

1 **A model-data intercomparison of CO₂ exchange during a large scale drought event:**

2 **Results from the NACP Site Synthesis**

3

4 Christopher Schwalm, Christopher Williams, Kevin Schaefer, NACP

5 contributors/collaborators

6

7 Manuscript for submission to JGR-Biogeosciences

8 Original version: October 7, 2009

9 Current version: November 24, 2009

10 Last saved: November 24, 2009

Abstract

Our current understanding of terrestrial carbon processes is represented in various models that are routinely used to integrate and scale observations of CO₂ exchange with remote sensing and other data in time and space. Yet rigorous accuracy assessments are rarely conducted to determine how well they mimic carbon processes across a range of vegetation types and environmental conditions. Here we compare observed and simulated integrals of CO₂ exchange using 44 eddy covariance flux towers in North America and model runs from 22 terrestrial biosphere models. The observational time record spans \approx 220 site-years over 10 biomes and overlaps with 2 large scale drought events in North America, providing a natural experiment to test model fidelity to observations as a function of drought and seasonality. We also evaluate the ability of the models to simulate the seasonal cycle of monthly CO₂ exchange using Taylor diagrams. Overall model performance was poor (< 0.50 predictive skill) with forested ecosystems better predicted than non-forested. Models achieved the highest levels of predictive skill in summer and winter in all biomes and in temperate evergreen forests during all climatic seasons. In contrast, model performance was consistently poor in spring and fall, especially in wetlands and mixed or deciduous forests, and during abnormally dry or wet periods. The best performing models were either highly specialized to specific biome types or used data assimilation techniques. Generalist models applied over a range of sites and biomes also performed well, with predictive skill in the upper quintile of all models.

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

Introduction

Drought is a reoccurring phenomenon in all climates (Larcher, 1995) and is characterized by a 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., 2004). 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 an intensification of the hydrological cycle; itself linked to global warming (Houghton et al., 2001; Huntington, 2006; Sheffield & Wood, 2008; Trenberth et al., 2007).

Forecasting the onset, amplitude, and duration of drought (Hwang & Carbone, 2009; Mishra & Desai, 2006; Sheffield & Wood, 2008) as well as CO₂ exchange under more frequent drought events related to climate change requires the use of models. Past validation studies of terrestrial biosphere models focused only on few models and sites typically in close proximity and primarily in forested biomes with differing objectives, levels of calibration to site data, and parameter values (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

1 model performance relative to drought requires high quality observed CO₂ exchange data,
2 a reliable drought metric as well as a natural experiment across sites and drought
3 conditions.

4
5 In this study we evaluate model performance using terrestrial CO₂ flux data and
6 simulated fluxes collected from 1991 to 2007. This timeframe included 2 widespread
7 droughts in North America: (i) the turn-of-the-century drought from 1998 to 2004 that
8 was centered in the interior West of North America (Seager et al., 2007a) and (ii) a
9 smaller-scaled drought event in the southern continental United States from winter of
10 2005/2006 through October 2007 (Seager et al., 2009). During these events low Palmer
11 Drought Severity Index values (Cook et al., 2007; Dai et al., 2004) and negative
12 precipitation anomalies (Seager et al., 2007a; 2009) occurred over broad geographic areas
13 which, in conjunction with ongoing eddy covariance measurements (Baldocchi et al.,
14 2001), provided flux data across gradients of time, space, seasonality, and drought. We
15 use this data to examine model performance relative to site-specific drought severity,
16 climatic season, and time and to link model behavior to model architecture and site-
17 specific attributes. Specifically, we address the following questions: Are current state-of-
18 the-art terrestrial biosphere models capable of simulating CO₂ exchange subject to
19 gradients in dryness and seasonality? Are these models able to reproduce the seasonal
20 variation of observed CO₂ exchange across sites? Are certain characteristics of model
21 structure coincident with better/worse model performance? Which biomes are simulated
22 poorly/well?

1 **Methods**

2 Modeled and observed net ecosystem productivity (NEP, net carbon balance including
3 soils where negative values indicate outgassing of CO₂ to the atmosphere) data were
4 analyzed from 22 terrestrial biosphere models (Table 1) and 44 eddy covariance (EC)
5 sites spanning \approx 220 site-years and 10 biomes in North America (Table 2). Simulated
6 NEP was based on model-specific runs using gap-filled observed weather at each site and
7 locally observed values of soil texture according to a standard protocol ([http://isynth-](http://isynth-site.pbworks.com)
8 [site.pbworks.com](http://isynth-site.pbworks.com)). Each model, except the mean model ensemble and an assimilated
9 model, was spun up to steady state initial conditions with a target NEP of zero integrated
10 over the simulation period.

11
12 Gaps in the meteorological data record occurred at EC sites due to data screening or
13 instrument failure. Missing values of air temperature, humidity, shortwave radiation, and
14 precipitation data, i.e., key model inputs, were filled using DAYMET (Thornton et
15 al.,1997) before 2003 or the nearest available climate station in the National Climatic
16 Data Center's Global Surface Summary of the Day (GSOD) database. Daily GSOD and
17 DAYMET data were temporally downscaled to hourly or half-hourly using the phasing
18 from observed mean diurnal cycles calculated from a 15-day moving window.

19
20 EC data were produced by AmeriFlux and Fluxnet Canada investigators and processed as
21 a synthesis product of the North American Carbon Program (NACP) Site Level Interim
22 Synthesis (<http://www.nacarbon.org/nacp/>). The observed NEP were corrected for
23 storage, despiked, filtered to remove conditions of low turbulence (friction velocity

1 filtered), and the gaps filled to create a continuous time series (Barr et al., 2004). The
2 time series included estimates of random uncertainty and uncertainty due to the friction
3 velocity filtering (Barr et al., 2004). In this analysis NEP was aggregated to monthly
4 integrals using only non-gap-filled data, i.e., observed values deemed spurious and
5 subsequently infilled were not considered. Coincident modeled NEP values were
6 similarly excluded. This removed the influence of gap-filling algorithms in the
7 comparison of observed and modeled NEP.

8
9 Drought level was quantified using the 3-month Standard Precipitation Index (SPI,
10 McKee et al., 1993). Monthly SPI values were taken from the U.S. Drought Monitor
11 (<http://drought.unl.edu/DM/>) whereby each tower was matched to nearby meteorological
12 station(s) indicative of local drought conditions given proximity, topography, and human
13 impact. This study used 3 drought levels: Dry required $SPI < -0.8$, wet corresponded to
14 $SPI > +0.8$, otherwise normal conditions existed. Climatic season was defined by 4
15 seasons of 3 months each with winter given by December, January, and February.

16
17 Model performance was evaluated using data groups defined by all possible intersections
18 of site, model, climatic season, and drought level with at least 3 pairs of observed and
19 simulated monthly NEP. All available months within each data group were
20 simultaneously evaluated for the presence or absence of predictive skill based on
21 correlation (ρ), root mean squared error (RMSE), and the standard error or observed NEP
22 (SE):

$$\rho_{ijkd}^{o,m} = \frac{(NEP_{ijkd}^o - \overline{NEP_{ijkd}^o})(NEP_{ijkd}^m - \overline{NEP_{ijkd}^m})}{\sqrt{\sigma_{NEP_{ijkd}^o} \sigma_{NEP_{ijkd}^m}}},$$

2

$$RMSE_{ijkd}^{o,m} = \left(\frac{\sum (NEP_{ijkd}^o - NEP_{ijkd}^m)^2}{N_{ijkd}} \right)^{0.5},$$

4

5 where the overbar indicates averaging across all values, σ is the standard deviation across
6 all months, N is the amount of paired observed and modeled monthly NEP integrals in a
7 data group, subscript i is for site, j is for model, k is for climatic season, d is for drought
8 level, and superscript o is for observations and m is for modeled estimates.

9

10 SE comprised 2 components calculated explicitly for monthly NEP integrals: random
11 uncertainty and uncertainty associated with the friction velocity threshold (u_*^{Th}). Random
12 uncertainty was estimated following Richardson & Hollinger (2007): (i) generate
13 synthetic NEP data using the gap-filling model (Barr et al., 2004), (ii) introduce gaps as
14 in the observed data with u_*^{Th} filtering, (iii) add noise, (iv) infill gaps in synthetic data,
15 (v) repeat the process 1000 times. The random uncertainty SE component was then the
16 standard deviation across all bootstrap replicates. The u_*^{Th} uncertainty component was
17 also estimated using the bootstrap. Here 1000 bootstrap replicates of NEP were generated
18 using a distribution of u_*^{Th} based on binning the raw flux data with respect to climatic
19 season, temperature, and site-year. The standard deviation across all replicates gave the
20 SE component associated with u_*^{Th} . These two components were then combined in
21 quadrature, i.e., $z = \sqrt{x^2 + y^2}$, to estimate the total SE of monthly observed NEP.

1

2 The presence of predictive skill (π_{ijkd}) required (i) that modeled and observed flux data
 3 were directly related (not anti-correlated) with at least a weak correlation serving as a
 4 threshold, (ii) that model error was less than the mean observed flux, i.e., that a
 5 discrimination between source and sink was possible, and (iii) that the variability of
 6 model-data agreement was scaled to the uncertainty of observed and simulated fluxes.
 7 Mathematically, predictive skill for a single data group was present ($\pi_{ijkd} = 1$) if the
 8 following three conditions were satisfied:

$$\begin{aligned}
 & (i) \quad \rho_{ijkd}^{o,m} \geq 0.2 \\
 & (ii) \quad \frac{RMSE_{ijkd}^{o,m}}{NEP_{ijkd}^o} \leq 1 \\
 & (iii) \quad \overline{NEP_{ijkd}^o} \pm 2\overline{SE}_{ijkd}^o \cap \overline{NEP_{ijkd}^m} \pm 2RMSE_{ijkd}^{o,m}
 \end{aligned}$$

10 Otherwise there was no predictive skill ($\pi_{ijkd} = 0$). This framework was used to aggregate
 11 across data groups as well. For example, calculating predictive skill for each model
 12 (π_{model}) required integration across all sites, climatic seasons, and drought levels by
 13 model. Similarly, biome-specific predictive skill (π_{biome}) was calculated by integrating
 14 across all sites, models, climatic seasons, and drought levels within a given biome. When
 15 aggregating in this manner the group-wise predictive skills were linearly combined
 16 weighted by N .

17

18 We defined the consistency of predictive skill by model (CV_π) as the coefficient of
 19 variation, in percent, of π across all sites simulated: $CV_n = \frac{\sigma_\pi}{\pi}$, where σ_π is the standard
 20 deviation of π across site by model. Consistency addresses the degree to which model

1 performance was related to model parameter set (site or biome unique constants) and
2 assumed that the effect of any errors in model structure and meteorological forcings on π
3 was randomly distributed. High consistency (relatively low CV_{π}) indicated less error in
4 model parameter set relative to a set level of predictive skill and vice versa.

5
6 The existence of predictive skill was then related to biome, climatic season, drought
7 level, site history and model structure using regression tree analysis (RTA) as a
8 supervised classification algorithm (Table 3). RTA is a form of binary recursive
9 partitioning (Breiman et al., 1984) that successively splits the data into subsets (nodes) by
10 minimizing within-subset variation. The result is a pruned tree-like topology whereby
11 predicted values (presence/absence of predictive skill) are derived by a top-to-bottom
12 traversal following the rules (branches) that govern subset membership until a predicted
13 value is reached (terminal node). The splitting rules at each node as well as its position
14 allow for a calculation of relative variable importance (Breiman et al., 1984) with the
15 most important variable given a score of 100.

16
17 A second characterization of model performance was done using Taylor diagrams
18 (Taylor, 2001); visual displays based on pattern matching, i.e., the degree to which
19 simulations matched the temporal evolution of monthly NEP. Taylor plots are polar
20 coordinate displays of ρ , RMSE, and σ of NEP where all 3 quantities were calculated
21 using the full data record by site and model (ranging from 7 to 178 months). Taylor
22 diagrams were constructed for the mean model ensemble and across-site mean model
23 performance. Each point representing a model was subsequently scored using a skill

metric: $S = 2(1 + \rho) / (\sigma_{norm} + 1/\sigma_{norm})^2$, where σ_{norm} is the ratio of simulated to observed σ (Taylor, 2001), bound by zero and unity where unity indicated perfect agreement. This scoring allowed a direct comparison with predictive skill (π), also bound by zero and unity. In sum, the calculation of predictive skill evaluated model performance relative to ecological controls on CO₂ exchange independent of time whereas the Taylor displays examined models' ability to mimic the monthly trajectory of observed NEP.

Results

Model performance based on predictive skill

Predictive skill was related to drought level, climatic season, model, and site ($p = 0.01$). Of the 26747 paired data-model months (across 3479 data groups) available for analysis 10530 ($\pi_{overall} = 0.39$) showed predictive skill. Overall performance was higher in forested than non-forest biomes with the highest performance levels occurring during periods of peak biological activity (climatic summer) under non-drought conditions (Table 4). While drought was a significant control on predictive skill the aggregated contrast was only significant for non-normal (abnormally wet and dry) vs. normal conditions.

Integrated across all models and sites predictive skill was, both within and across drought level, higher in summer and winter than in spring and fall (Table 4). A similar pattern, except for ENFT, was present at the biome level (Table 5) with severe losses in predictive skill occurring during climatic spring and fall, especially in DBF, MF, and WET. Across the 3 levels of dryness π_{biome} declined for non-normal hydrological conditions apart from CRO and ENFB where abnormally dry and wet conditions

1 respectively were better predicted. Overall biome-level predictive skill was loosely
2 ranked in 4 tiers: ENFT > ENFB, MF, DBF, CRO > GRA, WET, WSA > SHR. These
3 rankings were robust across models used in the majority of biomes although some
4 divergence was apparent for MF and WET (Figure 1).

5
6 By model, predictive skill showed, similar to the biome level, a sharp decline in spring
7 and fall with the best model performance achieved through data assimilation (LoTEC).
8 Also, most models showed the highest performance in normal hydrological conditions
9 (Figure 2). Overall performance by model (π_{model}) ranged from 0.04 to 0.68 (Figure 3);
10 CV_{π} had larger amplitude, ranging from 34% to 174%. These 2 characterizations of
11 model performance were inversely related ($r = -0.8$; $p < 0.05$) i.e., higher levels of
12 predictive skill were associated with higher levels of consistency in predictive skill.

13
14 Among crop models, SiBCrop and AgroIBIS performed well, especially in climatic
15 spring and during dry conditions. In contrast, DNDC exhibited 0.04 predictive skill
16 overall, $\pi_{model} = 0.55$ in climatic winter and zero otherwise. All crop only models ($n = 4$)
17 showed $\pi_{model} = 0$ for wet conditions. While the best model performance across all CRO
18 sites was achieved by a specialist model (AgroIBIS: $\pi_{biome} = 0.56$), the generalist Ecosys
19 model, used on 26 sites in 9 biomes, achieved a similar level of performance ($\pi_{biome} =$
20 0.54). Overall, i.e., across all sites, dryness levels, and climatic seasons, LoTEC had the
21 highest performance (0.68). This platform was optimized using a data assimilation
22 technique, unique among model runs evaluated here, and was applied at 10 sites.
23 Performance by individual models ($n = 10$) applied to a wider range of sites (at least 24

1 sites) was highest for SibCASA (0.51) and lowest for TECO (0.24). Model performance
2 in excess of 0.50 was achieved only with the mean model ensemble (MEAN), a generalist
3 model (SibCASA), an optimized platform (LoTEC), crop models (AgroIBIS, SiBcrop)
4 and 1 model used on 9 sites (ISOLSM) (Figure 3).

5
6 Site-level predictive skill also showed a high degree of variability. In 2 instances (1 forest
7 and croplands site) predictive skill across all models for a given site ranged from zero to
8 unity. Even for the best predicted site on average (US-Ho1) π_{model} ranged from 0.43 to 1
9 (Figure 4). Alternatively, no one site was predicted well ($\pi_{model} > 0.5$) by all models;
10 whereas 4 sites (CA-Let, US-ARM, US-SO2, and US-Ton) were consistently poorly
11 predicted ($\pi_{model} < 0.5$). In contrast, mean model performance was notably high at some
12 sites as a function of drought level and seasonality: US-Dk3 ($\pi_{site} = 0.88$), US-Ho1 (0.81),
13 and US-Ho1 (0.85) for dry, normal, and wet conditions respectively as well as US-Me2
14 (0.79), US-Ho1 (0.98), US-PFa and US-WCr (0.90), and US-Ho1 (0.94) for winter,
15 spring, summer, and fall respectively. Despite the wide range in model performance,
16 π_{model} was not related to the number of sites ($p = 0.9$) nor biomes ($p = 0.7$) simulated, i.e.,
17 using a more general rather than a specialized model did not result in a loss in model
18 performance. Similarly, model performance was not better at sites with longer term
19 records on average ($p > 0.16$), although ORCHIDEE did perform marginally better on
20 sites with more data ($p < 0.09$).

21
22 *Factors governing the distribution of predictive skill*

Biome and climatic season were the most important factors in the distribution of predictive skill (Figure 5) followed by 2 model structural attributes: the number of dynamic vegetation pools and overall model complexity (Table 3). Stand age was the most important site-specific attribute whereas only 3 of the 13 evaluated site disturbances achieved an importance score of at least 10 and were of tertiary importance overall. Models with lower complexity (the lowest 2 quintiles of overall complexity) performed better at ENFT sites and models with at least 5 dynamic vegetation pools fared better at non-ENFT sites. For SHR and WSA no model structural attribute was coincident with the presence of predictive skill; these biomes were always poorly predicted. While the regression tree algorithm achieved an accuracy of 79% for predicting the presence or absence of predictive skill, the site and model characteristics considered here did not explain the underlying cause of biome and seasonal differences in predictive skill.

Overall model performance using Taylor diagrams

Average model performance (both across-site and across-model) was evaluated using Taylor diagrams based on all simulated and observed NEP monthly integrals. Better model performance was indicated by higher S (Figure 3) values and proximity to the benchmark. 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 predictive skill, forested sites were better predicted than non-forested ones. The MEAN model, i.e., average model performance across all models, at all sites (excluding CA-SJ2 and US-Atq) showed $\rho \geq 0.2$ but generally (33 of 44 sites) underpredicted the variability associated with monthly NEP at forested (Figure 6) and non-forested (Figure 7) sites.

Similarly, 40 of 44 sites were predicted with $RMSE < \sigma_{obs}$. Also 8 (6 forested and 2 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 8). Most models underpredicted variability and showed $RMSE < \sigma_{obs}$. 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 one 1 site), ISOLSM (used on 9 sites), LoTEC (data assimilation model), SiBcrop and AgroIBIS (crop only models), EDCM (used on 10 sites), ecosys and SiBCASA (more generalist models used on 39 and 35 sites respectively), and MEAN (mean model ensemble). All of these “best” models had $\rho > 0.8$, $RMSE < 0.8$ 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. Finally, the general model rankings with regards to any single measure of model performance were preserved, i.e., higher π was coincident with higher S and lower CV_{π} ($p < 0.05$; $|r| \approx 0.5$).

Discussion

As expected no single model performed well across the full spectrum of dryness, seasonality, and sites simulated. Exact agreement between simulations and observations is not however a realistic expectation: models can match observations only within the observational uncertainty. While flux measurements are useful for validating and

1 improving the process-based simulation of CO₂ exchange (Friend et al., 2007),
2 accounting for the weaknesses of the observational record is central to model validation
3 (Grant et al., 2005; Hanson et al., 2004). In this study the uncertainty associated with
4 monthly NEP integrals was used to define predictive skill. This created an observational
5 interval, as opposed to a point estimate, allowing for a more accurate estimate of
6 predictive skill. Nonetheless model performance was low (there were few instances
7 where $\pi > 0.50$), underscoring the continued need for models to improve fidelity to
8 observations (Friedlingstein et al., 2006).
9
10 Contrasting predictive skill and the skill metric based on the Taylor diagrams (S vs. π)
11 highlighted difficulties in validating process-based models. Taylor diagrams were
12 indexed to time only whereas predictive skill was evaluated in the context of a factorial
13 design with 3 levels of dryness and 4 levels of seasonality. Overall, models were better in
14 capturing the temporal dynamics of monthly NEP than observed relationships between
15 carbon flux, dryness, and seasonality. However averaging over these ecologically
16 relevant factors can cause under- and overpredictions to cancel (Medlyn et al., 2005).
17 Alternatively, aggregate results can overestimate model fidelity to observations at shorter
18 timescales. It is noteworthy that metrics of performance skill were closely related. This
19 implies that evaluating model performance, especially relative to other models, is
20 possible using temporal evolution only but that doing so has limited diagnostic ability in
21 terms of highlighting model weaknesses and overestimates how well a model captures
22 key processes of carbon metabolism.
23

Effect of model structure and parameterization on model performance

The model characteristics considered here could not explain the underlying causes of the biome and seasonal differences in predictive skill (Figure 5). Lower variable importance scores revealed only a subsidiary link between predictive skill and the model structural attributes evaluated here: temporal resolution, types of algorithms for carbon balance terms, overall model complexity, number of vegetation and soil pools, soil layers, and diagnostic vs. prognostic leaf area index. In contrast, the effect of model parameterization was larger than model structure. 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 simulating CO₂ exchange than interannual climatic variability (Amthor et al., 2001). This is related to the use of biome-specific parameters relative to within-biome variability (Purves & Pacala, 2008). A corollary occurs in the context of validating EC observations as tower footprints can exhibit heterogeneity, particularly in soils, that is not reproduced in model site-specific parameters (Amthor et al., 2001).

The importance of model parameter sets was visible in this intercomparison in two ways. Firstly, biome had the highest variable importance score. Inasmuch as models rely on biome-specific parameter values this finding indicates model parameter sets as a key factor in the distribution of predictive skill. Furthermore, the high degree of variability in predictive skill across sets of biome-specific constants underscores that biomes may be too heterogeneous in time (Stoy et al., 2005; 2009) and space to assume constant parameter values relative to within-biome climate partitions (Hargrove et al. 2003).

1 Secondly, given that some model constants are site-specific scoring models for
2 consistency (calculating the coefficient of variation of predictive skill across all sites for
3 each model, Figure 1) highlighted the uniqueness of model results related to model
4 parameter sets using site-specific runs as a proxy. That is, the high degree of within-
5 model variation in model skill suggested that model parameter sets may need to be
6 refined to capture local, site-specific realities.

8 *Links between model performance and environmental factors*

9 Drought level was significantly related to predictive skill with the largest effect occurring
10 during climatic summer when dry (and wet) conditions were associated with a decline in
11 model performance. As summer corresponds to peak biological activity, the ability of
12 models to simulate NEP during anomalous climatic conditions influences carbon budgets
13 on longer time scales, i.e., process uncertainty with regards to dry (Sitch et al., 2008) and
14 wet conditions must be reduced to enhance model utility. While normal hydrological
15 conditions were generally better predicted exceptions to this rule occurred during climatic
16 fall as well as in ENFB and CRO where non-normal conditions showed the highest
17 degree of predictive skill. This indicated that the variability around normal conditions
18 was large and models' behavior insensitive, i.e., changes in driving variables within the
19 range of normal conditions were too subtle for the model but not for the simulated sites.
20 Finally, ecosystem response to longer-term drought can exhibit lags and positive
21 feedbacks (Arnone et al., 2008; Granier et al., 2007; Thomas et al., 2009; Williams et al.,
22 2009) that were not explicitly included in the drought metric used in this study.

1 In spring and fall, corresponding to leaf initiation and senescence for biomes with a
2 significant deciduous component, models showed the least predictive skill (see also
3 Morales et al., 2005). Springtime phenological cues influence the annual carbon balance
4 at point to continental scales (Barr et al., 2007; Delpierre et al., 2009; Keeling et al.,
5 1996). Similarly how well simulated fluxes fit observations has been linked to correctly
6 reproducing leaf initiation and senescence (Hanson et al., 2004). However, the observed
7 trend toward earlier leaf out does not appear coincident with an increase in the terrestrial
8 sink due to warmer/drier summers (Angert et al., 2005; Piao et al., 2008). This suggests
9 that a possible reason for the general loss of predictive skill in spring and fall was related
10 to how the initiation and cessation of canopy photosynthesis were modeled relative to soil
11 moisture and soil temperature (Waring & Running, 2007). In wetlands a similar decrease
12 in predictive skill was observed in the shoulder season. This was likely associated with
13 the general inability of models to simulate assimilation and respiration of ericaceous
14 shrubs and mosses common to wetlands (St-Hilaire et al., 2008).

15
16 ENFT and ENFB diverged in performance during fall and winter although these biomes
17 differ in climatic zone only. A similar divergence was reported using Biome-BGC, LPJ
18 and ORCHIDEE to simulate gross CO₂ uptake across a temperature gradient in Europe
19 (Jung et al., 2007a); average relative RMSE was higher for ENFB. This was linked to an
20 overestimation of leaf area index at the boreal sites and relationships between resource
21 availability and leaf area (Friedlingstein et al., 2006; Jung et al., 2007a; Stich et al.,
22 2008). Additionally, recent observations in the circumboreal region, where all ENFB sites
23 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). We speculate the loss of predictive skill in ENFB relative to ENFT 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.

Site history and constraints on model evaluation

Disturbance regime and how a model treats disturbance are known to impact model performance (Ito, 2008). In this study stand age and, to a lesser extent, planting, grazing, and fire (but not harvest) were controls on predictive skill. However comparing sites with and without fire activity showed no significant difference in mean predictive skill ($p > 0.10$), i.e., disturbance and age were confounded. Furthermore, sites with a recent stand replacement disturbance were poorly predicted, e.g., CA-SJ2, the worst predicted site, was harvested in 2000 and scarified in 2002, and US-SO2, a second poorly predicted site, suffered catastrophic wildfire during the analyzed data record. The poor model performance for recently disturbed sites follows from the steady state spinup used in simulation and the absence of an explicit disturbance layer. However, the distribution of site history metrics was skewed; only few sites were in the early stages of recovery from disturbance when NEP is more nonlinear relative to established stands. Furthermore, age class was biased toward older stands; of the 17 forested sites only 1 was classified as a young stand. Other site characteristics were also unbalanced; of the 29 sites evaluated for predictive skill 17 were forested with 4 or less CRO, GRA, SHR, WET, and WSA sites.

1 While regression trees are inherently robust additional observed and simulated fluxes in
2 rapidly growing young forested stands and undersampled biomes are desirable to better
3 characterize model performance.
4
5 Several aspects of the simulation protocol influenced the interpretation of our
6 performance analysis. Firstly, this analysis focused solely on non-gap-filled data to allow
7 the model-data intercomparison to inform model devolvement. However, the low
8 turbulence (friction velocity) filtering removes many more data at night than during the
9 day, so our analysis may be skewed towards daytime conditions. Secondly, each model
10 that used remotely sensed inputs (such as leaf area index) repeated an average seasonal
11 cycle calculated from site-specific time series based on all pixels within 1 km of the
12 tower site. This likely deflated relevant variable importance scores (Figure 5) and
13 precluded a full comparison of diagnostic vs. prognostic LAI. While only 6 models used
14 such inputs, including the best performing generalist model SiBCASA, removing the
15 inherent bias in using an invariant seasonal cycle over multiple years may improve model
16 performance. Incorporating disturbance information to recreate historical land use and
17 disturbance, especially for recent site entries, could also improve model performance.
18 Lastly, the models used here, except LoTEC and the mean model ensemble, were spun up
19 to steady state initial conditions. However, few, if any, of the sites are actually at steady
20 state, resulting in an inherent bias between simulated and observed NEP. Relaxing the
21 steady state assumption (Carvalhais et al., 2008) or initializing using observed wood
22 biomass and the quasi-steady state assumption (Schaefer et al., 2008) could improve
23 model performance.

1

2 Despite these constraints the distribution of predictive skill and overall skill metrics
3 indicated that the performance of comprehensive generalist models can equal that of
4 specialist site-specific models (see also Grant et al., 2005). This has important
5 implications for model choice and recasts model utility as a tradeoff between intensive
6 biome and site-specific parameterizations vs. generic parameter sets linked to gridded
7 datasets or remotely sensed data. Lastly the use of data assimilation techniques and mean
8 model ensembles greatly increased model predictive skill.

9

10 **Conclusion**

11 We used observed CO₂ exchange from 44 eddy covariance towers in North America with
12 simulations from 22 terrestrial biosphere models to validate model performance during a
13 large scale drought event. Overall forested sites were better predicted than non-forested
14 sites. Weaknesses in model performance concerned model parameter sets, growing
15 season length, especially in biomes with a significant deciduous component, and
16 abnormally dry or wet conditions. Undersampled biomes (grasslands, shrublands,
17 wetlands, woody savannah, and tundra) also showed a large divergence between
18 observations and simulations. The best model performance occurred in temperate
19 evergreen forests across all climatic seasons and during winter and summer across all
20 biomes. Models' ability to match observed monthly integrals of net ecosystem
21 productivity across levels of dryness and seasonality was generally poor (< 0.50
22 predictive skill). In contrast, performance was higher when evaluating month-to-month
23 trajectories. This indicated that the temporal evolution of NEP is better modeled than

1 responses to finer scaled changes in driving variables, albeit at the cost of under and
2 overpredictions that cancel out over longer-term simulations. Models with relatively high
3 performance metrics included generalist models applied over a wide range of sites and
4 biomes as well as highly specialized (used on croplands or a single site only), optimized
5 (parameters tuned using data assimilation), or ensembles (linear combinations of several
6 models).

8 **Acknowledgements**

9 CRS, CAW, and KS were supported by the U.S. National Science Foundation grant
10 ATM-0910766. We would like to thank the North American Carbon Program Site-Level
11 Interim Synthesis team, the Modeling and Synthesis Thematic Data Center, and the Oak
12 Ridge National Laboratory Distributed Active Archive Center for collecting, organizing,
13 and distributing the model output and flux observations required for this analysis. This
14 study was in part supported by the U.S. National Aeronautics and Space Administration
15 (NASA) grant NNX06AE65G, the U.S. National Oceanic and Atmospheric
16 Administration (NOAA) grant NA07OAR4310115, and the U.S. National Science
17 Foundation (NSF) grant OPP-0352957 to the University of Colorado at Boulder.

19 **References**

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

1

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

5

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

9

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

12

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

16

17 Baldocchi, D (2008) Breathing of the terrestrial biosphere: lessons learned from a global
18 network of carbon dioxide flux measurement systems. Australian Journal of Botany, 56,
19 1-26.

20

21 Baker IT, Prihodko L, Denning AS, Goulden M, Miller S, da Rocha HR (2008) Seasonal
22 drought stress in the Amazon: Reconciling models and observations. J. Geophys. Res.,
23 113, G00B01, doi:10.1029/2007JG000644.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

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.

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.

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

1 Cook ER, Woodhouse CA, Eakin CM, Meko DM, Stahle DW (2004) Long-term aridity
2 changes in the western United States. *Science*, 306, 1015-1018.
3

4 Cook ER, Seager R, Cane MA, Stahle DW (2007) North American droughts:
5 reconstructions, causes and consequences. *Earth-Sci. Rev.*, 81, 93-134.
6

7 Dai A, Trenberth KE, Qian T (2004) A global data set of Palmer Drought Severity Index
8 for 1870-2002: Relationship with soil moisture and effects of surface warming. *J.*
9 *Hydrometeorol.*, 5, 1117-1130.
10

11 Delpierre N, Soudani K, Francois C et al. (2009) Exceptional carbon uptake in European
12 forests during the warm spring of 2007: a data–model analysis. *Glob. Change Biol.*, 15,
13 1455-1474.
14

15 Friedlingstein P, Cox PM, Betts R et al. (2006) Climate-carbon cycle feedback analysis,
16 results from the C⁴MIP model intercomparison. *Journal of Climate*, 19, 3337-3353.
17

18 Friend AD, Arneeth A, Kiang NY et al. (2007) FLUXNET and modelling the global
19 carbon cycle. *Glob. Change Biol.*, 13, 610-633.
20

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

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

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.

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

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

Higuchi K, Shashkov A, Chan D, Saigusa N, Murayama S, Yamamoto S, Kondo H, Chen JM, Liu J, Chen B (2005) Simulations of seasonal and inter-annual variability of net CO₂ flux at Takayama with BEPS Ecosystem Model. *Agricultural and Forest Meteorology*, 134, 143-150.

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.

1 Huntington TG (2006) Evidence for intensification of the global water cycle: Review and
2 synthesis. *J. Hydrol.*, 319, 83-95.
3
4 Hwang Y, Carbone GJ (2009) Ensemble forecasts of drought indices using a conditional
5 residual resampling technique. *Journal of Applied Meteorology and Climatology*, 48,
6 1289-1301.
7
8 Ichii K, Suzuki T, Kato T et al (2009) Multi-model analysis of terrestrial carbon cycles in
9 Japan: reducing uncertainties in model outputs among different terrestrial biosphere
10 models using flux observations, *Biogeosciences Discuss.*, 6, 8455-8502.
11
12 Ito A (2008) The regional carbon budget of East Asia simulated with a terrestrial
13 ecosystem model and validated using AsiaFlux data. *Agr. Forest Meteorol.*, 148, 738-
14 747.
15
16 Jung M, Le Maire G, Zaehle S et al. (2007a) Assessing the ability of three land
17 ecosystem models to simulate gross carbon uptake of forests from boreal to
18 Mediterranean climate in Europe. *Biogeosciences*, 4, 647-656.
19
20 Jung M, Vetter M, Herold M et al. (2007b) Uncertainties of modeling gross primary
21 productivity over Europe: A systematic study on the effects of using different drivers and
22 terrestrial biosphere models. *Global Biogeochem. Cycles*, 21, GB4021,
23 doi:10.1029/2006GB002915.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

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

King AW, Post WM, Wullschleger SD, Ricciuto DM (2009) Sensitivity of modeled net ecosystem exchange to uncertainty in initial ecosystem carbon stocks (in preparation).

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.

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.

1 Liu S, Bliss N, Sundquist E, Huntington TG (2003) Modeling carbon dynamics in
2 vegetation and soil under the impact of soil erosion and deposition. Global
3 Biogeochemical Cycles, 17, doi:10.1029/2002GB002010.
4
5 Lokupitiya E, Denning S, Paustian K, Baker I, Schaefer K, Verma S, Meyers T.
6 Bernacchi CJ, Suyker A, Fischer M (2009) Incorporation of crop phenology in Simple
7 Biosphere Model (SiBcrop) to improve land-atmosphere carbon exchanges from
8 croplands. Biogeosciences, 6, 969-986.
9
10 McKee TB, Doeskin NJ, Kleist J (1993) The relationship of drought frequency and
11 duration to time scales. Proc. 8th Conf. on Applied Climatology, January 17-22, 1993,
12 American Meteorological Society, Boston, Massachusetts, 179-184.
13
14 Medlyn BE, Robinson AP, Clement R, McMurtrie RE (2005) On the validation of
15 models of forest CO₂ exchange using eddy covariance data: some perils and pitfalls. Tree
16 Physiology, 25, 839-857.
17
18 Medvigy D, Wofsy SC, Munger JW, Hollinger DY, Moorcroft PR (2009) Mechanistic
19 scaling of ecosystem function and dynamics in space and time: Ecosystem Demography
20 model version 2. J. Geophys. Res., 114, G01002, doi:10.1029/2008JG000812.
21
22 Meehl GA, Tebaldi C (2004) More intense, more frequent, and longer lasting Heat waves
23 in the 21st Century. Science, 305, 994-997.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

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.

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.

Piao SL, Friedlingstein P, Peylin P, Reichstein M, Luyssaert 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.

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

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

Richardson AD and 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.

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.

1 Schwalm CR, Williams CA, Schaefer KS et al. (2009) Assimilation exceeds respiration
2 sensitivity to drought: A FLUXNET synthesis. *Global Change Biology*, doi:
3 10.1111/j.1365-2486.2009.01991.x
4

5 Seager R (2007a) The turn of the century drought across North America: global context,
6 dynamics and past analogues. *Journal of Climate*, 20, 5527-5552.
7

8 Seager R, Ting MF, Held IM (2007b) Model projections of an imminent transition to a
9 more arid climate in Southwestern North America, *Science*, 316, 1181-1184.
10

11 Seager R, Tzanova A, Nakamura J (2009) Drought in the Southeastern United States:
12 Causes, variability over the last millennium and the potential for future hydroclimate
13 change, *Journal of Climate*, 22, 5021-5045.
14

15 Sheffield J, Wood EF (2008) Projected changes in drought occurrence under future
16 global warming from multi-model, multi-scenario, IPCC AR4 simulations. *Climate*
17 *Dynamics*, 13, 79-105.
18

19 Siqueira MB, Katul GG, Sampson DA, Stoy PC, Juang J-Y, McCarthy HR, Oren R
20 (2006) Multiscale model intercomparisons of CO₂ and H₂O exchange rates in a maturing
21 southeastern US pine forest. *Global Change Biology*, 12, 1189-1207.
22

1 Sitch S, Smith B, Prentice IC, Arneth A, Bondeau A, Cramer W, Kaplan JO, Levis S,
2 Lucht W, Sykes MT, Thonicke K, Venevsky S (2003) Evaluation of ecosystem
3 dynamics, plant geography and terrestrial carbon cycling in the LPJ dynamic global
4 vegetation model. *Global Change Biology*, 9, 161-185.
5
6 Sitch S, Huntingford C, Gedney N et al. (2008) Evaluation of the terrestrial carbon cycle
7 future plant geography and climate-carbon cycle feedbacks using five Dynamic Global
8 Vegetation Models (DGVMs). *Glob. Change Biol*, 14, 2015-2039.
9
10 Soja AJ, Tchebakova NM, French NHF, Flannigan MD, Shugart HH, Stocks BJ,
11 Sukhinin AI, Varfenova EI, Chapin FS, Stackhouse PW (2007). Climate induced boreal
12 forest change: Predictions versus current observations. *Global and Planetary Change*, 56,
13 274-296.
14
15 St-Hilaire F, Wu J, Roulet NT, Froking S, Lafleur PM, Humphreys ER, Arora V
16 (2008) McGill Wetland Model: evaluation of a peatland carbon simulator developed for
17 global assessments. *Biogeosciences Discussions*, 5, 1689-1725.
18
19 Stoy P, Katul G, Siqueira M, Juang J, McCarthy H, Kim H, Oishi A, Oren R (2005)
20 Variability in net ecosystem exchange from hourly to inter-annual time scales at adjacent
21 pine and hardwood forests: a wavelet analysis. *Tree Physiology*, 25, 887-902.
22

1 Stoy P, Richardson A, Baldocchi D, Katul G, Stanovick J, Mahecha M, Reichstein M,
2 Detto M, Law B, Wohlfahrt G, Arriga N, Campos J, McCaughey J, Montagnani L, Paw
3 U, KT, Sevanto S, Williams M (2009) Biosphere-atmosphere exchange of CO₂ in relation
4 to climate: a cross-biome analysis across multiple time scales. *Biogeosciences*, 6, 2297-
5 2312.

6

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

9

10 Thomas CK, Law BE, Irvine J, Martin JG, Pettijohn JC, Davis KJ (2009) Seasonal
11 hydrology explains interannual and seasonal variation in carbon and water exchange in a
12 semi-arid mature ponderosa pine forest in Central Oregon. *JGR Biogeosciences*. In press.

13

14 Thornton, PE, Running SW, White MA (1997) Generating surfaces of daily
15 meteorological variables over large regions of complex terrain. *J. of Hydrology*, 3-4, 214-
16 251

17

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

22

1 Tian, HQ, Chen G, Liu M, Zhang C, Sun G, Lu C, Xu X, Ren W, Pan P, Chappelka A
2 (2009) Model Estimates of Ecosystem Net Primary Productivity, Evapotranspiration, and
3 Water Use Efficiency in the Southern United States during 1895-2007. *Forest Ecology*
4 *and Management* (in press).

5
6 Trenberth KE, Jones PD, Ambenje P et al. (2007) Observations: Surface and atmospheric
7 climate change. In: *Climate Change 2007: The Physical Science Basis. Contribution of*
8 *Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on*
9 *Climate Change* (eds Solomon SD et al.). Cambridge University Press, Cambridge,
10 United Kingdom and New York, NY, USA.

11
12 Waring RH, Running SW (2007) *Forest ecosystems - analysis at multiple scales* (third
13 edition). Academic Press, San Diego, 440 pp.

14
15 Weng E, Luo Y (2008) Soil hydrological properties regulate grassland ecosystem
16 responses to multifactor global change: A modeling analysis. *J. Geophys. Res.*, 113,
17 G03003, doi:10.1029/2007JG000539.

18
19 Williams CA, Hanan NP, Scholes RJ, Kutsch WL (2009) Complexity in water and
20 carbon dioxide fluxes following rain pulses in an African savanna. *Oecologia*, 161, 469-
21 480.

- 1 Williamson TB, Price DT, Beverly JL, Bothwell PM, Frenkel B, Park J, Patriquin MN
2 (2008) Assessing potential biophysical and socioeconomic impacts of climate change on
3 forest-based communities: a methodological case study. Nat. Resour. Can., Can. For.
4 Serv., North. For. Cent., Edmonton, AB. Inf. Rep. NOR-X-415E.
5
6 Zhan XW, Xue YK, Collatz GJ (2003) An analytical approach for estimating CO₂ and
7 heat fluxes over the Amazonian region. Ecol. Model. 162, 97-117.
8
9 Zhou XL, Peng CH, Dang QL, Sun JF, Wu HB, Hua D (2008) Simulating carbon
10 exchange in Canadian Boreal forests I: model structure, validation, and sensitivity
11 analysis. Ecological Modelling, 219, 287-299.

Tables & Figures

Table 1. Models evaluated with months and sites sampled. Difference in sample sizes results from the use of standard errors of observed fluxes and site-specific drought metrics for predictive skill. Taylor diagrams included all available flux data.

Model	Timestep	Predictive skill		Taylor diagrams		Source
		Months	Sites	Months	Sites	
AgroIBIS [*]	half-hourly	165	4	192	5	Kucharik & Twine (2007)
BEPS	daily	860	10	945	10	Higuchi et al. (2004)
Biome-BGC	daily	1536	26	2001	36	Thornton et al. (2005)
Can-IBIS	half-hourly	1587	23	1978	27	Williamson et al. (2008)
CN-CLASS	half-hourly	1696	25	2082	31	Arain et al. (2006)
DLEM [%]	daily	1831	28	2246	33	Tian et al. (2009)
DNDC [*]	daily	122	4	192	5	Li et al. (2009)
Ecosys	hourly	1693	26	2450	39	Grant et al. (2005)
ED2	half-hourly	1430	22	1684	25	Medvigy et al. (2009)
EDCM	monthly	582	10	658	10	Liu et al. (2003)
EPIC ^{*,%}	daily	47	1	48	1	Causarano et al. (2007)
ISOLSM [%]	half-hourly	824	9	909	9	Riley et al. (2002)
LoTEC ^{\$}	hourly	759	10	825	10	King et al. (2009)
LPJ	daily	1711	25	2126	29	Sitch et al. (2003)
MEAN [#]	native	1940	30	2776	44	this study
ORCHIDEE	half-hourly	1873	30	2332	35	Krinner et al. (2005)
SiB3 [%]	half-hourly	1825	27	2258	31	Baker et al. (2008)
SiBCASA [%]	half-hourly	1952	30	2402	35	Schaefer et al. (2009)
SiBcrop [*]	half-hourly	165	4	192	5	Lokupitiya et al. (2009)
SSiB2 [%]	half-hourly	1952	30	2800	44	Zhan et al. (2003)
TECO	hourly	1952	30	2414	35	Weng & Luo (2008)
TRIPLEX-Flux	daily	246	7	291	7	Zhou et al. (2008)

^{*} Used only on croplands

[%] Used remotely sensed inputs of leaf area index, greenness, or fraction of absorbed fraction of photosynthetically active radiation, i.e., phenology was calculated in diagnostic mode using an average seasonal cycle by site for the full simulation period.

^{\$} Tuned using data assimilation

[#] Mean value across all modeled values

1 Table 2. Name, measurement period, and biome for all sites. Italicized site entries only
2 used in Taylor diagrams. CA sites located in Canada; US sites in the continental United
3 States except US-Atq and US-Brw (Alaska).

Site	Full Name	Period	Biome [#]
CA-Ca1	Campbell River – Mature Douglas-fir	1998-2006	ENFT
<i>CA-Ca2</i>	<i>Campbell River – Douglas-fir clearcut</i>	<i>2001-2006</i>	<i>ENFT</i>
<i>CA-Ca3</i>	<i>Campbell River – Douglas-fir juvenile</i>	<i>2002-2006</i>	<i>ENFT</i>
CA-Gro	Groundhog River Station	2004-2006	MF
CA-Let	Lethbridge Grassland	1997-2006	GRA
CA-Mer	Eastern Peatland – Mer Bleue	1999-2006	WET
CA-Oas	BERMS – Old Aspen	1997-2006	DBF
CA-Obs	BERMS – Old Black Spruce	2000-2006	ENFB
CA-Ojp	BERMS – Old Jack Pine	2000-2006	ENFB
CA-Qfo	Quebec – Mature Black Spruce	2004-2006	ENFB
<i>CA-SJ1</i>	<i>BERMS – Jack Pine, 1994 harvest</i>	<i>2002-2005</i>	<i>ENFB</i>
<i>CA-SJ2</i>	<i>BERMS – Jack Pine, 2002 harvest</i>	<i>2003-2006</i>	<i>ENFB</i>
<i>CA-SJ3</i>	<i>BERMS – Jack Pine, 1975 harvest</i>	<i>2004-2005</i>	<i>ENFB</i>
<i>CA-TP3</i>	<i>Turkey Point – Middle-aged^{\$}</i>	<i>2003-2007</i>	<i>ENFT</i>
CA-TP4	Turkey Point – Mature	2002-2007	ENFT
CA-WP1	Western Peatland – LaBiche River	2003-2007	WET
US-ARM	ARM – Southern Great Plains	2000-2006	CRO
<i>US-Atq</i>	<i>Atqasuk^{\$}</i>	<i>1999-2006</i>	<i>TUN</i>
<i>US-Brw</i>	<i>Barrow^{\$}</i>	<i>1998-2006</i>	<i>TUN</i>
<i>US-Dk2</i>	<i>Duke Forest – Hardwood^{\$}</i>	<i>2003-2005</i>	<i>DBF</i>
US-Dk3	Duke Forest – Loblolly Pine	1998-2005	ENFT
US-Ha1	Harvard Forest – EMS Tower	1991-2006	DBF
US-Ho1	Howland Forest – Main Tower	1996-2004	ENFT
<i>US-IB1</i>	<i>Fermi Lab – Maize/soybean rotation</i>	<i>2005-2007</i>	<i>CRO</i>
US-IB2	Fermi Lab – Prairie	2004-2007	GRA
US-Los	Lost Creek	2000-2006	WET
US-Me2	Metolius – Intermediate-aged Ponderosa Pine	2002-2007	ENFT
<i>US-Me3</i>	<i>Metolius – Ponderosa Pine, young #2</i>	<i>2004-2005</i>	<i>ENFT</i>
<i>US-Me4</i>	<i>Metolius – Ponderosa Pine, old-growth^{\$}</i>	<i>1996-2000</i>	<i>ENFT</i>
<i>US-Me5</i>	<i>Metolius – Ponderosa Pine, Young #1</i>	<i>1999-2002</i>	<i>ENFT</i>
US-MMS	Morgan Monroe State Forest	1999-2006	DBF
US-MOz	Missouri Ozark	2004-2007	DBF
US-Ne1	Mead – Irrigated maize	2001-2006	CRO
US-Ne2	Mead – Irrigated maize/soybean	2001-2006	CRO
US-Ne3	Mead – Rainfed maize/soybean	2001-2006	CRO
<i>US-NR1</i>	<i>Niwot Ridge</i>	<i>1998-2007</i>	<i>ENFT</i>
US-PFa	Park Falls / WLEF	1997-2005	MF

<i>US-Shd</i>	<i>Shidler</i>	<i>1997-2001</i>	<i>GRA</i>
US-SO2	Sky Oaks – Old	1998-2006	SHR
US-Syv	Sylvania Wilderness Area	2001-2006	MF
US-Ton	Tonzi Ranch	2001-2007	WSA
US-UMB	University of Michigan Biological Station (UMBS)	1998-2006	DBF
US-Var	Vaira Ranch	2001-2007	GRA
US-WCr	Willow Creek	1998-2006	DBF

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

^{\$} Sites used alternative post-processing protocol based on the La Thuile and Asilomar FLUXNET Synthesis dataset (<http://www.fluxdata.org/>; Moffat et al., 2007; Papale et al., 2006).

1 Table 3. Model-specific and site history predictants, in addition to biome, climatic
2 season, and drought level, used to classify the presence/absence of predictive skill.

Predictant*	Value type	Comment
Model temporal resolution	Factor	Value: daily, hourly, half-hourly
Prognostic vs. diagnostic LAI	Boolean	How was the canopy, i.e., LAI (leaf area index), calculated? Does model require LAI as an input? Note: A subset of models used remotely sensed LAI to calculate canopy biomass or for phenology (diagnostic mode). This represents partial assimilation of site data into model runs inasmuch as drought is reflected in the remotely sensed products used.
Number of vegetation pools	Numeric	Number of pools, both dynamic and static
Number of soil pools	Numeric	Number of pools, both dynamic and static
Number of soil layers	Numeric	Number of layers
Algorithm: Autotrophic respiration	Factor	Value: assumed fraction of instantaneous GPP, explicitly calculated, How is autotrophic respiration calculated
Algorithm: Canopy leaf biomass	Factor	Value: explicitly calculated, GPP/NPP fraction, LAI, nil How is canopy biomass calculated?
Algorithm: Ecosystem respiration	Factor	Value: autotrophic respiration + heterotrophic respiration, explicitly calculated, forced annual balance How is ecosystem respiration calculated?
Algorithm: Gross primary productivity	Factor	Value: enzyme kinetic model, light use efficiency model, nil, stomatal conductance model How is gross primary productivity calculated?
Algorithm: Heterotrophic respiration	Factor	Value: explicitly calculated, first or greater order model, zero-order model How is heterotrophic respiration calculated?
Algorithm: Net ecosystem productivity	Factor	Value: explicitly calculated, gross primary productivity - ecosystem respiration, net primary production - heterotrophic respiration How is net ecosystem productivity calculated?
Overall model complexity	Factor	Value: minimal, low, average, high, highest Total amount of first-order functional arguments for the following model-generated variables/outputs: Autotrophic respiration, canopy leaf biomass, ecosystem respiration, evapotranspiration, gross primary productivity, heterotrophic respiration, net ecosystem productivity, net primary production, soil moisture

Grazed	Boolean	Note: Factor values correspond to quintiles. Has listed management activity or disturbance or event occurred on site?
Fertilized		
Fire		
Harvest		
Herbicide		
Insects and pathogens		
Irrigation		
Natural regeneration		
Pesticide		
Planted		
Residue management		
Thinning		
Underburn		
Stand age class	Factor	Value: 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.

* Modeling architecture predictants were taken from the Metadata for Forward (Ecosystem) Model Intercomparison survey data collated by the NACP Site Synthesis (http://daac.ornl.gov/SURVEY8/survey_results.shtml). Of the 21 distinct modeling platforms used (excluding the mean model ensemble) 15 models contributed metadata. The predictive skill classification exercise excluded runs from: AgroIBIS, EDCM, LoTEC, SiBCASA, SSiB2, and TRIPLEX-Flux.

- 1 Table 4. Predictive skill and sample size by drought level and climatic season. Drought
- 2 level was based on monthly values of 3-month SPI, Dry required a Standard Precipitation
- 3 Index (SPI) value of < -0.8 ; Wet $> +0.8$. Otherwise Normal conditions existed. Sample
- 4 size is the total number of paired data-model site-months.

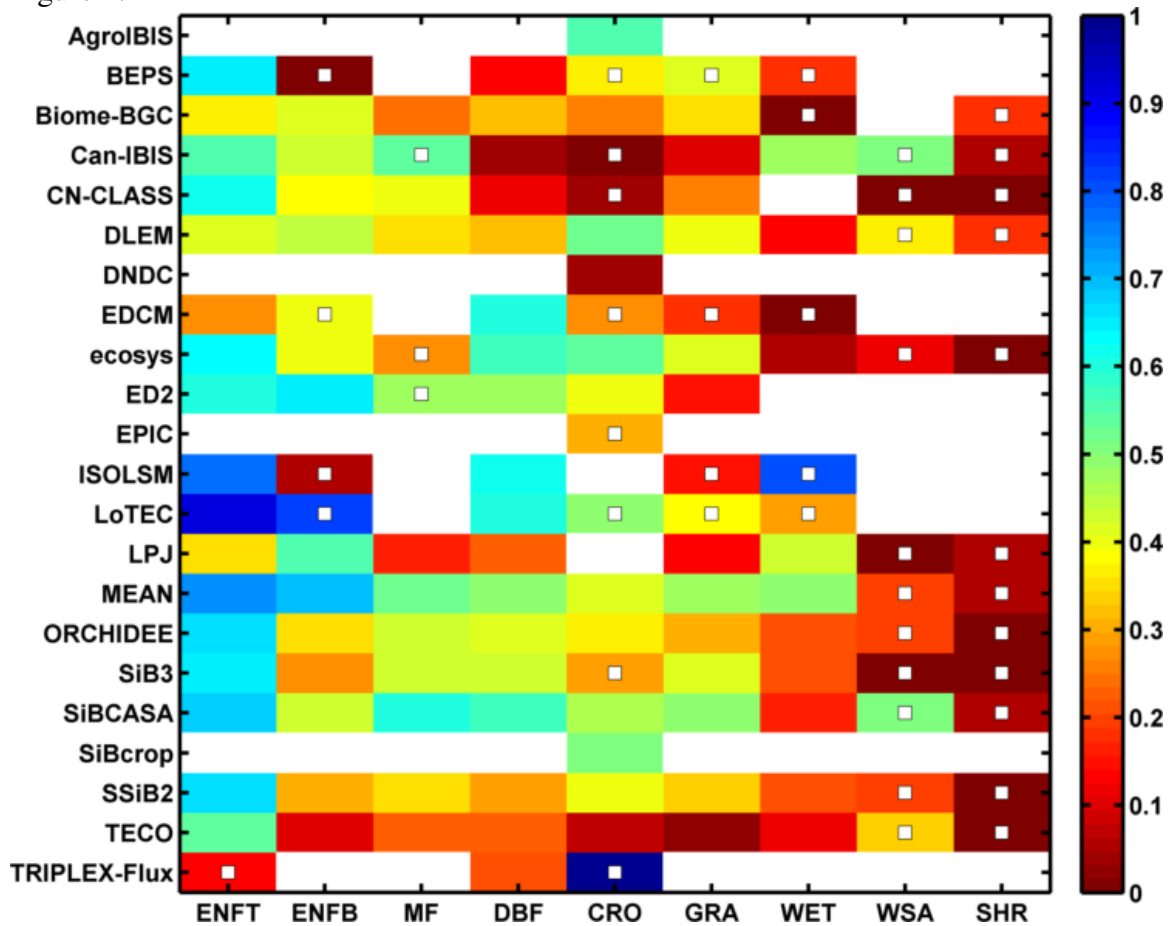
	Predictive skill				Sample size			
	Dry	Normal	Wet	Overall	Dry	Normal	Wet	Overall
Winter	0.42	0.49	0.47	0.47	1479	3996	1096	6571
Spring	0.27	0.35	0.22	0.31	1534	3990	1209	6733
Summer	0.45	0.54	0.42	0.49	1359	3910	1442	6711
Fall	0.28	0.30	0.33	0.30	1491	3936	1305	6732
Overall	0.35	0.42	0.36	0.39	5863	15832	5052	26747

- 1 Table 5. Predictive skill by climatic season, drought level, and biome. Drought level was
- 2 based on monthly values of 3-month SPI, Dry required a Standard Precipitation Index
- 3 (SPI) value of < -0.8 ; Wet $> +0.8$. Otherwise Normal conditions existed.

Biome [#]	Climatic season				Drought level			Overall
	Winter	Spring	Summer	Fall	Dry	Normal	Wet	
CRO	0.50	0.43	0.52	0.06	0.45	0.34	0.22	0.37
DBF	0.57	0.15	0.52	0.21	0.23	0.39	0.36	0.36
ENFB	0.32	0.49	0.58	0.23	0.30	0.37	0.52	0.41
ENFT	0.49	0.50	0.57	0.79	0.52	0.63	0.51	0.59
GRA	0.45	0.34	0.32	0.06	0.23	0.37	0.13	0.30
MF	0.61	0.07	0.68	0.16	0.35	0.42	0.30	0.38
SHR	0.02	0.16	0	0	0.04	0.04	0.05	0.04
WET	0.40	0.12	0.42	0.11	0.23	0.29	0.20	0.26
WSA	0.07	0.49	0.27	0.05	0.14	0.25	0.18	0.22
Overall	0.47	0.31	0.49	0.30	0.35	0.42	0.36	0.39

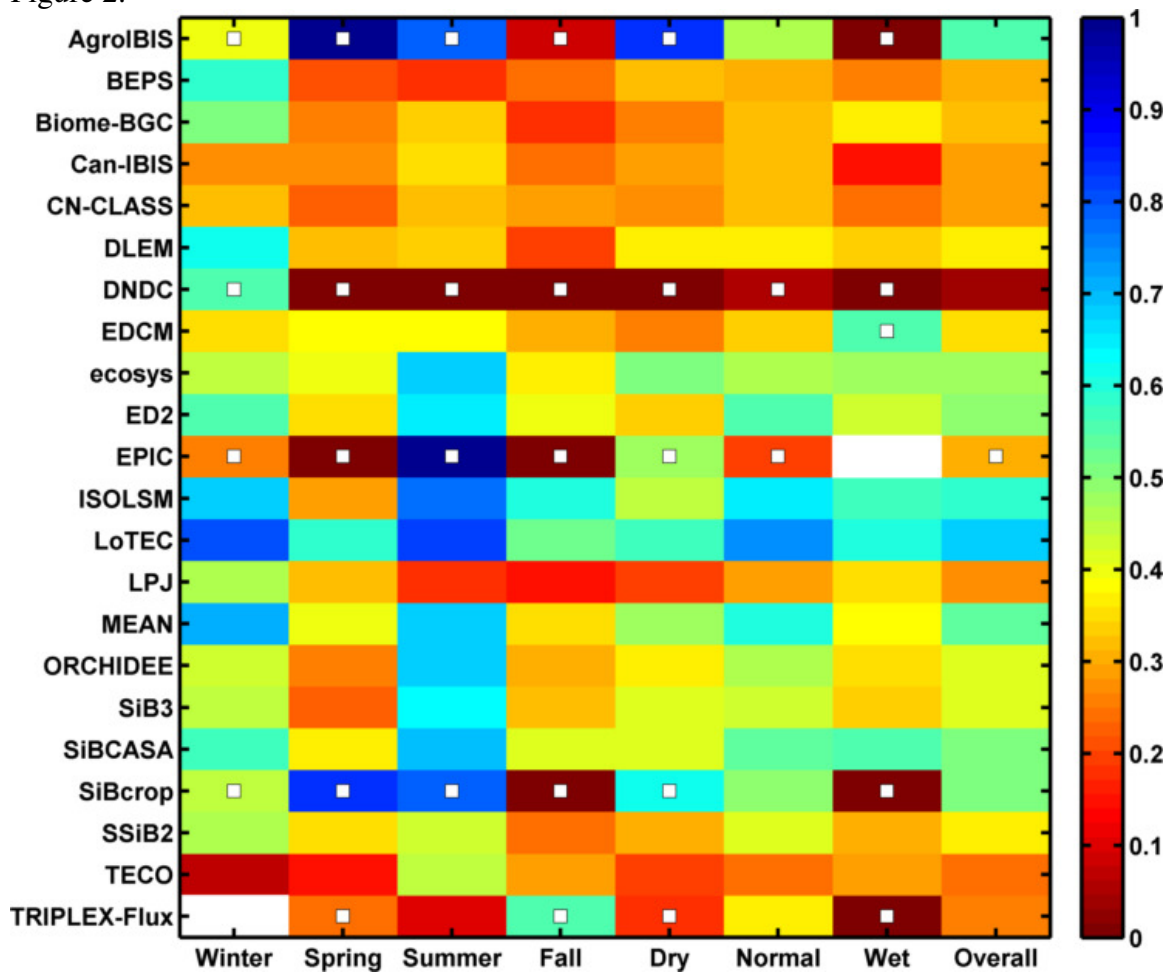
[#] Biome codes: CRO = cropland, GRA = grassland, ENFB = evergreen needleleaf forest – boreal climatic zone, ENFT = evergreen needleleaf forest – temperate climatic zone, DBF = deciduous broadleaf forest, MF = mixed (deciduous/evergreen) forest, WSA = woody savanna, SHR = shrubland, WET = wetland.

1 Figure 1.



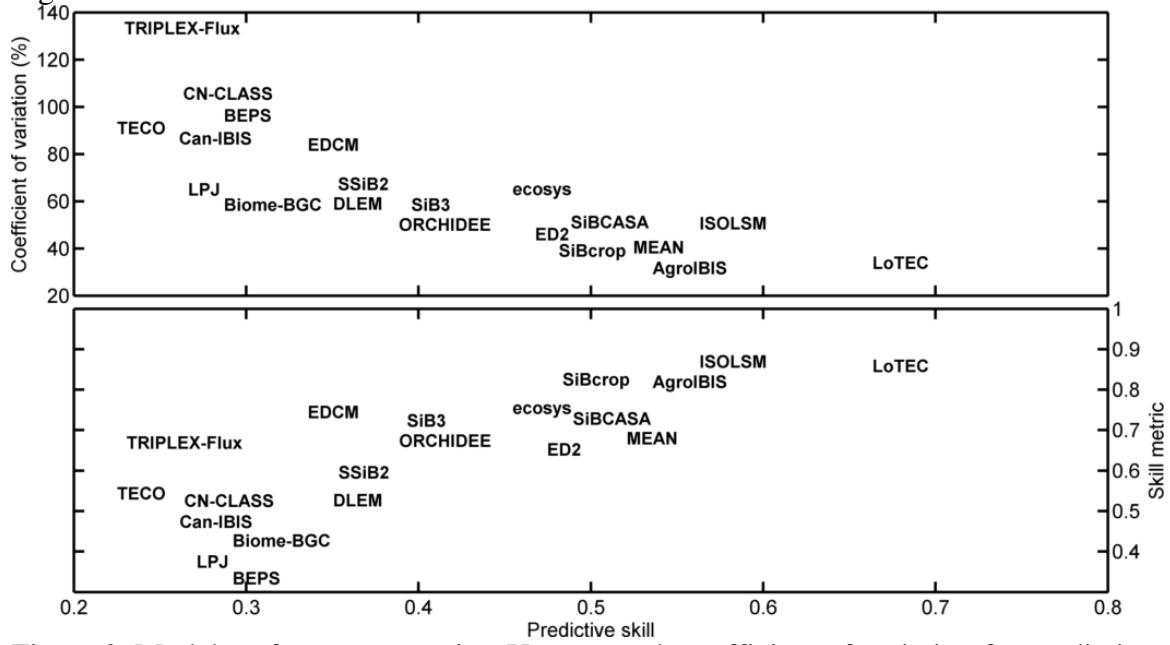
2
3 Figure 1. Summary of model predictive skill by biome for each model. Biomes ordered in
4 descending order based on across-model average predictive skill. White blocks: no
5 observations; gray squares: undersampled ($n < 100$ months).

1 Figure 2.



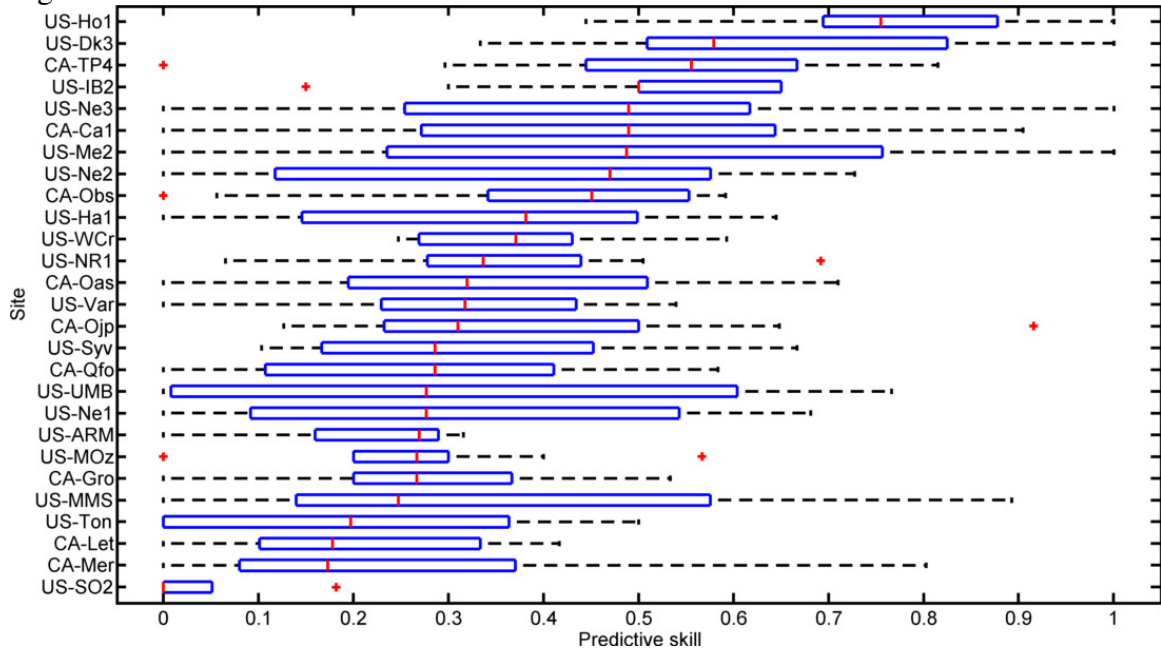
2
3 Figure 2. Summary of model predictive skill by climatic season, drought level, and
4 overall. White blocks: no observations; gray squares: undersampled ($n < 100$ months).

1 Figure 3.



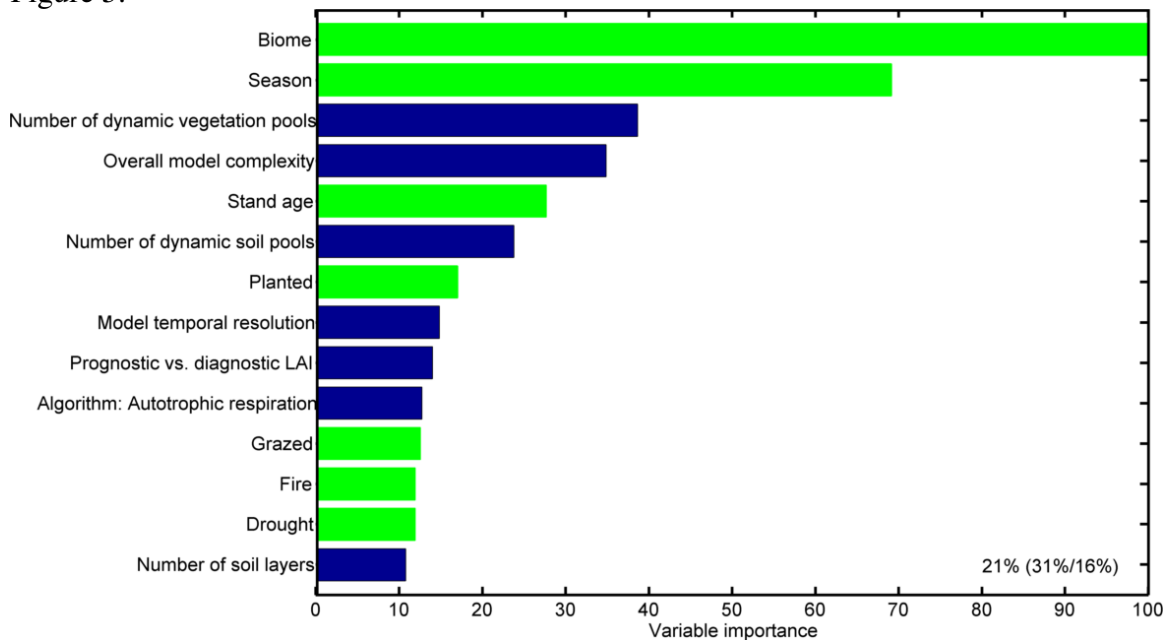
2
3 Figure 3. Model performance metrics. Upper panel: coefficient of variation for predictive
4 skill (CV_π) vs. predictive skill (π). Lower panel: overall model skill (S) from Taylor
5 diagrams vs. predictive skill. DNDC ($S = 0.34$, $\pi = 0.04$, $CV_\pi = 174\%$; outlier) and EPIC
6 ($S = 0.95$, $\pi = 0.30$, $CV_\pi = \text{undefined}$; only one site simulated) not shown. Model names
7 jittered to improve readability.

1 Figure 4.



2
3 Figure 4. Boxplot of model predictive skill by site. Panels show interquartile range (blue
4 box), median (solid red line), range (whiskers), and outliers (red cross; values more than
5 1.5 x interquartile range from the median). Only sites ($n = 27$) simulated with at least 10
6 unique models using steady state spinup shown. Sites sorted by median predictive skill.

1 Figure 5.



2
3 Figure 5. Variable importance scores for model-specific (blue) and site-specific (green)
4 predictants. Scores were generated from a regression tree with the presence/absence of
5 predictive skill by data group ($n = 2258$) as the response. *Planted*, *Grazed*, and *Fire* were
6 boolean variables related to the occurrence of the event on site. *Stand age* was young,
7 intermediate, or mature for forested stands; nil otherwise. For model-specific predictants
8 *Algorithm: Autotrophic respiration* was either assumed fraction of instantaneous gross
9 primary productivity, explicitly calculated, or nil (not calculated by the model). Numbers
10 in lower right corner indicate overall misclassification rate and, parenthetically, within
11 class misclassification rates for data groups with then without predictive skill. Only the
12 14 of 29 predictants with score > 10 shown; see Table 3 for complete listing of evaluated
13 model structural and site attributes.

1 Figure 6.

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

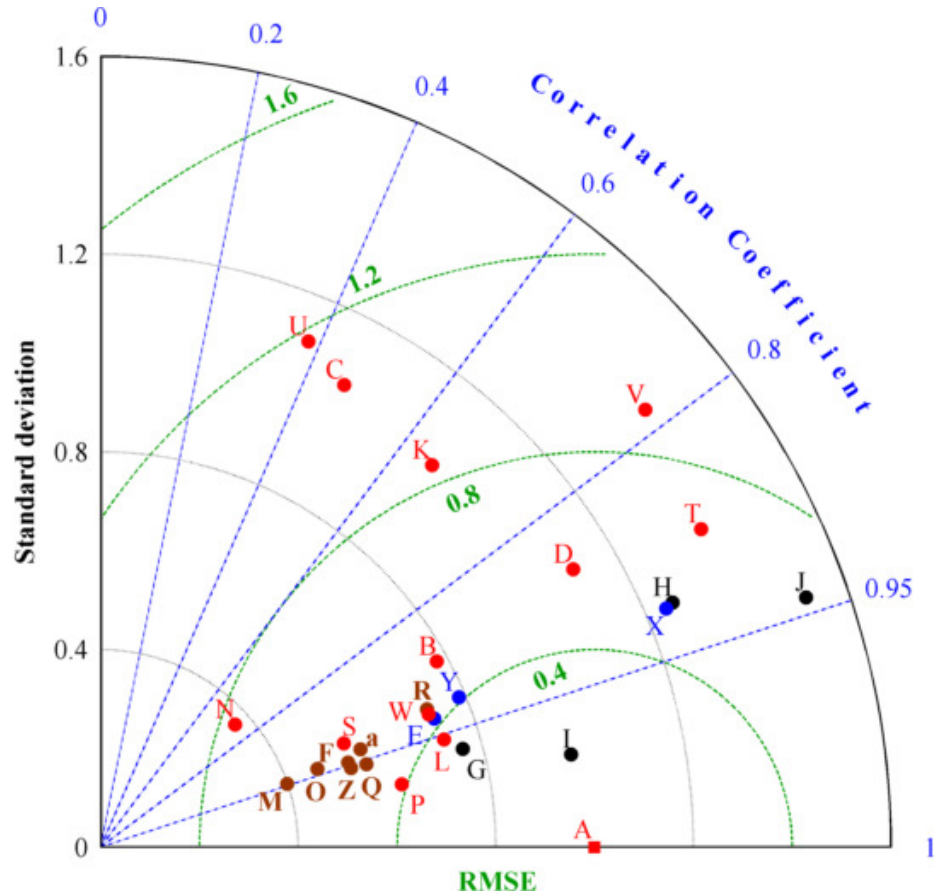
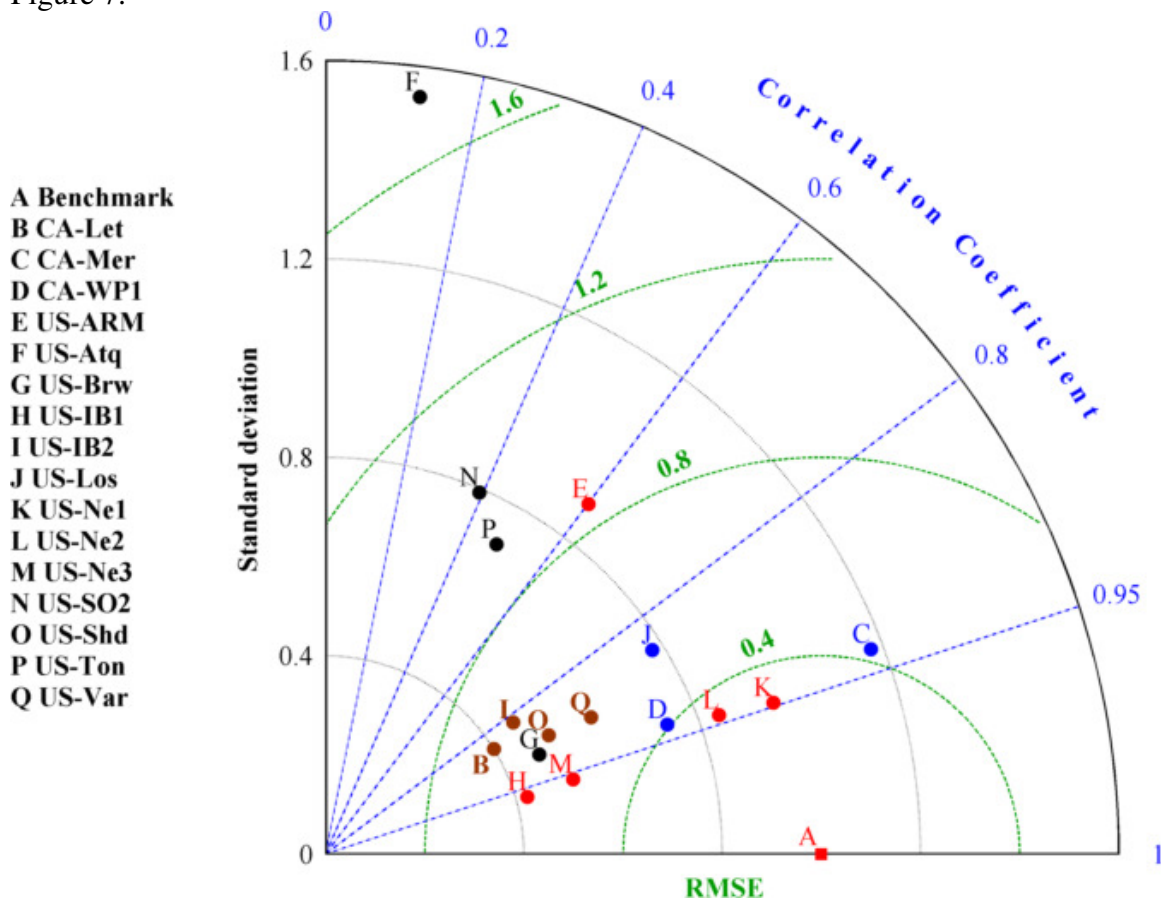


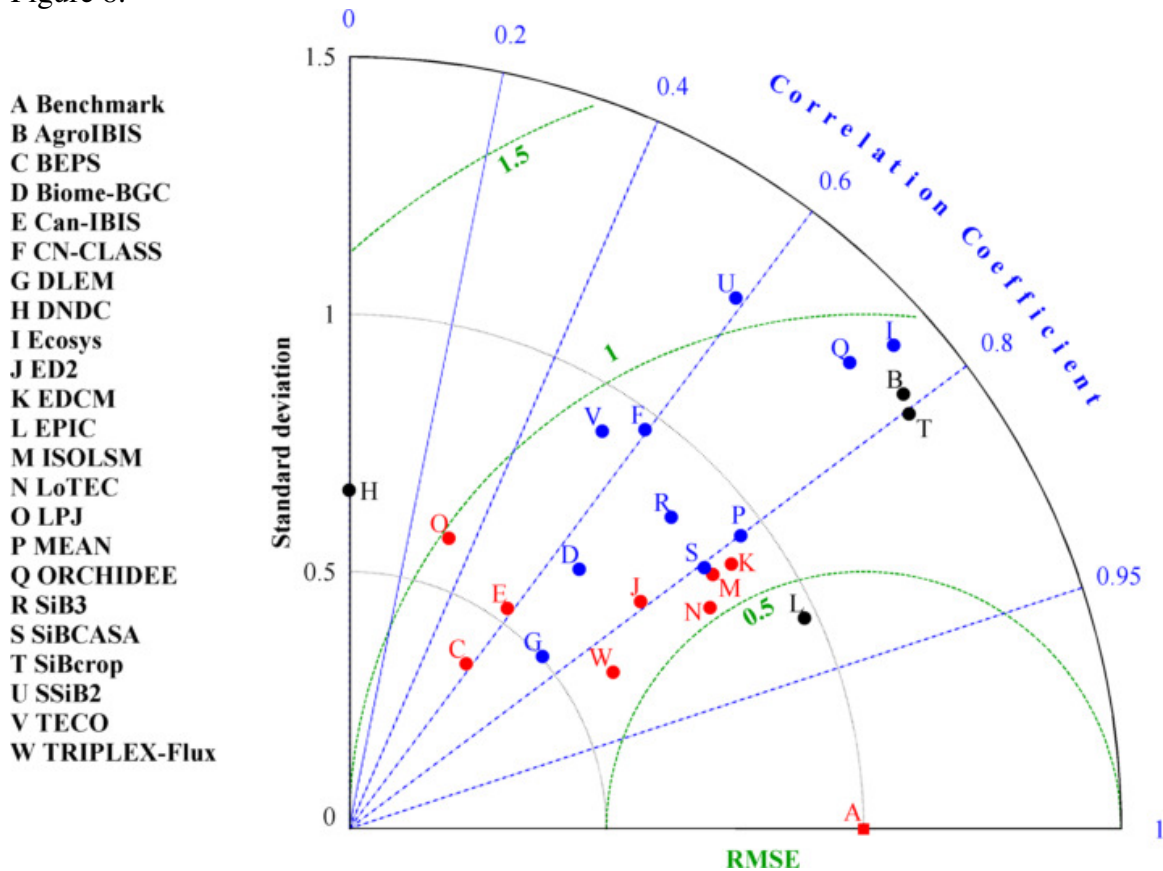
Figure 6. Taylor diagram of normalized mean model performance for forested sites. Each circle ($n = 26$ sites) is the site-specific mean model ensemble (MEAN). Benchmark (red square) corresponds to observed normalized monthly NEP; units of σ and RMSE are multiples of observed σ . Color coding of site letter and circles indicates biome: evergreen needleleaf forest – temperate climatic zone (red), deciduous broadleaf forest (brown), mixed (deciduous/evergreen) forest (blue), evergreen needleleaf forest – boreal climatic zone (black). Outlying sites (evergreen needleleaf forest – boreal climatic zone) not shown: CA-SJ1 ($\rho = 0.81$, $\sigma = 3.9$, RMSE = 3.1) and CA-SJ2 ($\rho = -0.67$, $\sigma = 4.3$, RMSE = 5.1).

1 Figure 7.



2
3 Figure 7. Taylor diagram of normalized mean model performance for non-forested sites.
4 Each circle ($n = 16$ sites) is the site-specific mean model ensemble (MEAN). Benchmark
5 (red square) corresponds to observed normalized monthly NEP; units of σ and RMSE are
6 multiples of observed σ . Color coding of site letter and circles indicates biome: croplands
7 (red), grasslands (brown), wetlands (blue), all other biomes (black).

1 Figure 8.



2
3 Figure 8. Taylor diagram of normalized across-site average model performance. Model σ
4 and RMSE were normalized by observed σ . Each circle ($n = 22$ models) corresponds to
5 the mean across all sites. Benchmark (red square) corresponds to observed normalized
6 monthly NEP; units of σ and RMSE are multiples of observed σ . Color coding of model
7 letter and circles indicates generality of model performance: specialist models used only
8 in croplands ($n \leq 5$ sites; black), generalist models used across range of biomes and sites
9 ($n \geq 30$ sites, blue), all other models (red). The correlation for DNDC ($\rho = -0.13$) is
10 displayed as zero for readability.

11
12 [End of document]