduced χ^2 statistic. This is the squared difference between paired model and data points over observational error normalized by degrees of freedom:

$$\chi^2 = \frac{1}{n} \sum_{ijkl} \left(\frac{NEE_{obs} - NEE_{sim}}{4\delta_{NEE}} \right)^2, \quad (2)$$

where δ_{NEE} is uncertainty of monthly *NEE* (see below), 4 normalizes the uncertainty in observed *NEE* to correspond to a 95% confidence interval, the summation is across any arbitrary data group (denoted by subscripts on the summation operator). χ^2 values are linked to model-data mismatch where a value of unity indicates that model and data are in agreement relative to data uncertainty.

A final characterization of model performance used Taylor diagrams [Taylor, 2001]; visual displays based on pattern matching, i.e., the degree to which simulations matched the temporal evolution of monthly *NEE*. Taylor plots are polar coordinate displays of the linear correlation coefficient (ρ), centered root mean squared error (*RMSE*; pattern error without considering bias), and the standard deviation of *NEE* (σ). Taylor diagrams were constructed for the mean model ensemble (MEAN) and across-site mean model performance using the full data record for each combination of site and model (ranging from 7 to 178 months). More generally, each polar coordinate point for any arbitrary data group can be scored:

$$S = \frac{2(1 + \rho)}{(\sigma_{norm} + 1/\sigma_{norm})^2}, \quad (3)$$

where S is the model skill metric bound by zero and unity where unity indicates perfect agreement, and σ_{norm} is the ratio of simulated to observed standard deviation [Taylor, 2001].

To scale model skill metrics across gradients of site, biome, model, seasonality, and dryness level we aggregated across data groups weighting each by sample size. For example, χ^2 for model I , denoted by subscript $j = I$, is given by:

$$\chi_{j=I}^2 = \sum_{ikl} \frac{n_{ikl} \chi_{ikl}^2}{n_{j=I}} \quad (4)$$

where the summation is over all sites, seasons, and levels of dryness where model I was used as denoted by subscripts i , k , and l , respectively; $n_{j=I}$ is the total site-months simulated with model I ; and $\chi_{j=I}^2$ is aggregated χ^2 for model I . We did not evaluate model performance for any data group with $n < 3$. In sum, Taylor displays and skill examined models' ability to mimic the monthly trajectory of observed NEE , the calculation of $NMAE$ quantified bias in units of mean observed NEE , and χ^2 values quantified how well model-data mismatch scales with flux uncertainty.

Observational flux uncertainty

We calculated the standard error of monthly NEE (δ_{NEE}) [Barr et al., 2009] by combining random uncertainty and uncertainty associated with the friction velocity threshold (u_*^{Th}), a value used to identify and reject spurious nighttime NEE measurements. Random uncertainty was estimated following Richardson & Hollinger [2007]: (i) generate synthetic NEE data using the gap-filling model [Barr et al., 2004; 2009] for a given site-year, (ii) introduce gaps as in the observed data with u_*^{Th} filtering, (iii) add noise, (iv) infill gaps using the gap-filling model, (v) repeat the process 1000 times for each site.

The random uncertainty component of δ_{NEE} was then the standard deviation across all 1000 realizations aggregated to months.

The u_*^{Th} uncertainty component of δ_{NEE} was also estimated using Monte Carlo methods. Here 1000 realizations of NEE were generated using 1000 draws from a distribution of u_*^{Th} . This distribution was based on binning the raw flux data with respect to climatic season, temperature, and site-year and estimating u_*^{Th} in each bin [Papale et al., 2005]. The standard deviation across all realizations gave the u_*^{Th} uncertainty component of δ_{NEE} . Both components were combined in quadrature to one standard error of monthly NEE ($= \delta_{NEE}$) [Barr et al., 2009].

Relating model skill to model structure and site history

The models evaluated here range widely in their emphasis and structure (Table 1). Some focus on biophysical calculations (SiB3, BEPS), where others emphasize biogeochemistry (DLEM), or ecosystem dynamics (ED2). However, as terrestrial biosphere models simulate carbon cycling with hydrological variables, most models contain both biophysics and biogeochemistry. This motivated characterizing model structure with definite attributes, e.g., prognostic vs. prescribed canopy phenology, number of soil pools, and type of NEE algorithm (Table 3). To resolve how such characteristics and site history impacted model skill we calculated S for all observed combinations of site, model, seasonality, and drought level and cross-referenced these with 13 site history variables and 14 model attributes (Table 3). Only 20 models were available for this exercise, MEAN and the optimized LoTEC were excluded. We used S

as it is bound by zero (no agreement) and unity (perfect agreement) in contrast to *NMAE* and χ^2 which are unbound. The Taylor skill metric (*S*) was first discretized into three classes based on terciles. These classes, representing three tiers of model-data agreement, were then related to biome, climatic season, drought level, site history, and model structure using regression tree analysis (RTA) as a supervised classification algorithm. RTA is a form of binary recursive partitioning [Breiman et al., 1984] that successively splits the data (Taylor skill classes as the response; all other attributes as predictors) into subsets (nodes) by minimizing within-subset variation. The result is a pruned tree-like topology whereby predicted values (Taylor skill metric class) are derived by a top-to-bottom traversal following the rules (branches) that govern subset membership until a predicted value is reached (terminal node). The splitting rules at each node as well as its position allow for a calculation of relative variable importance [Breiman et al., 1984] with the most important variable given a score of 100. Variables of high importance were further analyzed using conditional means, i.e., comparing mean values for each predictor value, with statistical differences determined using Bonferroni corrections for multiple comparisons [Hochberg & Tamhane, 1987].

Results

Model-data agreement relative to climatic season, dryness, and biome

Overall agreement across $n = 31025$ months was better in forested than non-forested biomes; both *NMAE* (Table 4) and χ^2 values (Table 5) were closer to zero and unity respectively. At the biome level, model skill was loosely ranked in five tiers: evergreen needleleaf forests in the temperate zone, mixed forests > deciduous broadleaf forests,

evergreen needleleaf forests in the boreal zone > grasslands, woody savannahs > croplands, shrublands, wetlands > tundra. These rankings were robust across models used in the majority of biomes, although some divergence was apparent for croplands and shrublands (Figure 1). Relative to seasonality and drought level models were most consistent with observations during periods of peak biological activity (climatic summer) and under dry conditions (Figure 2). However, across the three levels of dryness, changes in model-data agreement were negligible for *NMAE* (~4% change, Table 4) but more pronounced for χ^2 (from 8.10 to 12.72, Table 5). Averaged over just the warm season (excluding climatic winter) dry conditions were coincident with worse model-data agreement, e.g., *NMAE* was -0.99, -0.91, and -0.84 for dry, normal, and wet, respectively. In biomes with a clear seasonal cycle in leaf area index (LAI) a loss of model skill occurred during climatic spring and fall (Table 4 & 5), especially for *NMAE*.

Skill metrics by model

Regardless of metric, model skill was highly variable. Of the three model skill metrics, *NMAE* was related to both Taylor skill and χ^2 ($\rho = -0.65$; $p < 0.0001$). Jointly, high Taylor skill co-occurred with *NMAE* and χ^2 values closer to zero and unity respectively (Figure 3). Across models *NMAE* ranged from -0.42 of the overall mean observed flux to -2.18 for LoTEC and DNDC, respectively. Values of χ^2 varied from 2.17 to 29.87 for LoTEC and CN-CLASS, respectively. Alternatively, the degree of model-data mismatch (the distance between observations and simulations) was at least 2.17 times the observational flux uncertainty. Similarly, Taylor skill showed a high degree of scatter (Figure 4),

although two crop only models (SiBCrop and AgroIBIS), LoTEC, and ISOLSM were more conservative and showed a general high degree of consistency with observations.

Among crop models, SiBCrop and AgroIBIS performed well, especially in climatic spring and during wet conditions. In contrast, the crop only DNDC model exhibited poor model-data agreement with $\chi^2 > 15$ in climatic spring and summer as well as across all drought levels. Although four crop only simulators were analyzed, the best agreement in croplands ($NMAE$ and χ^2 closer to zero and unity respectively) was achieved by SiB3 and Ecosys, models used in multiple biomes. Based on all three skill metrics the LoTEC model ($NMAE = -0.42$, $\chi^2 = 2.17$, $S = 0.95$) was most consistent with observations across all sites, dryness levels, and climatic seasons. This platform was optimized using a data assimilation technique, unique among model runs evaluated here, and was applied at 10 sites. In addition, the mean model ensemble (MEAN) performed well ($NMAE = -0.74$, $\chi^2 = 3.35$, $S = 0.80$). For individual models ($n = 12$) used at a wider range of sites (at least 24 sites), model consistency with observations was highest for Ecosys ($NMAE = -0.69$, $\chi^2 = 7.71$, $S = 0.94$) and lowest for CN-CLASS ($NMAE = -1.50$, $\chi^2 = 29.87$, $S = 0.48$).

Site-level model-data agreement also showed a high degree of variability (Figure 4). At three croplands sites (US-Ne1, US-Ne2, and US-Ne3) Taylor skill ranged from zero to unity. Both $NMAE$ and χ^2 exhibited similar scatter by site (not shown). Even for the best predicted site (US-Syv), S ranged from 0.19 to 0.95. Only two forested sites (CA-Qfo and CA-TP4) were predicted well ($S > 0.5$) by all models; whereas only one tundra site (US-Atq) was consistently poorly predicted ($S < 0.5$). Despite the wide range in model

performance, model skill ($NMAE$, χ^2 , and S) was not correlated with the number of sites ($p > 0.5$) or biomes ($p > 0.3$) simulated, i.e., using a more general rather than a specialized model did not result in a loss in model performance. Also, model-data agreement was not better at sites with longer data records ($p > 0.1$).

The steady-state protocol had negligible effect on model skill. Long-term simulated NEE by site and model varied from -2904 to 2227 g C m⁻² yr⁻¹ with 90% of all values between -600 and 100 g C m⁻² yr⁻¹. The extreme values were primarily croplands simulated outside of crop only models. Overall, only 5 models achieved steady-state (simulated $NEE \rightarrow 0$) over the full simulation: Biome-BGC, LPJ, SiBCASA, SiB3, and TECO. Similar to simulated values, observed annual integrals at the 44 sites examined did not show steady state (Table 1) and varied from -718 to 571 g C m⁻² yr⁻¹. Nonetheless, model skill was not related to how close model spinup and initial conditions approximated steady state or how close a given site was to an observed NEE of zero. All three skill metrics were uncorrelated with long-term observed or simulated average annual NEE ($p > 0.05$). However, two models did show significant relationships: For Ecosys, χ^2 increased (decrease in model skill) and S decreased as observed or simulated NEE approached zero; a system closer to steady state was coincident with less model-data agreement. BEPS was similar, showing lower S and more negative $NMAE$ (decrease in model skill) for sites closer to steady state.

Model and site-specific consistency with observations using Taylor diagrams

Average model performance (both across-site and across-model) was evaluated using Taylor diagrams based on all simulated and observed monthly *NEE* data. Better model performance was indicated by proximity to the benchmark, representing the observed state. The benchmark was normalized by observed standard deviation such that the distance of σ and *RMSE* from the benchmark was in observed σ units. Similar to model skill metrics, forested sites were better predicted than non-forested ones. The MEAN model showed $\rho \geq 0.2$, apart from CA-SJ2 and US-Atq, but generally (33 of 44 sites) underpredicted the variability associated with monthly *NEE* at forested (Figure 5) and non-forested (Figure 6) sites. Similarly, 40 of 44 sites were predicted with $RMSE < \sigma$. Also 8 (6 forested and two croplands sites: CA-Obs, CA-Qfo, CA-TP4, US-Ho1, US-IB1, US-MMS, US-Ne3, US-UMB) of the 44 sites were predicted with $\rho \geq 0.95$ and $RMSE < 1$. The worst predicted site was CA-SJ2 with $\rho = -0.67$, $\sigma = 4.3$, and $RMSE = 5.1$.

Overall model performance, aggregated across sites, was similar (Figure 7). Most models underpredicted variability and showed $RMSE < \sigma$. Of all 22 models only DNDC exhibited $\rho < 0.2$. Based on proximity to the benchmark, i.e., a high *S* value (Figure 3), the best models were: EPIC (crop only model used on one site), ISOLSM (used on 9 sites), LoTEC (data assimilation model), SiBcrop and AgroIBIS (crop only models), EDCM (used on 10 sites), Ecosys and SiBCASA (models used on most sites, 39 and 35 respectively), and MEAN (mean model ensemble for all 44 sites). All of these “best” models had $\rho > 0.75$, $RMSE < 0.75$ and slightly underpredicted variability; except the

crop only models and Ecosys where variability was overpredicted. Models whose average behavior was furthest away from the benchmark were DNDC followed by BEPS.

Links between model skill, model structure, and site history

Biome classification was the most important factor in the distribution of model skill (Figure 8) sampled across all combinations of site, model, climatic season, and drought ($n = 3132$ groups). Climatic season and stand age, the highest scored site-specific attribute, followed biome as lead determinants of model skill. Of the 12 evaluated site disturbances (Table 3) only grazing, which occurred on croplands, grasslands, and woody savannahs, achieved an importance score of at least 25. Apart from drought and grazing activity, the remaining determinants were model-specific: the number of soil layers, vegetation pools, canopy phenology, and soil pools. Two carbon flux calculations also had a variable score > 25 , with *NEE* being the highest.

Comparing mean S for these relatively important model attributes (Figure 9) revealed three instances where model structure showed a statistically significant relationship with model skill: prescribed canopy phenology, a daily time step, and calculating *NEE* as the difference between *GPP* and ecosystem respiration. Models using canopy characteristics and phenology prescribed from remotely sensed products achieved higher skill ($S = 0.54$) than either prognostic or semi-prognostic models ($S = 0.43$; $p < 0.05$). Using a daily time step showed lower model skill ($S = 0.40$) relative to non-daily time steps ($S = 0.50$; $p < 0.05$). Finally, calculating *NEE* as the difference between *GPP* and total ecosystem respiration showed greater skill ($S = 0.50$) than other calculation methods ($S = 0.42$; $p <$

0.05). None of the other model attributes we studied showed statistically significant relationships between model structure and skill.

While not statistically significant, both vegetation pools and soil layers exhibited a weak pattern whereby the simplest and most complex models showed higher skill than models of intermediate complexity (Figure 9). Models with no soil model (zero soil layers) or no vegetation pools showed greater skill than models with the simplest soil model or smallest number of vegetation pools. As the number of soil layers or pools increased, so did model skill, indicating that a more comprehensive treatment of biological and physical processes can improve model skill. For vegetation pools, there was a limit where increased complexity beyond eight pools did not improve model-data agreement.

Despite these effects, model attributes were of secondary importance. The change in S relative to biome varied from 0.28 to 0.55; a much larger range than seen for model attributes. Similarly, the high variable importance scores for biome and climatic season, as well as the lower score for drought level, corroborated the relationships between these factors and model skill as seen with $NMAE$ and χ^2 . While the regression tree algorithm achieved an accuracy of 68.5% for predicting Taylor skill class, the site history and model characteristics considered here did not explain the underlying cause of biome and seasonal differences in model skill.

Discussion

Effect of parameter sets on model performance

Model parameter sets are a large source of variability in terms of model performance [Jung et al., 2007b]. They influence output and accuracy [Grant et al., 2005] and are more important for accurately simulating CO₂ exchange than capturing effects of interannual climatic variability [Amthor et al., 2001]. For at least some of the models studied here this can be related to the use of biome-specific parameters relative to within-biome variability [Purves & Pacala, 2008]. A corollary occurs in the context of EC observations as tower footprints can exhibit heterogeneity, particularly in soils, that is not reproduced in site-specific parameters [Amthor et al., 2001].

The importance of model parameter sets was visible in this intercomparison in two ways. First, biome had the highest variable importance score. Inasmuch as models rely on biome-specific parameter values, this finding indicates that model parameter sets are a key factor in the distribution of model skill. This extends to plant functional types due to the high degree of overlap between both. Furthermore, the variability (Figure 4) in model skill across parameter sets, i.e., across models, underscores that biomes may be too heterogeneous in time [Stoy et al., 2005; 2009] and space to be well-represented by constant parameters relative to, e.g., within-biome climate variability [Hargrove et al. 2003]. Second, the general high degree of site-specific variation in model skill (Figure 4) suggested that model parameter sets may need to be refined to capture local, site-specific realities.

Effect of model structure on model performance

In general, models with the highest model-data agreement all used prescribed canopy phenology, calculated *NEE* as the difference between *GPP* and ecosystem respiration, and did not use a daily time step. Models that exhibited all of these structural characteristics (SiBCASA, SiB3, and ISOLSM) showed high degrees of model-data agreement across all three skill metrics. Similarly, Ecosys, which used a prognostic canopy but otherwise had similar structural characteristics as SiBCASA, also performed well. Relative to model complexity, consistency with observations was highest in those models with either the simplest structure (e.g., one soil carbon pool in ISOLSM) or the most complex (e.g., SiBCASA with 13 carbon pools). Models with a prognostic canopy seem to perform better with more carbon pools and soil layers (e.g., Ecosys). No model with a prognostic canopy and a low number of carbon pools and soil layers placed in the top tercile of model skill for any skill metric, except SiBcrop and AgroIBIS for Taylor skill in croplands. Using multi-model ensembles (MEAN) or data assimilation to optimize model parameter sets (LoTEC) can compensate for differences in model structure to improve model skill.

The relationships between model structure and model skill were consistent across all biomes. As a whole, the models performed better at forested sites than non-forested sites, but the same models showed the highest consistency with observations in each biome (Ecosys and SiB3). This is true even for agriculture sites, where Ecosys and SiB3 scored as high as crop only models. This suggests that any model with requisite structural attributes can successfully simulate carbon flux in all types of ecosystems.

Links between model performance and environmental factors

Model skill was only weakly linked to drought, showing high variability across dryness level by biome and model. Only during the warm season (all climatic seasons excluding winter) did aggregate model skill decline under drought conditions. While this points to process uncertainty [Sitch et al., 2008], ecosystem response to longer-term drought can exhibit lags and positive feedbacks [Arnone et al., 2008; Granier et al., 2007; Thomas et al., 2009; Williams et al., 2009] that were not explicitly included in the drought metric used here but did influence simulation behavior through model structure, e.g., soil moisture model and soil resolution.

In spring and fall, especially for biomes with a significant deciduous component, models showed a decline in model skill (Table 4) relative to periods of peak biological activity (climatic summer) (see also Morales et al. [2005]). While this was more pronounced for *NMAE* (Table 4) than χ^2 (Table 5), phenological cues are known to influence the annual carbon balance at multiple scales [Barr et al., 2007; Delpierre et al., 2009; Keeling et al., 1996]. The loss of model skill seen in this study during spring and fall was likely linked to poor treatment of leaf initiation and senescence as well as season-specific effects of soil moisture and soil temperature on canopy photosynthesis [Hanson et al., 2004]. In this study seasonality was second only to biome in driving model skill (Figure 8). This and the lack of link between model skill and site history strongly implicate phenology as a needed refinement of terrestrial biosphere simulators.

The evergreen needleleaf forest biome diverged in performance based on whether the sites were located in the temperate or boreal zones. A similar divergence was reported using Biome-BGC, LPJ, and ORCHIDEE to simulate gross CO₂ uptake across a temperature gradient in Europe [Jung et al., 2007a]; average relative *RMSE* was higher for evergreen needleleaf forests in the boreal zone. This was linked to an overestimation of LAI at the boreal sites and relationships between resource availability and leaf area [Friedlingstein et al., 2006; Jung et al., 2007a; Sitch et al., 2008]. Additionally, recent observations in the circumboreal region, where all boreal evergreen needleleaf forested sites are located, suggest that transient effects of climate change, e.g., increased severity and intensity of natural disturbances (fire, pest outbreaks) and divergence from climate normals in temperature, have already occurred [Soja et al., 2007] and influence resource availability. We speculate the loss of model skill in boreal relative to temperate evergreen needleleaf forests was linked to insufficient characterization of cold temperature sensitivity of metabolic processes and water flow in plants as well as freeze-thaw dynamics [Schaefer et al., 2007; 2009] and that this was exacerbated by the effects of transient climate change.

Effects of site history and protocol on model evaluation

Disturbance regime and how a model treats disturbance are known to impact model performance [Ito, 2008]. In this study, stand age impacted model skill whereas site history was of marginal importance (Figure 8). However, CA-SJ2, the worst predicted site (Figure 5), was harvested in 2000 and scarified in 2002, and US-SO2, a second poorly predicted shrubland site (Figure 6), suffered catastrophic wildfire during the

analysis period. The poor model performance for recently disturbed sites followed from assumed steady state as used in some simulations and the absence of modeling logic to accommodate disturbance. However, the distribution of site history metrics was skewed; only few sites were burned, harvested, or in the early stages of recovery from disturbance when *NEE* is more nonlinear relative to established stands. Furthermore, age class was biased toward older stands; of the 17 forested sites only one was classified as a young stand. Other site characteristics were also unbalanced; all non-forested biomes occurred on five or less sites; with only one site each for shrublands and woody savannahs. While regression trees are inherently robust, additional observed and simulated fluxes in rapidly growing young forested stands, recently burned or harvested sites, and undersampled biomes are desirable to better characterize model performance.

Aspects of the NACP site synthesis protocol and analysis framework also influenced the interpretation of our results. First, this analysis focused solely on non-gap-filled data to allow the model-data intercomparison to inform model development. However, the low turbulence (friction velocity) filtering removed more data at night than during the day. Average data coverage across all sites was 82% for daytime and 39% at night, respectively (Table 2), so our analysis is skewed towards daytime conditions. Secondly, each model that used remotely sensed inputs (such as LAI) repeated an average seasonal cycle calculated from site-specific time series based on all pixels within 1 km of the tower site. This likely deflated relevant variable importance scores (Figure 8) and precluded a full comparison of prescribed vs. prognostic LAI. While only few models used such inputs (Table 1), removing the inherent bias of an invariant seasonal cycle over

multiple years may improve model performance. Incorporating disturbance information to recreate historical land use and disturbance, especially for recent site entries, could also improve model performance. Lastly, despite the model simulation protocol's emphasis on steady state, this condition was not achieved for most sites (Table 2), even when discounting observational uncertainty, or most models. None of the four crop only models achieved steady state. This followed from site history of croplands in general where active management precluded any system steady state, e.g., DNDC allowed for prescribed initial soil carbon pools. For those models (five of the 21 evaluated) that achieved steady state in initialization this resulted in an inherent bias between simulated and observed *NEE* for all sites regardless of site history. However, as biome and seasonality largely governed the distribution of model skill, this bias was too small to manifest itself in this study. Relaxing the steady state assumption [Carvalhais et al., 2008] or initializing using observed wood biomass and the quasi-steady state assumption [Schaefer et al., 2008] could improve these models' performance.

Conclusion

We used observed CO₂ exchange from 44 eddy covariance towers in North America with simulations from 21 terrestrial biosphere models and a mean model ensemble to examine model skill across gradients in dryness, seasonality, biome, site history, and model structure. Models' ability to match observed monthly net ecosystem exchange was generally poor; the mean squared distance between observations and simulations was ~10-times observational error. Overall, forested sites were better predicted than non-forested sites. Weaknesses in model performance concerned model parameter sets and

phenology, especially for biomes with a clear seasonal cycle in leaf area index. Drought was weakly linked to model skill with abnormally dry conditions during the growing season showing marginally worse model-data agreement compared to non-dry conditions. Sites with disturbances during the analysis period and undersampled biomes (grasslands, shrublands, wetlands, woody savannah, and tundra) also showed a large divergence between observations and simulations. The highest degree of model-data agreement occurred in temperate evergreen forests in all climatic seasons and during summer across all biomes. Overall skill was higher for models that estimated net ecosystem exchange as the difference between gross primary productivity and ecosystem respiration, used prescribed canopy phenology, and did not use a daily time step. The model ensemble (mean simulated value across all models) and an optimized model (parameters tuned using data assimilation) also performed well. Models with preferred structural attributes included generalist models (models used at multiple sites and biomes, e.g., SiB3, Ecosys) that exhibited high degrees of model-data agreement across all biomes, indicating that a single model can successfully simulate carbon flux in all types of ecosystems. That is, different model architectures were not needed for different types of ecosystems and model choice is recast as a function of ease of parameterization and initialization.

Acknowledgements

CRS, CAW, and KS were supported by the U.S. National Science Foundation grant ATM-0910766. We would like to thank the North American Carbon Program Site-Level Interim Synthesis team, the Modeling and Synthesis Thematic Data Center, and the Oak Ridge National Laboratory Distributed Active Archive Center for collecting, organizing,

and distributing the model output and flux observations required for this analysis. This study was in part supported by the U.S. National Aeronautics and Space Administration (NASA) grant NNX06AE65G, the U.S. National Oceanic and Atmospheric Administration (NOAA) grant NA07OAR4310115, and the U.S. National Science Foundation (NSF) grant OPP-0352957 to the University of Colorado at Boulder.

References

Amthor, JS et al. (2001) Boreal forest CO₂ exchange and evapotranspiration predicted by nine ecosystem process models: Intermodel comparisons and relationships to field measurements. *J. Geophys. Res.*, 106(D24), 33,623-33,648.

Angert A, Biraud S, Bonfils C et al. (2005) Drier summers cancel out the CO₂ uptake enhancement induced by warmer springs. *Proceedings of the National Academy of Sciences of the United States of America*, 102, 10823-10827.

Arain MA, Yaun F, Black TA (2006) Soil-plant nitrogen cycling modulated carbon exchanges in a western temperate conifer forest in Canada. *Agricultural and Forest Meteorology*, 140, 171-192.

Arnone, JA, Verburg PSJ, Johnson DW et al. (2008) Prolonged suppression of ecosystem carbon dioxide uptake after an anomalously warm year. *Nature* 455:383-386.

Baldocchi D, Falge E, Gu LH et al. (2001) FLUXNET: A new tool to study the temporal and spatial variability of ecosystem-scale carbon dioxide, water vapor, and energy flux densities. *Bulletin of the American Meteorological Society*, 82, 2415-2434.

Baldocchi, D (2008) Breathing of the terrestrial biosphere: lessons learned from a global network of carbon dioxide flux measurement systems. *Australian Journal of Botany*, 56, 1-26.

Baker IT, Prihodko L, Denning AS, Goulden M, Miller S, da Rocha HR (2008) Seasonal drought stress in the Amazon: Reconciling models and observations. *J. Geophys. Res.*, 113, G00B01, doi:10.1029/2007JG000644.

Barr AG, Black TA, Hogg EH et al. 2004. Inter-annual variability in the leaf area index of a boreal aspen-hazelnut forest in relation to net ecosystem production. *Agricultural Forest Meteorology*, 126: 237-255.

Barr AG, Black TA, Hogg EH et al. (2007) Climatic controls on the carbon and water balances of a boreal aspen forest, 1994–2003. *Global Change Biology*, 13, 561-576.

Barr A, Hollinger D, Richardson, AD (2009) CO₂ Flux Measurement Uncertainty Estimates for NACP. *Eos Trans. AGU*, 90(52), Fall Meet. Suppl., Abstract B54A-04.

Bergeron O, Margolis HA, Black TA, Coursolle C, Dunn AL, Barr AG, Wofsy SC (2007) Comparison of CO₂ fluxes over three boreal black spruce forests in Canada. *Global Change Biol.* 13, 89-107.

Bradford JB, Birdsey RA, Joyce LA, Ryan MG (2008) Tree age, disturbance history, and carbon stocks and fluxes in subalpine Rocky Mountain forests. *Global Change Biology*, 14, 2882-2897.

Breiman L, Friedman JH, Olshen RA, Stone CJ (1984) *Classification and regression trees*. Wadsworth, Belmont, CA, 358 pp.

Carvalhais N et al. (2008) Implications of the carbon cycle steady state assumption for biogeochemical modeling performance and inverse parameter retrieval. *Global Biogeochem. Cycles*, 22, GB2007, doi:10.1029/2007GB003033.

Causarano HJ, Shaw JN, Franzluebbers AJ, Reeves DW, Raper RL, Balkcom KS, Norfleet ML, Izaurralde RC (2007) Simulating field-scale soil organic carbon dynamics using EPIC. *Soil Science Society of America Journal*, 71, 1174-1185.

Chen JM, Govind A, Sonnentag O, Zhang Z, Barr A, Amiro B (2006) Leaf area index measurements at Fluxnet-Canada forest sites. *Agric. For. Meteorol.*, 140, 257-268.

Ciais P et al (2005) Europe-wide reduction in primary productivity caused by the heat and drought in 2003. *Nature* 437:529-533

Cook BD, Davis KJ, Wang W et al. (2004a) Carbon exchange and venting anomalies in an upland deciduous forest in northern Wisconsin, USA. *Agricultural and Forest Meteorology*, 126, 271-295.

Cook ER, Woodhouse CA, Eakin CM, Meko DM, Stahle DW (2004b) Long-term aridity changes in the western United States. *Science*, 306, 1015-1018.

Cook ER, Seager R, Cane MA, Stahle DW (2007) North American droughts: reconstructions, causes and consequences. *Earth-Sci. Rev.*, 81, 93-134.

Dai A, Trenberth KE, Qian T (2004) A global data set of Palmer Drought Severity Index for 1870-2002: Relationship with soil moisture and effects of surface warming. *J. Hydrometeorol.*, 5, 1117-1130.

Davis KJ, Bakwin PS, Yi CX, Berger NW, Zhao CL, Teclaw RM, Isebrands JG (2003) The annual cycles of CO₂ and H₂O exchange over a northern mixed forest as observed from a very tall tower. *Global Change Biology*, 9, 1278-1293.

Delpierre N, Soudani K, Francois C et al. (2009) Exceptional carbon uptake in European forests during the warm spring of 2007: a data–model analysis. *Glob. Change Biol.*, 15, 1455-1474.

Desai AR, Bolstad PV, Cook BD, Davis KJ, Carey EV (2005) Comparing net ecosystem exchange of carbon dioxide between an old-growth and mature forest in the upper midwest, USA. *Agricultural and Forest Meteorology*, 128, 33-55.

Flanagan LB, Wever LA, Carlson PJ (2002) Seasonal and interannual variation in carbon dioxide exchange and carbon balance in a northern temperate grassland. *Global Change Biol.*, 8, 599-615.

Friedlingstein P, Cox PM, Betts R et al. (2006) Climate-carbon cycle feedback analysis, results from the C⁴MIP model intercomparison. *Journal of Climate*, 19, 3337-3353.

Granier A, Reichstein M, Breda N et al. (2007) Evidence for soil water control on carbon and water dynamics in European forests during the extremely dry year: 2003. *Agricultural and Forest Meteorology*, 143, 123–145.

Grant RF et al. (2005) Intercomparison of techniques to model high temperature effects on CO₂ and energy exchange in temperate and boreal coniferous forests. *Ecol. Model.*, 188, 217–252.

Griffis TJ, Black TA, Morgenstern K, Barr AG, Nestic Z, Drewitt GB, Gaumont-Guay D, McCaughey JH (2003) Ecophysiological controls on the carbon balances of three southern boreal forests. *Agricultural and Forest Meteorology*, 117, 53-71.

Gu et al. (2006) Direct and indirect effects of atmospheric conditions and soil moisture on surface energy partitioning revealed by a prolonged drought at a temperate forest site. *Journal of Geophysical Research*, Vol. 111, D16102, doi:10.1029/2006JD007161.

Hanson PJ, Amthor JS, Wullschlegel SD et al. (2004) Oak forest carbon and water simulations: model intercomparisons and evaluations against independent data. *Ecological Monographs*, 74, 443-489.

Harazono Y, Mano M, Miyata A, Zulueta RC, Oechel WC (2003) Inter-annual carbon dioxide uptake of a wet sedge tundra ecosystem in the Arctic. *Tellus* 55B, 2, 215-231.

Hargrove WW, Hoffman FM, Law BE (2003) New Analysis Reveals Representativeness of AmeriFlux Network. *Earth Observing System Transactions, American Geophysical Union* 84(48):529.

Hochberg Y, Tamhane AC (1987) *Multiple Comparison Procedures*. Wiley, New York, 480 pp.

Houghton JT, Ding Y, Griggs DJ, Noguer M, van der Linden PJ, Xia D, Maskell K, Johnson CA (eds) (2001) *Climate Change 2001: The Scientific Basis: Contributions of Working Group I to the Third Assessment Report of the Intergovernmental Panel on Climate Change*. Cambridge University Press, New York, 881 pp.

Huntington TG (2006) Evidence for intensification of the global water cycle: Review and synthesis. *J. Hydrol.*, 319, 83-95.

Hwang Y, Carbone GJ (2009) Ensemble forecasts of drought indices using a conditional residual resampling technique. *Journal of Applied Meteorology and Climatology*, 48, 1289-1301.

Ichii K, Suzuki T, Kato T et al (2009) Multi-model analysis of terrestrial carbon cycles in Japan: reducing uncertainties in model outputs among different terrestrial biosphere models using flux observations, *Biogeosciences Discuss.*, 6, 8455-8502.

Irvine J, Law BE, Kurpius M, Anthoni PM, Moore D, Schwarz P (2004) Age related changes in ecosystem structure and function and the effects on carbon and water exchange in ponderosa pine. *Tree Physiology*, 24, 753-763.

Ito A (2008) The regional carbon budget of East Asia simulated with a terrestrial ecosystem model and validated using AsiaFlux data. *Agr. Forest Meteorol.*, 148, 738-747.

Jung M, Le Maire G, Zaehle S et al. (2007a) Assessing the ability of three land ecosystem models to simulate gross carbon uptake of forests from boreal to Mediterranean climate in Europe. *Biogeosciences*, 4, 647-656.

Jung M, Vetter M, Herold M et al. (2007b) Uncertainties of modeling gross primary productivity over Europe: A systematic study on the effects of using different drivers and terrestrial biosphere models. *Global Biogeochem. Cycles*, 21, GB4021, doi:10.1029/2006GB002915.

Keeling CD, Chin JF, Whorf TP (1996) Increased activity of northern vegetation inferred from atmospheric CO₂ measurements. *Nature*, 382, 146-149.

Krinner G, Viovy N, de Noblet-Ducoudré N, Ogée J, Polcher J, Friedlingstein P, Ciais P, Sitch S, Prentice IC (2005) A dynamic global vegetation model for studies of the coupled atmosphere-biosphere system, *Global Biogeochem. Cycles*, 19, GB1015, doi:10.1029/2003GB002199.

Kucharik CJ, Twine TE (2007) Residue, respiration, and residuals: Evaluation of a dynamic agroecosystem model using eddy flux measurements and biometric data. *Agricultural and Forest Meteorology*, 146, 134-158.

Lafleur PM, Roulet NT, Bubier JL, Moore TR, Frolking S (2003) Interannual variability in the peatland-atmosphere carbon dioxide exchange at an ombrotrophic bog. *Global Biogeochemical Cycles*, 17, 5.1-5.13, 1036, doi:10.1029/2002GB001983.

Larcher W (1995) *Physiological Plant Ecology*. Springer Verlag, Berlin, 506 pp.

Li et al. (2009) Modeling impacts of alternative farming management practices on greenhouse gas emissions from a winter wheat-maize rotation system in China. *Agriculture, Ecosystems and Environment*, 135, 24-33.

Liu J, Chen JM, Cihlar J, Chen W (1999) Net primary productivity distribution in the BOREAS region from a process model using satellite and surface data. *J. Geophys. Res.*, 104, 27735–27754.

Liu S, Bliss N, Sundquist E, Huntington TG (2003) Modeling carbon dynamics in vegetation and soil under the impact of soil erosion and deposition. *Global Biogeochemical Cycles*, 17, doi:10.1029/2002GB002010.

Lokupitiya E, Denning S, Paustian K, Baker I, Schaefer K, Verma S, Meyers T.

Bernacchi CJ, Suyker A, Fischer M (2009) Incorporation of crop phenology in Simple Biosphere Model (SiBcrop) to improve land-atmosphere carbon exchanges from croplands. *Biogeosciences*, 6, 969-986.

Loveland TR, Reed BC, Brown JF, Ohlen DO, Zhu J, Yang L, Merchant JW (2001) Development of a global land cover characteristics database and IGBP DISCover from 1-km AVHRR data. *International Journal of Remote Sensing*, 21, 1,303-1,330.

Luo H, Oechel WC, Hastings SJ, Zulueta R, Qian Y, Kwon H (2007) Mature semiarid chaparral ecosystems can be a significant sink for atmospheric carbon dioxide. *Global Change Biology*, 13, 386-396.

Ma S, Baldocchi DD, Xu L, Hehn T (2007) Inter-annual variability in carbon dioxide exchange of an oak/grass savanna and open grassland in California. *Agricultural and Forest Meteorology*, 147, 157-171

McCaughey JH, Pejam MR, Arain MA, Cameron DA (2006) Carbon dioxide and energy fluxes from a boreal mixedwood forest ecosystem in Ontario, Canada. *Agric. For. Meteorol.* 140, 79-96.

McKee TB, Doeskin NJ, Kleist J (1993) The relationship of drought frequency and duration to time scales. *Proc. 8th Conf. on Applied Climatology*, January 17-22, 1993, American Meteorological Society, Boston, Massachusetts, 179-184.

Medlyn BE, Robinson AP, Clement R, McMurtrie RE (2005) On the validation of models of forest CO₂ exchange using eddy covariance data: some perils and pitfalls. *Tree Physiology*, 25, 839-857.

Medvigy D, Wofsy SC, Munger JW, Hollinger DY, Moorcroft PR (2009) Mechanistic scaling of ecosystem function and dynamics in space and time: Ecosystem Demography model version 2. *J. Geophys. Res.*, 114, G01002, doi:10.1029/2008JG000812.

Meehl GA, Tebaldi C (2004) More intense, more frequent, and longer lasting Heat waves in the 21st Century. *Science*, 305, 994-997.

Mishra VR, Desai AK (2006) Drought forecasting using feed-forward recursive neural network. *Ecol. Modell.*, 198, 127-138.

Moffat A, Papale D, Reichstein M et al. (2007) Comprehensive comparison of gap-filling techniques for eddy covariance net carbon fluxes. *Agricultural and Forest Meteorology*, 147, 209-232.

Morales P, Sykes MT, Prentice IC et al (2005) Comparing and evaluating process-based ecosystem model predictions of carbon and water fluxes in major European forest biomes. *Global Change Biol.*, 11, 2211-2233.

Oberbauer SF, Tweedie CE, Welker JM et al. (2007) Tundra CO₂ fluxes in response to experimental warming across latitudinal and moisture gradients. *Ecol. Monogr.*, 72, 221-38.

Papale D, Reichstein M, Aubinet M et al. (2006) Towards a standardized processing of net ecosystem exchange measured with eddy covariance technique: algorithms and uncertainty estimation. *Biogeosciences*, 3, 571-583.

Peel MC, Finlayson BL, McMahon TA (2007) Updated world map of the Köppen–Geiger climate classification. *Hydrol. Earth Syst. Sci.*, 11, 1633-1644.

Peichl M, Arain MA (2007) Allometry and partitioning of above- and below-ground tree biomass in an age-sequence of white pine forests. *Forest Ecology and Management*, 253, 68-80.

Piao SL, Friedlingstein P, Peylin P, Reichstein M, Luysaert S, Margolis H, Fang JY, Barr AG, Chen A, Grelle A, Hollinger DY, Laurila T, Lindroth A, Richardson A, Vesala T (2008) Net carbon dioxide losses of northern ecosystems in response to autumn warming. *Nature* 451, 49-53.

Post WM, Izaurralde RC, Jastrow JD et al. (2004) Enhancement of carbon sequestration in U. S. soils. *Bioscience*, 54, 10, 895-908.

Purves DW, Pacala S (2008) Predictive models of forest dynamics. *Science*, 320, 1452-1453.

Ricciuto DM, King AW, Gu L, Post WM (2008) Estimates of terrestrial carbon cycle model parameters by assimilation of FLUXNET data: Do parameter variations cause bias in regional flux estimates? *Eos Trans. AGU*, 89(53), Fall Meet. Suppl., Abstract B54A-03.

Ricciuto DM, Thornton PE, Schaefer K, Cook RB, Davis KJ (2009) How uncertainty in gap-filled meteorological input forcing at eddy covariance sites impacts modeled carbon and energy flux. *Eos Trans. AGU*, 90(52), Fall Meet. Suppl., Abstract B54A-03.

Richardson AD, Hollinger DY (2007) A method to estimate the additional uncertainty in gap-filled NEE resulting from long gaps in the CO₂ flux record. *Agricultural and Forest Meteorology*, 147, 199-208.

Richardson AD, Hollinger DY, Dail DB, Lee JT, Munger W, O'Keefe J (2009) Influence of spring phenology on seasonal and annual carbon balance in two contrasting New England forests. *Tree Physiology*, 29, 321-331.

Riley WJ, Still CJ, Torn MS, Berry JA (2002) A mechanistic model of H₂O and C₁₈O fluxes between ecosystems and the atmosphere: Model description and sensitivity analyses, *Global Biogeochem. Cycles*, 16, 1095, doi:10.1029/2002GB001878.

Schaefer K, Zhang T, Tans P, Stöckli R (2007) Temperature anomaly reemergence in seasonally frozen soils. *J. Geophys. Res.*, 112, D20102, doi:10.1029/2007JD008630.

Schaefer K, Collatz GJ, Tans P et al. (2008) Combined Simple Biosphere/Carnegie-Ames-Stanford Approach terrestrial carbon cycle model, *J. Geophys. Res.*, 113, G03034, doi:10.1029/2007JG000603.

Schaefer K, Zhang T, Slater AG, Lu L, Etringer A, Baker I (2009) Improving simulated soil temperatures and soil freeze/thaw at high-latitude regions in the Simple Biosphere/Carnegie-Ames-Stanford Approach model. *J. Geophys. Res.*, 114, F02021, doi:10.1029/2008JF001125.

Schmid HP, Grimmer CSB, Cropley F, Offerle B, Su HB (2000) Measurements of CO₂ and energy fluxes over a mixed hardwood forest in the mid-western United States. *Agric. For. Meteorol.* 103, 357-374.

Schmid HP, Su HB, Vogel CS, Curtis PS (2003) Ecosystem-atmosphere exchange of carbon dioxide over a mixed hardwood forest in northern lower Michigan. *J. Geophys. Res.* Vol. 108, No. D14, 4417, doi:10.1029/2002JD003011.

Schwalm CR, Black TA, Amiro BD et al. (2006) Photosynthetic light use efficiency of three biomes across an east–west continental-scale transect in Canada. *Agric. Forest Meteorol.*, 140, 269-286.

Schwalm CR, Black TA, Morgenstern K, Humphreys ER (2007) A method for deriving net primary productivity and component respiratory fluxes from tower-based eddy covariance data: a case study using a 17-year data record from a Douglas-fir chronosequence. *Global Change Biol.*, 13, 370-385.

Schwalm CR, Williams CA, Schaefer KS et al. (2009) Assimilation exceeds respiration sensitivity to drought: A FLUXNET synthesis. *Global Change Biology*, 16, 657-670.

Seager R (2007a) The turn of the century drought across North America: global context, dynamics and past analogues. *Journal of Climate*, 20, 5527-5552.

Seager R, Ting MF, Held IM (2007b) Model projections of an imminent transition to a more arid climate in Southwestern North America, *Science*, 316, 1181-1184.

Seager R, Tzanova A, Nakamura J (2009) Drought in the Southeastern United States: Causes, variability over the last millennium and the potential for future hydroclimate change, *Journal of Climate*, 22, 5021-5045.

Sheffield J, Wood EF (2008) Projected changes in drought occurrence under future global warming from multi-model, multi-scenario, IPCC AR4 simulations. *Climate Dynamics*, 13, 79-105.

Siqueira MB, Katul GG, Sampson DA, Stoy PC, Juang J-Y, McCarthy HR, Oren R (2006) Multiscale model intercomparisons of CO₂ and H₂O exchange rates in a maturing southeastern US pine forest. *Global Change Biology*, 12, 1189-1207.

Sims PL, Bradford JA (2001) Carbon dioxide fluxes in a southern plains prairie. *Agricultural and Forest Meteorology*, 109, 117-134.

Sitch S, Smith B, Prentice IC, Arneth A, Bondeau A, Cramer W, Kaplan JO, Levis S, Lucht W, Sykes MT, Thonicke K, Venevsky S (2003) Evaluation of ecosystem dynamics, plant geography and terrestrial carbon cycling in the LPJ dynamic global vegetation model. *Global Change Biology*, 9, 161-185.

Sitch S, Huntingford C, Gedney N et al. (2008) Evaluation of the terrestrial carbon cycle future plant geography and climate-carbon cycle feedbacks using five Dynamic Global Vegetation Models (DGVMs). *Glob. Change Biol*, 14, 2015-2039.

Soja AJ, Tchebakova NM, French NHF, Flannigan MD, Shugart HH, Stocks BJ, Sukhinin AI, Varfenova EI, Chapin FS, Stackhouse PW (2007). Climate induced boreal forest change: Predictions versus current observations. *Global and Planetary Change*, 56, 274-296.

Stoy P, Katul G, Siqueira M, Juang J, McCarthy H, Kim H, Oishi A, Oren R (2005) Variability in net ecosystem exchange from hourly to inter-annual time scales at adjacent pine and hardwood forests: a wavelet analysis. *Tree Physiology*, 25, 887-902.

Stoy P, Richardson A, Baldocchi D, Katul G et al. (2009) Biosphere-atmosphere exchange of CO₂ in relation to climate: a cross-biome analysis across multiple time scales. *Biogeosciences*, 6, 2297-2312.

Sulman BN, Desai AR, Cook BD, Saliendra N, Mackay DS (2009) Contrasting carbon dioxide fluxes between a drying shrub wetland in northern Wisconsin, USA, and nearby forests. *Biogeosciences*, 6, 1115-1126.

Suyker AE, Verma SB, Burba GG (2003) Interannual variability in net CO₂ exchange of a native tallgrass prairie. *Global Change Biology*, 9, 255-265.

Syed HK, Flanagan LB, Carlson PJ, Glenn AJ, Van Gaalen KE (2006) Environmental control of net ecosystem CO₂ exchange in a treed, moderately rich fen in northern Alberta. *Agric. For. Meteorol.* 140, 97-114.

Taylor JR (1996) An introduction to error analysis. University Science Books, Mill Valley, California, USA, 327 pp.

Taylor KE (2001) Summarizing multiple aspects of model performance in a single diagram. *J. Geophys. Res.*, 106, 7183-7192.

Thomas CK, Law BE, Irvine J, Martin JG, Pettijohn JC, Davis KJ (2009) Seasonal hydrology explains interannual and seasonal variation in carbon and water exchange in a semi-arid mature ponderosa pine forest in Central Oregon. *J. Geophys. Res.*, 114, G04006, doi:10.1029/2009JG001010.

Thornton, PE, Running SW, White MA (1997) Generating surfaces of daily meteorological variables over large regions of complex terrain. *J. of Hydrology*, 3-4, 214-251

Thornton PE, Running SW, Hunt ER (2005) Biome-BGC: Terrestrial Ecosystem Process Model, Version 4.1.1. Model product. Available on-line [<http://daac.ornl.gov>] from Oak Ridge National Laboratory Distributed Active Archive Center, Oak Ridge, Tennessee, U.S.A. doi:10.3334/ORNLDAAAC/805.

Tian, HQ, Chen G, Liu M, Zhang C, Sun G, Lu C, Xu X, Ren W, Pan P, Chappelka A (2010) Model estimates of ecosystem net primary productivity, evapotranspiration, and water use efficiency in the Southern United States during 1895-2007. *Forest Ecology and Management*, 259, 1311-1327.

Trenberth KE, Jones PD, Ambenje P et al. (2007) Observations: Surface and atmospheric climate change. In: *Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change* (eds Solomon SD et al.). Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA.

Urbanski S, Barford C, Wofsy S et al. (2007), Factors controlling CO₂ exchange on timescales from hourly to decadal at Harvard Forest. *J. Geophys. Res.*, 112, G02020, doi:10.1029/2006JG000293.

Verma SB, Dobermann A, Cassman KG et al. (2005) Annual carbon dioxide exchange in irrigated and rainfed maize-based agroecosystems. *Agricultural and Forestry Meteorology*, 131, 77-96.

Vickers D, Thomas C, Law BE (2009) Random and systematic CO₂ flux sampling errors for tower measurements over forests in the convective boundary layer. *Agricultural and Forest Meteorology*, 149, 73-83.

Weng E, Luo Y (2008) Soil hydrological properties regulate grassland ecosystem responses to multifactor global change: A modeling analysis. *J. Geophys. Res.*, 113, G03003, doi:10.1029/2007JG000539.

Williams CA, Hanan NP, Scholes RJ, Kutsch WL (2009) Complexity in water and carbon dioxide fluxes following rain pulses in an African savanna. *Oecologia*, 161, 469-480.

Williamson TB, Price DT, Beverly JL, Bothwell PM, Frenkel B, Park J, Patriquin MN (2008) Assessing potential biophysical and socioeconomic impacts of climate change on forest-based communities: a methodological case study. *Nat. Resour. Can., Can. For. Serv., North. For. Cent., Edmonton, AB. Inf. Rep. NOR-X-415E*.

Zha T, Barr AG, Black TA et al. (2009) Carbon sequestration in boreal jack pine stands following harvesting. *Global Change Biology*, 15, 1475-1487.

Zhan XW, Xue YK, Collatz GJ (2003) An analytical approach for estimating CO₂ and heat fluxes over the Amazonian region. *Ecol. Model.* 162, 97-117.

Zhou XL, Peng CH, Dang QL, Sun JF, Wu HB, Hua D (2008) Simulating carbon exchange in Canadian Boreal forests I: model structure, validation, and sensitivity analysis. *Ecological Modelling*, 219, 287-299.

Tables

1 Table 1. Summary of model characteristics.

Model Attribute	Model									
	AgroIBIS	BEPS	Biome-BGC	Can-IBIS	CN-CLASS	DLEM	DNDC	Ecosys	ED2	EDCM
Temporal Resolution	Half-hourly	Daily	Daily	Half-hourly	Half-hourly	Daily	Daily	Hourly	Half-hourly	Monthly
Vegetation Pools	4	4	7	3	4	6	3	9	9	8
Soil Pools	7	9	4	7	3	3	9	9	4	5
Soil Layers	11	3	1	7	3	2	10	15	9	10
Canopy Phenology	Prognostic	Semi-Prognostic	Prognostic	Prognostic	Prognostic	Semi-Prognostic	Prognostic	Prognostic	Prognostic	Prognostic
Nitrogen Cycle	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Gross Primary Productivity (<i>GPP</i>)	Enzyme Kinetic Model	Enzyme Kinetic Model	Stomatal Conductance Model	Enzyme Kinetic Model	Enzyme Kinetic Model	Stomatal Conductance Model	Light Use Efficiency Model	Enzyme Kinetic Model	Enzyme Kinetic Model	Light Use Efficiency Model
Heterotrophic Respiration (<i>HR</i>)	First or Greater Order Model	Air Temperature Soil Temperature Precipitation Soil Moisture Evaporation Soil Carbon Soil Nitrogen	Soil Temperature Soil Moisture Soil Carbon	First or Greater Order Model	First or Greater Order Model	Decay Methane Air Temperature Soil Temperature Litter and Soil Carbon Soil Nitrogen Soil Moisture	Decay Methane Soil Temperature Precipitation Soil Moisture Soil Carbon Vegetation Carbon Soil Nitrogen	Decay Methane CO ₂ Diffusion Dissolved Carbon Soil Temperature Soil Moisture Surface Incident Shortwave Radiation Surface Incident Longwave Radiation Soil Carbon Vegetation Carbon Soil Nitrogen Leaf Nitrogen	Soil Temperature Soil Moisture Soil Carbon Soil Nitrogen	Soil Temperature Soil Moisture Carbon Loss Vegetation Carbon Soil Nitrogen

Model Attribute	AgroIBIS	BEPS	Biome-BGC	Can-IBIS	CN-CLASS	DLEM	DNDC	Ecosys	ED2	EDCM
Autotrophic Respiration (<i>AR</i>)	Air Temperature Soil Temperature Precipitation Soil Moisture Surface Incident Shortwave Radiation Surface Incident Longwave Radiation Vegetation Carbon	Air Temperature <i>GPP</i>	Air Temperature Vegetation Carbon Leaf Nitrogen	Air Temperature Soil Temperature Precipitation Soil Moisture Surface Incident Shortwave Radiation Surface Incident Longwave Radiation Vegetation Carbon	Fraction of Instantaneous <i>GPP</i>	Air Temperature Vegetation Carbon Leaf Nitrogen <i>GPP</i>	Soil Temperature	Air Temperature Soil Temperature Vegetation Carbon Leaf Nitrogen	Air Temperature Soil Temperature Vegetation Carbon Leaf Nitrogen <i>GPP</i>	Proportional to Growth
Ecosystem Respiration (<i>R</i>)	<i>AR + HR</i>	<i>AR + HR</i>	Air Temperature Soil Temperature Soil Moisture Soil Carbon Vegetation Carbon LAI	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>
Net Primary Production (<i>NPP</i>)	<i>GPP - AR</i>	<i>GPP - AR</i>	Surface Incident Shortwave Radiation Vapor Pressure Deficit CO ₂ Vegetation Carbon Leaf Nitrogen LAI	<i>GPP - AR</i>	Fraction of Instantaneous <i>GPP</i>	<i>GPP - AR</i>	Air Temperature Precipitation Soil Moisture Potential Evaporation Vegetation Carbon Soil Nitrogen Leaf Nitrogen fPAR	<i>GPP - AR</i>	<i>GPP - AR</i>	Air Temperature Precipitation on Soil Carbon Soil Nitrogen Soil Moisture Vegetation Carbon Leaf Nitrogen LAI

Model Attribute	AgroIBIS	BEPS	Biome-BGC	Can-IBIS	CN-CLASS	DLEM	DNDC	Ecosys	ED2	EDCM
Net Ecosystem Exchange (<i>NEE</i>)	<i>NPP - HR</i>	<i>NPP - HR</i>	Soil Temperature Soil Moisture Surface Incident Shortwave Radiation Vapor Pressure Deficit	<i>NPP - HR</i>	<i>GPP - R</i>	<i>NPP - HR</i>	<i>NPP - HR</i>	<i>GPP - R</i>	<i>NPP - HR</i>	<i>NPP - HR</i>
Biomes Simulated	Croplands	6	8	10	9	10	Croplands	10	6	6
Sites Simulated	5	10	36	27	31	33	5	39	25	10
Months Simulated	192	945	2001	1978	2082	2246	192	2450	1684	658
Source	Kucharik & Twine [2007]	Liu et al. [1999]	Thornton et al. [2005]	Williamson et al. [2008]	Arain et al. [2006]	Tian et al. [2010]	Li et al. [2009]	Grant et al. [2005]	Medvigy et al. [2009]	Liu et al. [2003]

1
2 Table 1 continued.

Model Attribute	Model										
	EPIC	ISOLSM	LoTEC	LPJ	ORCHIDEE	SiB3	SiBCASA	SiBcrop	SSiB2	TECO	Triplex-FLUX
Temporal Resolution	Daily	Half-hourly	Half-hourly	Daily	Half-hourly	Half-hourly	10 minutes	Half-hourly	Half-hourly	Hourly	Half-hourly
Vegetation Pools	3	0	4	3	8	0	8	4	0	3	0
Soil Pools	0	1	5	2	8	0	5	1	0	5	0
Soil Layers	15	0	14	2	0	10	15	10	3	10	0
Canopy Phenology	Prognostic	Prescribed	Prognostic	Prognostic	Prognostic	Prescribed	Prescribed	Prognostic	Prescribed	Prognostic	Prescribed
Nitrogen Cycle	Yes	No	No	No	No	Yes	No	Yes	No	No	No
Gross Primary Productivity (<i>GPP</i>)	Nil	Stomatal Conductance Model	Enzyme Kinetic Model	Stomatal Conductance Model	Enzyme Kinetic Model	Enzyme Kinetic Model	Enzyme Kinetic Model	Enzyme Kinetic Model	Stomatal Conductance Model	Stomatal Conductance Model	Stomatal Conductance Model

Model Attribute	EPIC	ISOLSM	LoTEC	LPJ	ORCHIDEE	SiB3	SiBCASA	SiBcrop	SSiB2	TECO	Triplex-FLUX
Heterotrophic Respiration (<i>HR</i>)	CO ₂ Diffusion Dissolved Carbon Loss Air Temperature Soil Temperature Precipitation Soil Moisture	First or Greater Order Model	Soil Temperature Soil Moisture Soil Carbon	Soil Temperature Soil Moisture Soil Carbon	Soil Temperature Soil Moisture Soil Carbon	Zero-order Model	Soil Temperature Soil Moisture Soil Carbon	Soil Temperature Soil Carbon	Zero-order Model	First or Greater Order Model	First or Greater Order Model
Autotrophic Respiration (<i>AR</i>)	Nil	Fraction of Instantaneous <i>GPP</i>	Air Temperature Soil Temperature Soil Moisture Vegetation Carbon <i>GPP</i>	Air Temperature Soil Moisture Vegetation Carbon	Air Temperature Vegetation Carbon	Fraction of Instantaneous <i>GPP</i>	Air Temperature Soil Moisture Vegetation Carbon	Air Temperature Vegetation Carbon <i>GPP</i>	Air Temperature Soil Moisture Surface Incident Shortwave Radiation Relative Humidity LAI fPAR CO ₂	Air Temperature Vegetation Carbon	Fraction of Annual <i>GPP</i>
Ecosystem Respiration (<i>R</i>)	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>	<i>AR + HR</i>	Forced Annual Balance	<i>AR + HR</i>	Forced Annual Balance	Forced Annual Balance	<i>AR + HR</i>	<i>AR + HR</i>

Model Attribute	EPIC	ISOLSM	LoTEC	LPJ	ORCHIDEE	SiB3	SiBCASA	SiBcrop	SSiB2	TECO	Triplex- FLUX
Net Primary Production (<i>NPP</i>)	Light Use Efficiency Model	Nil	<i>GPP - AR</i>	<i>GPP - AR</i>	<i>GPP - AR</i>	<i>GPP - AR</i>	Air Temperature Soil Moisture CO ₂ Relative Humidity	<i>GPP - AR</i>	<i>GPP - AR</i>	<i>GPP - AR</i>	Fraction of Instantane ous GPP
Net Ecosystem Exchange (<i>NEE</i>)	<i>NPP - HR</i>	<i>GPP - R</i>	<i>NPP - HR</i>	<i>NPP - HR</i>	<i>GPP - R</i>	<i>GPP - R</i>	<i>GPP - R</i>	<i>GPP - R</i>	<i>GPP - R</i>	<i>GPP - R</i>	<i>GPP - R</i>
Biomes Simulated	Croplands	5	6	9	10	10	10	Croplands	10	10	3
Sites Simulated	US-Ne3	9	10	29	35	31	35	5	44	35	7
Months Simulated	48	909	825	2126	2332	2258	2402	192	2800	2414	291
Source	Causarano et al. [2007]	Riley et al. [2002]	Hanson et al. [2004]	Sitch et al. [2003]	Krinner et al. [2005]	Baker et al. [2008]	Schaefer et al. [2009]	Lokupitiya et al. [2009]	Zhan et al. [2003]	Weng & Luo [2008]	Zhou et al. [2008]

1 Table 2. Summary of site characteristics.

Site ID	Name	Priority	Country	Latitude	Longitude	Elevation (m a.s.l.)	IGBP Class	Köppen–Geiger Climate Classification
CA-Ca1	British Columbia - Campbell River - Mature Forest Site	1	Canada	49.87	-125.33	300	ENF	Maritime temperate
CA-Ca2	British Columbia - Campbell River - Clearcut Site	2	Canada	49.87	-125.29	180	ENF	Maritime temperate
CA-Ca3	British Columbia - Campbell River - Young Plantation Site	2	Canada	49.53	-124.90	165	ENF	Maritime temperate
CA-Gro	Ontario - Groundhog River - Mature Boreal Mixed Wood	1	Canada	48.22	-82.16	300	MF	Warm summer continental
CA-Let	Lethbridge	1	Canada	49.71	-112.94	960	GRA	Warm summer continental
CA-Mer	Eastern Peatland - Mer Bleue	1	Canada	45.41	-75.52	70	WET	Warm summer continental
CA-Oas	Sask. - SSA Old Aspen	1	Canada	53.63	-106.20	530	DBF	Continental subarctic
CA-Obs	Sask. - SSA Old Black Spruce	1	Canada	53.99	-105.12	629	ENF	Continental subarctic
CA-Ojp	Sask. - SSA Old Jack Pine	1	Canada	53.92	-104.69	579	ENF	Continental subarctic
CA-Qfo	Quebec Mature Boreal Forest Site	1	Canada	49.69	-74.34	382	ENF	Continental subarctic
CA-SJ1	Sask. - 1994 Harvested Jack Pine	2	Canada	53.91	-104.66	580	ENF	Continental subarctic
CA-SJ2	Sask. - 2002 Harvested Jack Pine	2	Canada	53.94	-104.65	518	ENF	Continental subarctic
CA-SJ3	Sask. - SSA 1975 Harvested Young Jack Pine	2	Canada	53.88	-104.64	511	ENF	Continental subarctic
CA-TP3	Ontario - Turkey Point Middle-aged White Pine	2	Canada	42.71	-80.35	219	ENF	Warm summer continental
CA-TP4	Ontario - Turkey Point Mature White Pine	1	Canada	42.71	-80.36	219	ENF	Warm summer continental
CA-WP1	Western Peatland - LaBiche-Black Spruce/Larch Fen	1	Canada	54.95	-112.47	540	MF	Continental subarctic
US-ARM	OK - ARM Southern Great Plains Site - Lamont	1	USA	36.61	-97.49	310	CRO	Humid subtropical
US-Atq	AK - Atkasuk	1	USA	70.47	-157.41	16	WET	Tundra
US-Brw	AK - Barrow	1	USA	71.32	-156.63	1	WET	Tundra
US-Dk2	NC - Duke Forest - Hardwoods	1	USA	35.97	-79.10	160	DBF	Humid subtropical
US-Dk3	NC - Duke Forest - Loblolly Pine	1	USA	35.98	-79.09	163	ENF	Humid subtropical
US-Ha1	MA - Harvard Forest EMS Tower (HFR1)	1	USA	42.54	-72.17	303	DBF	Warm summer continental
US-Ho1	ME - Howland Forest (Main Tower)	1	USA	45.20	-68.74	60	ENF	Warm summer continental
US-IB1	IL - Fermi National Accelerator Laboratory - Batavia (Agricultural Site)	1	USA	41.86	-88.22	227	CRO	Hot summer continental
US-IB2	IL - Fermi National Accelerator Laboratory - Batavia (Prairie Site)	1	USA	41.84	-88.24	227	GRA	Hot summer continental
US-Los	WI - Lost Creek	1	USA	46.08	-89.98	480	CSH	Warm summer continental
US-MMS	IN - Morgan Monroe State Forest	1	USA	39.32	-86.41	275	DBF	Humid subtropical
US-MOz	MO - Missouri Ozark Site	1	USA	38.74	-92.20	219	DBF	Humid subtropical
US-Me2	OR - Metolius - Intermediate Aged Ponderosa Pine	1	USA	44.45	-121.56	1253	ENF	Dry-summer subtropical
US-Me3	OR - Metolius - Second Young Aged Pine	2	USA	44.32	-121.61	1005	ENF	Dry-summer subtropical

Site ID	Name	Priority	Country	Latitude	Longitude	Elevation (m a.s.l.)	IGBP Class	Köppen–Geiger Climate Classification
US-Me4	OR - Metolius - Old Aged Ponderosa Pine	2	USA	44.50	-121.62	915	ENF	Dry-summer subtropical
US-Me5	OR - Metolius - First Young Aged Pine	2	USA	44.44	-121.57	1183	ENF	Dry-summer subtropical
US-NR1	CO - Niwot Ridge Forest (LTER NWT1)	1	USA	40.03	-105.55	3050	ENF	Continental subarctic
US-Ne1	NE - Mead - Irrigated Continuous Maize Site	1	USA	41.17	-96.48	361	CRO	Hot summer continental
US-Ne2	NE - Mead - Irrigated Maize - Soybean Rotation Site	1	USA	41.16	-96.47	361	CRO	Hot summer continental
US-Ne3	NE - Mead - Rainfed Maize - Soybean Rotation Site	1	USA	41.18	-96.44	361	CRO	Hot summer continental
US-PFa	WI - Park Falls/WLEF	1	USA	45.95	-90.27	485	MF	Warm summer continental
US-SO2	CA - Sky Oaks - Old Stand	1	USA	33.37	-116.62	1392	CSH	Dry-summer subtropical
US-Shd	OK - Shidler- Oklahoma	1	USA	36.93	-96.68	350	GRA	Humid subtropical
US-Syv	MI - Sylvania Wilderness Area	1	USA	46.24	-89.35	540	MF	Warm summer continental
US-Ton	CA - Tonzi Ranch	1	USA	38.43	-120.97	177	WSA	Dry-summer subtropical
US-UMB	MI - University of Michigan Biological Station	1	USA	45.56	-84.71	234	DBF	Warm summer continental
US-Var	CA - Vaira Ranch - Ione	1	USA	38.41	-120.95	129	GRA	Dry-summer subtropical
US-WCr	WI - Willow Creek	1	USA	45.81	-90.08	520	DBF	Warm summer continental

1
2 Table 2 continued.

Site ID	Annual NEE (g C m ²)	Annual NEE Error (g C m ²)	Daytime Data Coverage (%)	Nighttime Data Coverage (%)	LAI	Annual Precipitation (mm)	Mean Annual Air Temperature (°C)	Measurement Period	Biome	Source
CA-Ca1	-244.3	61.1	99	26	6.1	1256	8.7	1998-2006	ENFT	Schwalm et al. [2007]
CA-Ca2	571.7	31.5	96	23	4.4	1222	8.8	2001-2006	ENFT	Schwalm et al. [2007]
CA-Ca3	91.2	37.9	91	27	2.2	1554	9.5	2002-2006	ENFT	Schwalm et al. [2007]
CA-Gro	-36.5	33.5	93	34	4.1	427	3.3	2004-2006	MF	McCaughey et al. [2006]
CA-Let	-132.9	14.3	96	46	0.7	335	6.5	1997-2006	GRA	Flanagan et al. [2002]
CA-Mer	-68.5	21.6	79	56	1.3	935	6.2	1999-2006	WET	Lafleur et al. [2003]
CA-Oas	-158.0	28.5	94	56	3.8	460	2.3	1997-2006	DBF	Barr et al. [2004]
CA-Obs	-56.3	16.1	89	45	5.6	470	1.6	2000-2006	ENFB	Griffis et al. [2003]
CA-Ojpp	-29.9	16.6	91	50	3.4	461	1.5	2000-2006	ENFB	Griffis et al. [2003]
CA-Qfo	-13.7	21.0	93	40	4	819	2.7	2004-2006	ENFB	Bergeron et al. [2007]
CA-SJ1	28.0	15.3	87	31	0.8	344	0.6	2002-2005	ENFB	Zha et al [2009]
CA-SJ2	117.0	6.1	89	47	1.3	537	0.1	2003-2006	ENFB	Zha et al [2009]
CA-SJ3	-82.0	17.7	92	34	4.3	694	0.8	2004-2005	ENFB	Zha et al [2009]

Site ID	Annual NEE (g C m ²)	Annual NEE Error (g C m ²)	Daytime Data Coverage (%)	Nighttime Data Coverage (%)	LAI	Annual Precipitation (mm)	Mean Annual Air Temperature (°C)	Measurement Period	Biome	Source
CA-TP3	-385.3	-	24	25	3.5	1011	8.7	2003-2007	ENFT	Peichl & Arain [2007]
CA-TP4	-133.2	29.5	95	43	3.5	959	8.6	2002-2007	ENFT	Peichl & Arain [2007]
CA-WP1	-195.8	16.4	96	50	2.7	481	1.7	2003-2007	WET	Syed et al. [2006]
US-ARM	-128.4	74.4	89	36	3.1	629	15.6	2000-2006	CRO	Sims & Bradford [2001]
US-Atq	-12.8	-	50	22	1.1	118	-10.6	1999-2006	TUN	Oberbauer et al. [2007]
US-Brw	-72.0	-	49	29	1.5	108	-10.9	1999-2002	TUN	Harazono et al. [2003]
US-Dk2	-718.1	-	48	1	7	1091	15.1	2003-2005	DBF	Siqueira et al. [2006]
US-Dk3	-350.0	139.0	75	37	5.6	1126	14.7	1998-2005	ENFT	Siqueira et al. [2006]
US-Ha1	-217.4	65.9	78	34	3.38	1122	7.9	1991-2006	DBF	Urbanski et al. [2007]
US-Ho1	-223.0	33.4	70	47	5.2	818	6.6	1996-2004	ENFT	Richardson et al. [2009]
US-IB1	-269.0	31.3	92	46	1.29	718	10.1	2005-2007	CRO	Post et al. [2004]
US-IB2	-86.0	42.0	80	49	5.38	818	10.4	2004-2007	GRA	Post et al. [2004]
US-Los	-78.0	19.2	82	54	4.24	666	3.8	2000-2006	WET	Sulman et al. [2009]
US-MMS	-346.1	66.3	97	46	4.9	1109	12.4	1999-2006	DBF	Schmid et al. [2000]
US-MOz	-305.7	48.9	94	33	3.91	730	13.3	2004-2007	DBF	Gu et al. [2006]
US-Me2	-536.0	65.8	63	46	3.62	434	7.6	2002-2007	ENFT	Thomas et al. [2009]
US-Me3	-198.0	32.7	83	28	0.52	423	8.5	2004-2005	ENFT	Vickers et al. [2009]
US-Me4	-612.3	-	55	41	2.1	641	8.3	1996-2000	ENFT	Irvine et al. [2004]
US-Me5	-206.0	10.6	97	48	1.1	350	7.6	1999-2002	ENFT	Irvine et al. [2004]
US-NR1	-37.2	27.0	89	44	4.2	663	2.5	1998-2007	ENFT	Bradford et al. [2008]
US-Ne1	-424.0	41.8	93	42	6.5	832	11.1	2001-2006	CRO	Verma et al. [2005]
US-Ne2	-382.0	41.8	96	51	6.5	823	10.8	2001-2006	CRO	Verma et al. [2005]
US-Ne3	-258.0	43.3	94	55	6.2	627	10.9	2001-2006	CRO	Verma et al. [2005]
US-PFa	45.0	41.1	85	30	4.05	736	5.1	1997-2005	MF	Davis et al. [2003]
US-SO2	22.4	25.6	87	30	3	695	13.8	1998-2006	SHR	Luo et al. [2007]
US-Shd	-75.5	22.0	96	49	5.9	1179	14.8	1997-2001	GRA	Suyker et al. [2003]
US-Syv	48.5	34.7	53	51	4.1	700	4.4	2001-2006	MF	Desai et al. [2005]
US-Ton	-67.8	52.0	77	25	0.6	549	16.4	2001-2007	WSA	Ma et al. [2007]
US-UMB	-132.0	42.4	86	39	4.23	629	7.4	1998-2006	DBF	Schmid et al. [2003]
US-Var	7.3	110.6	80	22	2.4	563	15.9	2001-2007	GRA	Ma et al. [2007]
US-WCr	-222.6	54.1	48	55	5.36	712	5.3	1998-2006	DBF	Cook et al. [2004a]

1 Sources: IGBP classification: Loveland et al. [2001]; Köppen–Geiger: Peel et al. [2007]; LAI for USA sites:
2 <http://public.ornl.gov/ameriflux/>; LAI for Canadian sites: Chen et al. [2006] & Schwalm et al. [2006]
3 Annual precipitation and mean annual air temperature are measurement period averages of meteorological inputs used to drive model
4 simulations.
5 *NEE* values show yearly integrals and associated error: one standard deviation based on uncertainty due to random noise and the
6 friction velocity threshold aggregated to yearly values and summed in quadrature [Barr et al., 2009].
7 Data coverages are percentages of half-hourly *NEE* measurements that satisfied quality control standards (friction velocity threshold)
8 for day- and nighttime separately.
9 Priority: 1 - Primary sites with complete (includes ancillary and biological data templates) records; 2 - Secondary chronosequence
10 sites
11 Standardized Precipitation Index available only for Priority 1 sites excluding US-Atq, US-Brw, US-Dk2, US-IB1 & US-Shd.
12 CA-TP3, US-Atq, US-Brw, US-Dk2 & US-Me4 sites used post-processing protocol from the La Thuile and Asilomar FLUXNET
13 Synthesis dataset [<http://www.fluxdata.org/>; Moffat et al., 2007; Papale et al., 2006] and lack *NEE* uncertainties.
14 Biome is combination of IGBP class and Köppen–Geiger climate.
15 US-Atq & US-Brw - Artic wetlands classified as tundra biome
16 CA-WP1 - treed fen (IGBP mixed forest) grouped with wetlands biome
17 US-Los - shrub wetland site (IGBP closed shrublands) grouped with wetlands biome
18 US-SO2 - closed shrublands grouped with shrublands (open or closed) biome
19 IGBP class and biome codes: CRO = croplands, CSH - closed shrublands, GRA = grasslands, ENF = evergreen needleleaf forest,
20 ENFB = evergreen needleleaf forest - boreal zone, ENFT = evergreen needleleaf forest - temperate zone, DBF = deciduous broadleaf
21 forest, MF = mixed (deciduous/evergreen) forest, WSA = woody savanna, SHR = shrublands, TUN = tundra, WET = wetlands.

1 Table 3. Model structural and site history predictors used to classify Taylor skill with
 2 regression tree analysis. Taylor skill (S ; Eq 3) was divided into three classes using
 3 terciles. Model structural predictants are from the Metadata for Forward (Ecosystem)
 4 Model Intercomparison survey collated by the NACP Site Synthesis
 5 (http://daac.ornl.gov/SURVEY8/survey_results.shtml). Site history data are from
 6 <http://public.ornl.gov/ameriflux/>, www.fluxnet.org, and Schwalm et al. (2006).

Predictor	Value
Model temporal resolution	Daily, half-hourly or less, hourly, monthly
Canopy	Prognostic, semi-prognostic, prescribed. Prescribed canopy from remote sensing, semi-prognostic has some prescribed input into canopy leaf biomass but calculates phenology with other prognostic variables.
Number of vegetation pools	Number of pools, both dynamic and static
Number of soil pools	Number of pools, both dynamic and static
Number of soil layers	Number of layers
Nitrogen	True if the model has a nitrogen cycle; otherwise false.
Steady state	True if the simulated long-term NEE integral approaches zero; otherwise false.
Autotrophic respiration (AR)	Fraction of annual GPP , fraction of instantaneous GPP , explicitly calculated, nil, proportional to growth
Ecosystem respiration (R)	$AR + HR$, explicitly calculated, forced annual balance
Gross primary productivity (GPP)	Enzyme kinetic model, light use efficiency model, nil, stomatal conductance model
Heterotrophic respiration (HR)	Explicitly calculated, first or greater order model, zero-order model
Net ecosystem exchange (NEE)	Explicitly calculated, $GPP - R$, $NPP - HR$
Net primary productivity (NPP)	Explicitly calculated, fraction of instantaneous GPP , $GPP - AR$, light use efficiency model
Overall model complexity	Low, average, high Values correspond to terciles of the total amount of first-order functional arguments for the following model-generated variables/outputs: AR , canopy leaf biomass, R , evapotranspiration, GPP , HR , NEE , NPP , soil moisture.
Site history	True if the below listed management activity or disturbance or event occurred on site; otherwise false. Grazed, fertilized, fire, harvest, herbicide, insects and pathogens, irrigation, natural regeneration, pesticide, planted, residue management, thinning
Stand age class	Young, intermediate, nil, mature, multi-cohort. Values based on stand age in forested sites; stands without a clear dominant stratum are treated as multi-cohort; non-forest types have nil.

1 Table 4. Normalized mean absolute error (*NMAE*) by climatic season, drought level, and
 2 biome. Drought level was based on monthly values of 3-month Standard Precipitation
 3 Index (SPI): dry value were < -0.8; wet > +0.8. Otherwise normal conditions existed.

Biome [#]	Climatic Season				Drought Level			Overall
	Winter	Spring	Summer	Fall	Dry	Normal	Wet	
CRO	1.90	4.64	-0.79	12.73	-1.43	-1.54	-1.59	-1.55
DBF	0.81	93.7	-0.52	-2.14	-1.01	-1.00	-0.95	-1.00
ENFB	1.52	-1.12	-0.69	-1.92	-0.87	-1.15	-3.43	-1.12
ENFT	-6.34	-0.66	-0.50	-0.76	-0.63	-0.72	-0.63	-0.68
GRA	-25.46	-0.84	-1.11	5.19	-1.52	-1.32	-3.07	-1.51
MF	1.10	-7.48	-0.47	57.70	-1.42	-1.04	-1.15	-1.12
SHR	-87.37	-1.37	-3.03	-140.17	-1.82	-2.18	-41.13	-2.88
TUN	-1.43	-11.07	-20.63	6.38	19.22	-24.06	-1.81	-20.15
WET	1.80	-5.07	-0.59	-4.72	-1.21	-1.20	-2.38	-1.27
WSA	-2.73	-0.75	-1.47	10.56	-1.39	-1.32	-1.51	-1.37
Overall	2.42	-1.35	-0.61	-1.94	-0.97	-1.01	-1.00	-1.00

[#] Biome codes: CRO = cropland, GRA = grassland, ENFB = evergreen needleleaf forest – boreal zone, ENFT = evergreen needleleaf forest – temperate zone, DBF = deciduous broadleaf forest, MF = mixed (deciduous/evergreen) forest, WSA = woody savanna, SHR = shrubland, TUN = tundra, WET = wetland.

4
 5 Table 5. Reduced χ^2 statistic by climatic season, drought level, and biome. Drought level
 6 was based on monthly values of 3-month Standard Precipitation Index (SPI): dry value
 7 were < -0.8; wet > +0.8. Otherwise normal conditions existed.

Biome [#]	Climatic Season				Drought Level			Overall
	Winter	Spring	Summer	Fall	Dry	Normal	Wet	
CRO	3.22	10.66	39.75	49.71	14.43	23.54	32.75	25.8
DBF	5.29	10.74	8.77	4.55	5.58	7.86	8.67	7.34
ENFB	21.25	17.75	4.98	6.61	11.64	12.02	18.51	12.61
ENFT	4.39	7.90	3.27	2.26	4.71	4.29	4.60	4.45
GRA	10.89	11.38	25.01	17.22	13.97	10.99	26.01	16.07
MF	3.74	4.67	2.05	2.02	2.92	3.24	2.98	3.08
SHR	13.34	27.98	12.52	11.2	9.26	21.31	10.31	16.26
WET	23.65	27.27	11.74	7.54	21.51	17.36	12.91	17.47
WSA	0.61	5.81	11.88	3.39	6.73	4.64	6.35	5.37
Overall	8.18	11.95	11.27	9.45	8.10	9.98	12.72	10.26

[#] Biome codes: CRO = cropland, GRA = grassland, ENFB = evergreen needleleaf forest – boreal zone, ENFT = evergreen needleleaf forest – temperate zone, DBF = deciduous broadleaf forest, MF = mixed (deciduous/evergreen) forest, WSA = woody savanna, SHR = shrubland, WET = wetland.

1 Figure Legends

2
3 Figure 1. Normalized mean absolute error ($NMAE$) by biome for each model. Biomes in
4 ascending order based on model-specific $NMAE$; biomes on the left show better average
5 agreement with observations. $NMAE$ is normalized by mean observed flux. Across all
6 sites, seasons, and drought levels within a given biome this value is negative ($NEE < 0$),
7 indicating a sink. $NMAE$ values closer to zero coincide with a higher degree of model-
8 data agreement. Woody savannahs and shrublands not shown: only one site each. Tundra
9 ($n = 2$ sites) has $NMAE < -10$ for all models. CN-CLASS croplands value is off-scale (= -
10 8.98). Black cross: no observations; white circle: undersampled ($n < 100$ months).

11
12 Figure 2. Normalized mean absolute error ($NMAE$) by climatic season and drought level.
13 $NMAE$ is normalized by mean observed flux such that most values are negative ($NEE <$
14 0), indicating a sink. Positive values indicate a source ($NEE > 0$). These occur in winter
15 for all models as well as spring and fall for all crop only models: AgroIBIS, DNDC,
16 EPIC, SiBcrop. Such values are displayed on the same color bar but with opposite sign.
17 Off-scale values: AgroIBIS and SiBcrop in fall are -7.1 and -11.1 respectively. DNDC in
18 fall and spring is -11.4 and -8.7 respectively. Black cross: no observations; white circle:
19 undersampled ($n < 100$ months).

20
21 Figure 3. Model skill metrics for all 22 models. Skill metrics are Taylor skill (S ; Eq. 3),
22 normalized mean absolute error ($NMAE$), and reduced χ^2 statistic (χ^2). Better model-data
23 agreement corresponds to the upper left corner. Benchmark represents perfect model-data
24 agreement: $S = 1$, $NMAE = 0$ and $\chi^2 = 1$. Gray interpolated surface added and model
25 names jittered to improve readability.

26
27 Figure 4. Boxplots of Taylor skill by model and site. Taylor skill (S ; Eq. 3) is a single
28 value summary of a Taylor diagram where unity indicates perfect agreement with
29 observations. Panels show interquartile range (blue box), median (solid red line), range
30 (whiskers), and outliers (red cross; values more than 1.5 x interquartile range from the
31 median). Top panel: only models ($n = 21$) used on at least two sites shown. Bottom panel:
32 only sites ($n = 32$) simulated with at least 10 unique models, excluding the mean model
33 ensemble (MEAN) and the assimilated LoTEC, shown. Models and sites sorted by
34 median Taylor skill.

35
36 Figure 5. Taylor diagram of normalized mean model performance for forested sites. Each
37 circle ($n = 26$ sites) is the site-specific mean model ensemble (MEAN). Benchmark (red
38 square) corresponds to observed normalized monthly NEE ; units of σ and $RMSE$ are
39 multiples of observed σ . Color coding of site letter and circles indicates biome: evergreen
40 needleleaf forest – temperate zone (red), deciduous broadleaf forest (brown), mixed
41 (deciduous/evergreen) forest (blue), evergreen needleleaf forest – boreal zone (black).
42 Outlying sites (evergreen needleleaf forest – boreal zone) not shown: CA-SJ1 ($\rho = 0.81$, σ
43 = 3.9, $RMSE = 3.1$) and CA-SJ2 ($\rho = -0.67$, $\sigma = 4.3$, $RMSE = 5.1$).

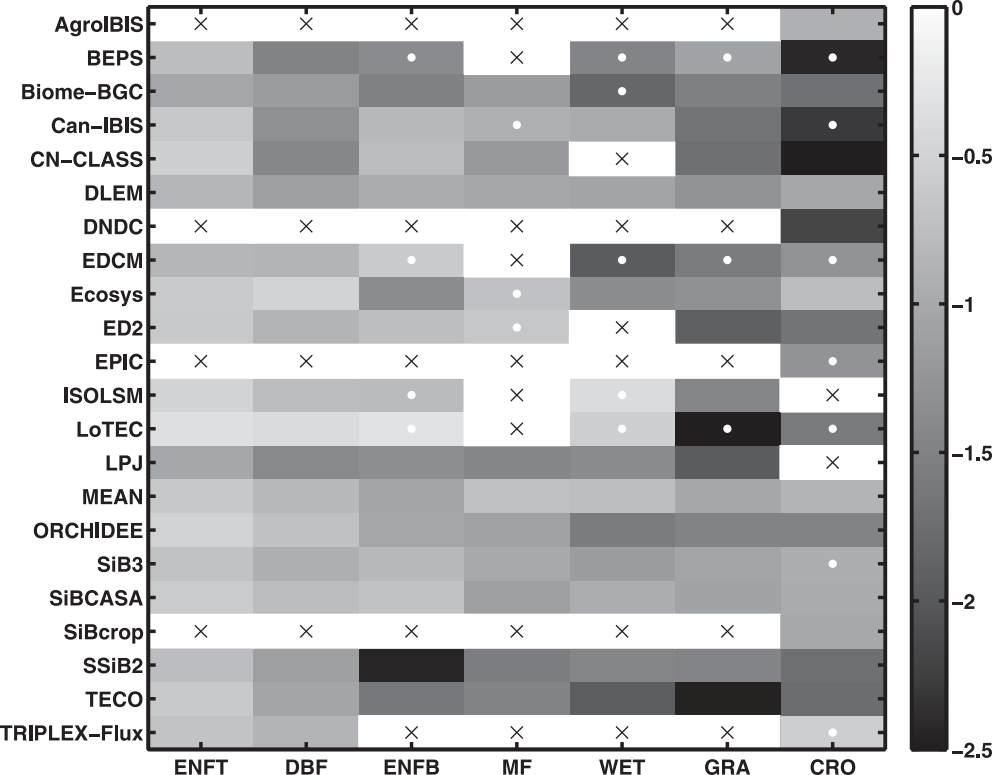
44
45 Figure 6. Taylor diagram of normalized mean model performance for non-forested sites.
46 Each circle ($n = 16$ sites) is the site-specific mean model ensemble (MEAN). Benchmark

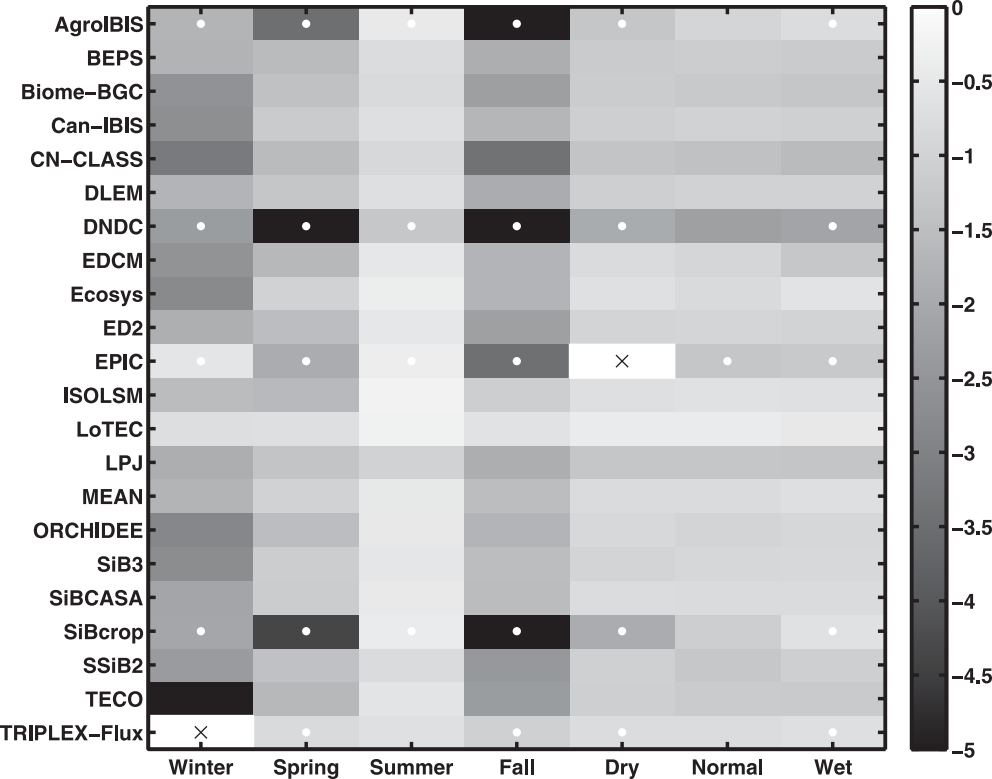
1 (red square) corresponds to observed normalized monthly *NEE*; units of σ and *RMSE* are
2 multiples of observed σ . Color coding of site letter and circles indicates biome: croplands
3 (red), grasslands (brown), wetlands (blue), all other biomes (black).

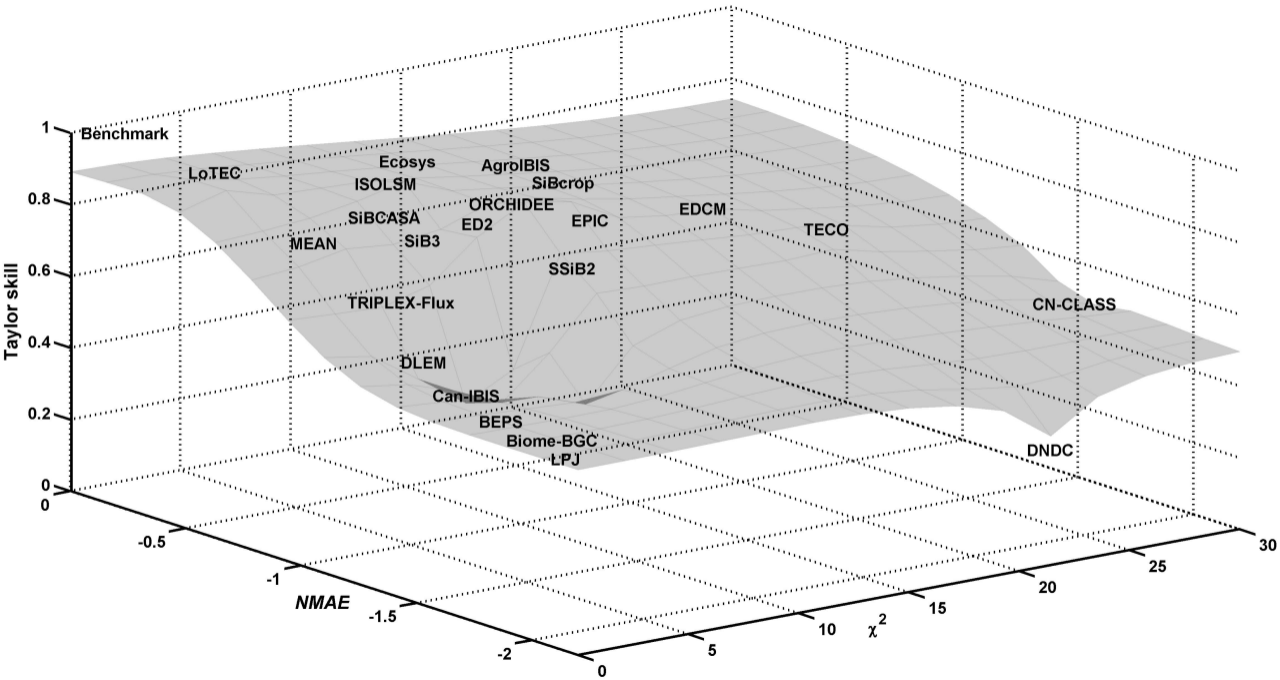
4
5 Figure 7. Taylor diagram of normalized across-site average model performance. Model σ
6 and *RMSE* were normalized by observed σ . Each circle ($n = 22$ models) corresponds to
7 the mean across all sites. Benchmark (red square) corresponds to observed normalized
8 monthly *NEE*; units of σ and *RMSE* are multiples of observed σ . Color coding of model
9 letter and circles indicates generality of model performance: specialist models used only
10 in croplands ($n \leq 5$ sites; black), generalist models used across a range of biomes and
11 sites ($n \geq 30$ sites, blue), all other models (red). The correlation for DNDC ($\rho = -0.13$) is
12 displayed as zero for readability.

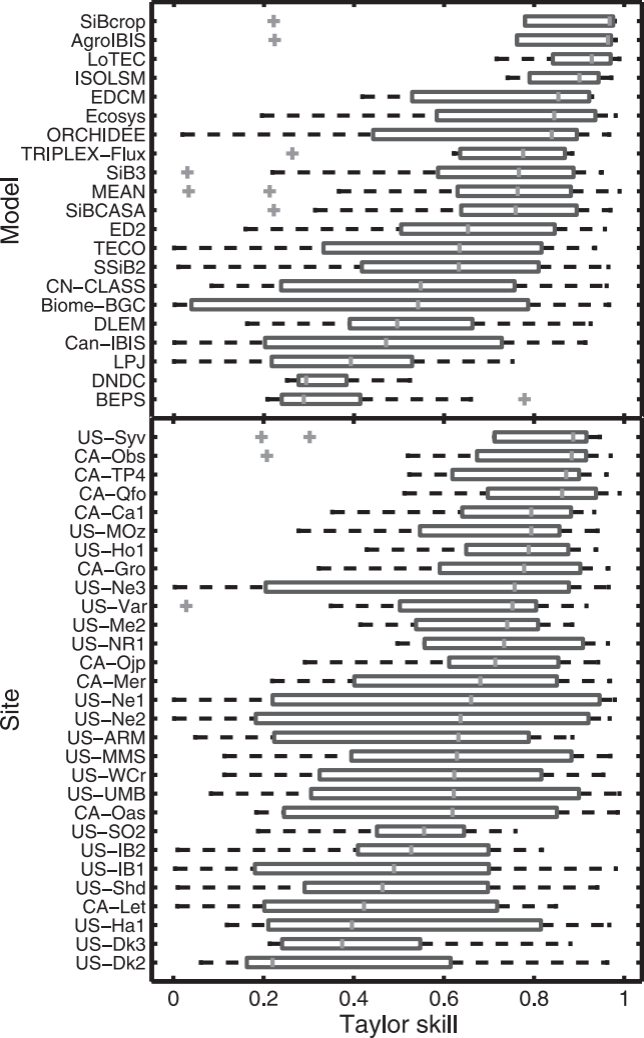
13
14 Figure 8. Variable importance scores for model-specific (blue) and site-specific (green)
15 predictors. Scores were generated from a regression tree with the Taylor skill classes
16 based on terciles ($n = 3132$) as the response. Only the 12 of 28 predictants with score $>$
17 25 shown; see Table 3 for complete listing of evaluated model structural and site
18 attributes.

19
20 Figure 9. Bar graphs of mean Taylor skill by model attribute. Whiskers represent one
21 standard error of the mean. Only model-specific attributes with variable important scores
22 > 25 shown. Note y-axis on right panels starts at 0.4.



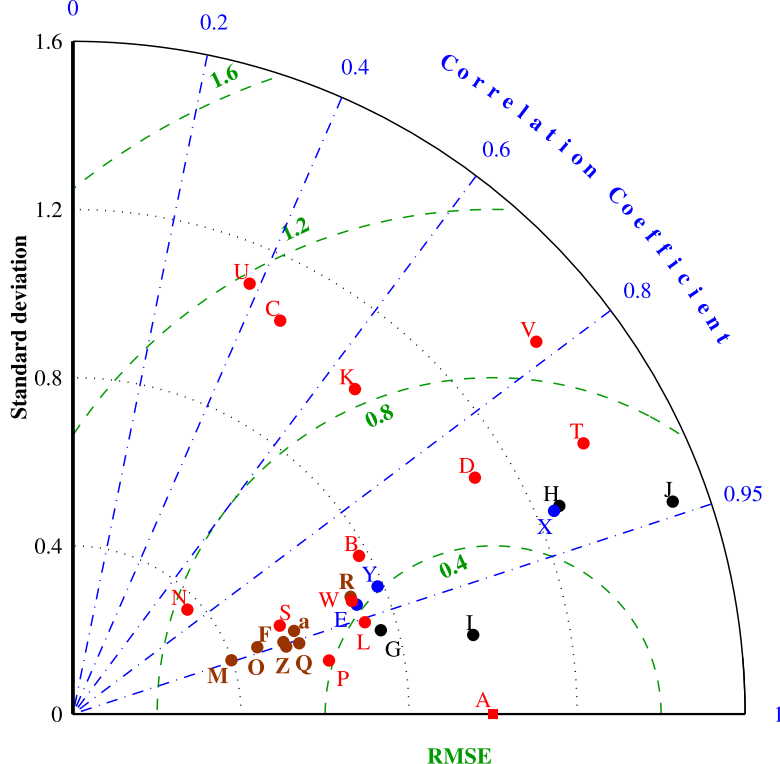




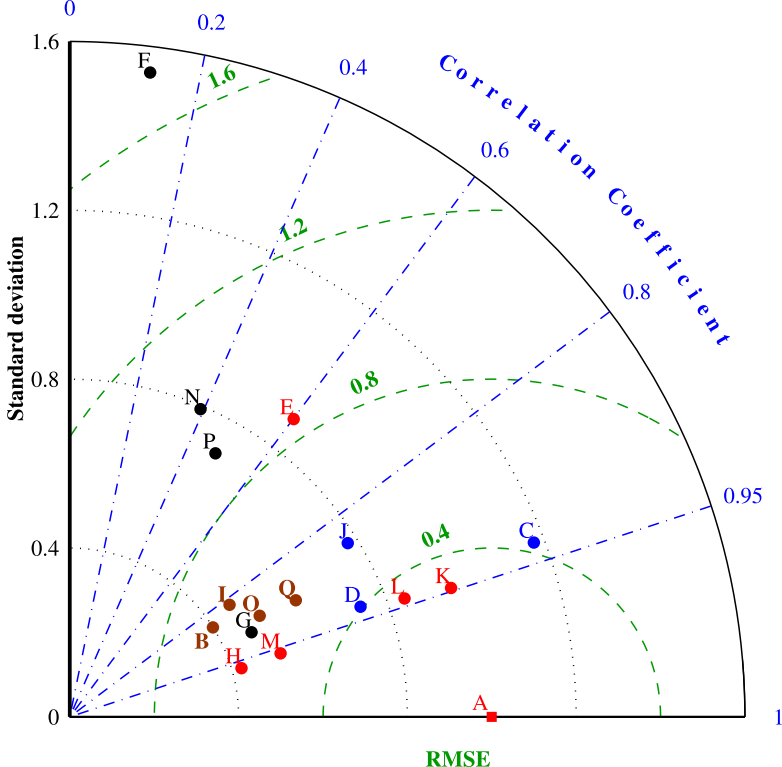


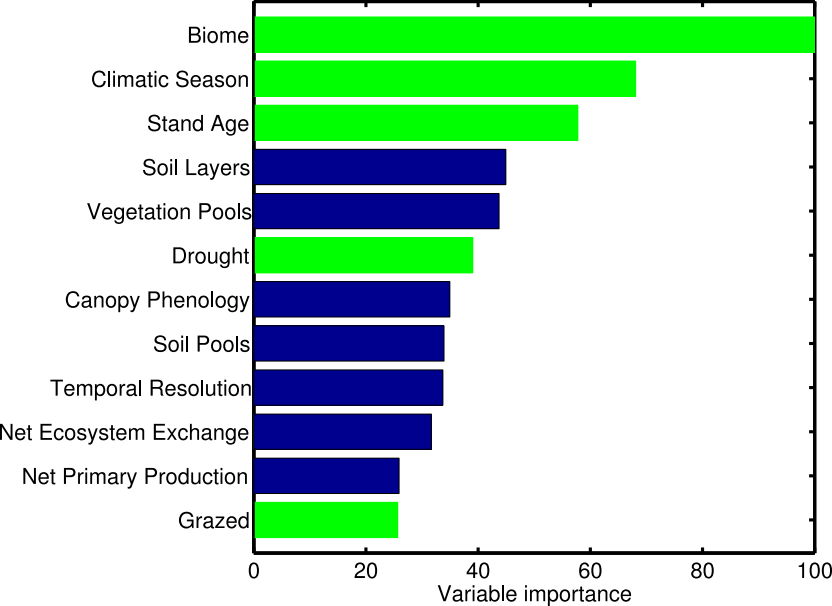
A Benchmark

- B CA-Ca1
- C CA-Ca2
- D CA-Ca3
- E CA-Gro
- F CA-Oas
- G CA-Obs
- H CA-Ojp
- I CA-Qfo
- J CA-SJ3
- K CA-TP3
- L CA-TP4
- M US-Dk2
- N US-Dk3
- O US-Ha1
- P US-Ho1
- Q US-MMS
- R US-MOz
- S US-Me2
- T US-Me3
- U US-Me4
- V US-Me5
- W US-NR1
- X US-PFa
- Y US-Syv
- Z US-UMB
- a US-WCr



- A Benchmark
- B CA-Let
- C CA-Mer
- D CA-WP1
- E US-ARM
- F US-Atq
- G US-Brw
- H US-IB1
- I US-IB2
- J US-Los
- K US-Ne1
- L US-Ne2
- M US-Ne3
- N US-SO2
- O US-Shd
- P US-Ton
- Q US-Var





Net Ecosystem Exchange

