

¹ Title: **Facilitating feedbacks between field**
² **measurements and ecosystem models**

³ Running Title: **Feedbacks between measurements and models**

⁴ David S. LeBauer^{1,*} Dan Wang¹ Katherine T. Richter² Carl C. Davidson² Michael C.
⁵ Dietze^{1,2}

⁶ ¹ Energy Biosciences Institute, ² Department of Plant Biology, University of Illinois,
⁷ Urbana, Illinois, USA

⁸ * Corresponding author

⁹ Telephone: (949) 433-7410

¹⁰ Fax: (217) 244-3637

¹¹ E-mail: dlebauer@illinois.edu

Abstract

Ecological models help us understand how ecosystems function, predict responses to global change, and identify future research needs. However, widespread use of models is limited by the technical challenges of model-data synthesis and information management. To address these challenges, we present a ecoinformatic workflow, the Predictive Ecosystem Analyzer (PEcAn), that facilitates model analysis. Herein we describe the PEcAn modules that synthesize plant trait data to estimate model parameters, propagate parameter uncertainties through to model output, and evaluate the contribution of each parameter to model uncertainty. We illustrate a comprehensive approach to the estimation of parameter values, starting with a statement of prior knowledge that is refined by species level data using Bayesian meta-analysis; this is the first use of a rigorous meta-analysis to inform the parameters of a mechanistic ecosystem model. Parameter uncertainty is propagated using ensemble methods to estimate model uncertainty. Variance decomposition allows us to quantify the contribution of each parameter to model uncertainty; this information can be used to prioritize subsequent data collection. By streamlining the use of models and focusing efforts to identify and constrain the dominant sources of uncertainty in model output, the approach used by PEcAn can speed scientific progress.

We demonstrate PEcAn's ability to incorporate data to reduce uncertainty in productivity of a perennial grass monoculture (*Panicum virgatum* L.) modeled by the Ecosystem Demography model. Prior estimates were specified for fifteen model parameters, and species-level data were available for seven of these. Meta-analysis of species-level data substantially reduced the contribution of three parameters (specific leaf area [SLA], maximum carboxylation rate, and stomatal slope) to overall model uncertainty. By contrast, root turnover rate, root respiration rate, and leaf width had

37 little effect on model output, therefore trait data had little impact on model uncertainty.
38 For fine root allocation the decrease in parameter uncertainty was offset by an increase in
39 model sensitivity. Remaining model uncertainty is driven by growth respiration, fine root
40 allocation, leaf turnover rater, and SLA. By establishing robust channels of feedback
41 between data collection and ecosystem modeling, PEcAn provides a framework for more
42 efficient and integrative science.

43 **keywords:** traits, ecoinformatics, ecophysiology, Ecosystem Demography, sensitivity
44 analysis, variance decomposition, ecological forecast, Bayesian, meta-analysis, ecosystem
45 model

46 Introduction

47 In the face of unprecedented global change there is growing demand for predictions of
48 ecosystem responses that provide actionable information for policy and management
49 (Clark et al., 2001). Currently, the response of the terrestrial biosphere remains one of
50 the largest sources of uncertainty in projections of climate change (Denman et al., 2007).
51 This uncertainty comes from a combination of the uncertainties about our conceptual
52 understanding of ecological systems, as captured by the structure and assumptions of the
53 models used to make ecological forecasts, the uncertainties in the parameters of these
54 models, and the uncertainties associated with the underlying data itself (McMahon et al.,
55 2009). Reducing these uncertainties requires that we be able to synthesize existing
56 information, efficiently identify the dominant sources of model uncertainty and target
57 them with further field research.

58 Despite the acknowledged importance of these activities, there is often a disconnect
59 between model simulation and data collection. Both model-data synthesis and the
60 investigation of uncertainty remain challenging, while the use of models to quantitatively
61 inform data collection is extremely rare. Most modeling uses a single point estimate for
62 each parameter, effectively treating each parameter value as completely certain. However,
63 such point estimates do not account for the degree to which we understand a parameter
64 based on observations. Furthermore, the rationale for a particular estimate is often
65 unclear, as is the degree to which the estimate represents the process being observed or
66 its representation in a model. In many cases, parameter values are chosen iteratively to
67 “tune”, or “calibrate” the model output to observations. A first step toward constraining
68 model uncertainty is to account for uncertainty in model parameters instead of relying on
69 point estimates.

70 More rigorous approaches to estimating parameter values include model optimization and

71 data assimilation (Medvigy et al., 2009; Reichstein et al., 2003), as well as Bayesian
72 model-data fusion (Luo et al., 2011). However, these approaches have generally started
73 with uninformative or vague prior estimates of model parameters. These vague priors
74 ignore available data that could directly inform parameter values; the most commonly
75 used vague prior distribution is a uniform. A uniform prior assigns equal probability to
76 parameter values over its entire range, in many cases over many orders of magnitude. The
77 use of such vague priors often exacerbates problems with equifinality (Richardson and
78 Hollinger, 2005; Williams et al., 2009; Luo et al., 2009) which can produce unidentifiable
79 parameters, as well as biologically unrealistic parameter sets that generate the right model
80 output for the wrong reasons (Beven and Freer, 2001; Beven, 2006; Williams et al., 2009).
81 Another reason to use informed priors is to take advantage of one of the key strengths of
82 the Bayesian paradigm: the ability to synthesize multiple sources of information in a
83 rigorous and consistent framework. For example, plant traits related to leaf stoichiometry
84 and photosynthetic capacity are often well constrained by previous research (Skillman,
85 2008; Reich and Oleksyn, 2004; Wright et al., 2004; Wullschleger, 1993), while other traits,
86 such as root respiration rate, are more difficult to measure and data are sparse. Informed
87 priors allow existing information to be formally integrated into model parameterization,
88 even if there is no data for the particular species or plant functional type (PFT) being
89 measured; the level of confidence in a parameter value is reflected in its variance.
90 Models have rarely been used to quantify the value of data with respect to reducing
91 uncertainty. Instead, data collection is often focused on answering specific questions in
92 specific spatial, temporal, and taxonomic contexts. In these contexts, the value of a
93 particular data set is based on the ability to answer a particular question. However, the
94 same data set may have a very different value in the context of reducing model
95 uncertainty. For example, a single data point used to inform a poorly understood but
96 influential model parameter can reduce model uncertainty more than a large collection of

data on a trait that is relatively well studied. In a modeling context, the value of an additional data point depends both on how much it constrains parameter uncertainty and the sensitivity of model output to the parameter. Thus, the ability to comprehensively utilize available data in model parametrization can help to identify gaps in existing knowledge, improve the ability of models to account for current understanding, and inform data collection efforts by identifying the knowledge gaps most responsible for uncertainty.

While the increasing sophistication of model-data fusion and uncertainty accounting is a critical step in the right direction, the complexity of such analyses can make models even less accessible. One of the reasons to make models more accessible, and to make them better at synthesizing existing data, is that they are fundamentally a formal, quantitative distillation of our current understanding of how a system works. As such, models can be used to identify gaps in our understanding and target further research. This feedback between models and data could be improved if models were routinely evaluated in a way that quantifies the value of data with respect to reducing uncertainty. We fundamentally believe that streamlining the informatics of modeling – the need to track, process, and synthesize data and model output – will make the development and application of ecological data and models more accessible, transparent, and relevant.

In this paper we present the Predictive Ecosystem Analyzer (PEcAn) as a step toward meeting these objectives. PEcAn is a scientific workflow that manages the flows of data used and produced by ecological models, and that assists with model parametrization, error propagation, and error analysis. PEcAn accomplishes two goals: first, it synthesizes data and propagates uncertainty through an ecosystem model; second, it places an information value on subsequent data collection that enables data collection that efficiently reduces uncertainty. In addition to quantifying the information content of any prediction or assessment, these techniques also help identify the gaps in our knowledge of

ecological and biogeochemical processes (Saltelli et al., 2008).

PEcAn addresses the challenge of synthesizing plant trait data from the literature in a way that accounts for the different scales and sources of uncertainty. Available data is synthesized using a Bayesian meta-analysis, and the meta-analysis posterior estimates of plant traits are used as parameters in an ecosystem model.

A model ensemble is a set of model runs with parameter values drawn from the meta-analysis posteriors estimate of plant traits. Output from a model ensemble represents the posterior predictive distributions of ecosystem responses that account for trait parameter uncertainty (hereafter “model posterior” refers to the “model ensemble output”). Sensitivity analysis and variance decomposition help to determine which traits (model parameters) drive uncertainty in ecosystem response (model posterior) (Saltelli et al., 2008; Larocque et al., 2008). These analyses help target parameters for further constraint with trait data, forming a critical feedback loop that drives further field research and provides an informative starting point for data assimilation. Here we illustrate an application of PEcAn to the assessment of aboveground yield in a perennial grass monoculture.

In the sections below, we provide an overview of the components of PEcAn’s integrated framework for data synthesis and ecological prediction. We start with a description of the methods implemented in the workflow (Implementation). This includes descriptions of the database, Bayesian meta-analysis, ensemble analysis, sensitivity analysis, and variance decomposition. Finally we present an example of the application of the system (Application) to analyze the aboveground biomass of switchgrass (*Panicum virgatum* L.), by the Ecosystem Demography model, version 2.1 (ED2) (Medvigy et al., 2009; Moorcroft et al., 2001).

Implementation

PEcAn workflow

The Predictive Ecosystem Analyzer (PEcAn) manages the flow of information into and out of ecosystem models. PEcAn is not a model itself, it is a scientific workflow consisting of discrete steps, or modules. Individual modules are building blocks of the workflow, represented by the rectangles in Figure 1, while flows of information are represented by arrows. This makes PEcAn an encapsulated, semi-automated system for model parametrization, error propagation, and analysis.

A central objective of the PEcAn workflow is to make the entire modeling process transparent, reproducible, and adaptable to new questions (*sensu* Stodden et al., 2010; Ellison, 2010). To achieve this objective, PEcAn's adheres to "best practice" guidelines for ecological data management and provenance tracking (Jones et al., 2006; Michener and Jones, 2012).

PEcAn uses a database to track data provenance and a settings file to control workflow analyses and model runs. The database records the site, date, management, species, and treatment information for each trait observation used in the meta-analysis. Settings related to the experimental design and computation are set and recorded in a separate file for each analysis.

The PEcAn source code, as well as the inputs and output used in the analysis described below (see Application) are provided as an appendix. However, new users are encouraged to utilize the latest release available on the project web site (www.pecanproject.org).

This site also provides a virtual machine and a web-interface that minimize the effort required to run PEcAn and begin using an ecosystem model. The PEcAn "virtual machine" provides all of the required software dependencies in a pre-configured desktop environment that can be run on any standard operating system using a freely available

virtualization software such as VirtualBox (www.virtualbox.org/) or VmWare Player (www.vmware.com). The virtual machine minimizes the installation time and pre-requisite knowledge required to perform analyses, and can be used to support investigation, development, and education. The web interface is even easier to use, but does not provide access to all of PEcAn’s functionality.

The PEcAn software is primarily written in R and developed in a Linux environment. It also relies on a MySQL database, bash, JAGS, and specialized R packages. PEcAn has a family of model-specific functions that manage the details of launching of model runs and reading model output.

Although PEcAn does not depend on any specific model, it was developed to support ecosystem models that run in high-performance computing environments; for this reason, it is capable of running models locally, remotely, or through queuing systems regardless of whether PEcAn is compiled locally or run as a virtual machine. The PEcAn 1.1 release described herein runs with the Ecosystem Demography model, in addition, the current release also supports SIPNET (Moore et al., 2008) and near term support for IBIS (Kucharik et al., 2000), DayCent (Parton et al., 1998), and BioCro (Miguez et al., 2009) is under development.

Trait Database

Model parameters are associated with corresponding prior distributions, and in many cases, with species-level data. Both prior distributions and data are stored in a relational database (Appendix B). PEcAn directly accesses the database, which contains additional meta-data for each data set, including site descriptions, measurement conditions, experimental details, and citations.

Trait Priors

A fundamental component of the Bayesian approach to parameter estimation is the use of priors. Priors formally incorporate knowledge of a parameter based on previous studies into a new analysis. In the current study, we leverage previously collected data from non-target species to place biologically informed constraint on the distribution of a plant trait parameter. When additional data for a specific species or plant functional type is available, priors are further constrained before being used as model parameters. When no additional data are available, these priors are used directly to parameterize the model. For the *P. virgatum* example described below, priors were set using data from all plant species, from only grass species, or from just C4 grass species depending on available data. Sources of this prior information included data from previous and ad-hoc syntheses, expert knowledge, and biophysical constraints (Table 1).

Prior distributions used in the meta-analysis were fit to one of four types of information: 1) data from multiple species, 2) the posterior predictive distribution for a new species from a meta-analysis of data (when error estimates were available), 3) a central tendency informed by data with expert constraint on the confidence interval, or 4) expert constraints on both the central tendency and confidence intervals. In case number 2, the across-species meta-analysis “posterior” informs the prior for the species-level meta-analysis. In all cases, maximum likelihood estimation was used to fit a prior distribution. When more than one candidate distribution was considered, Aikake’s Information Criterion (AIC) was used to select the best fit distribution. The choice of prior was confirmed by visually inspecting the prior density functions overlain by data or expert constraints (Figure 2).

Meta-analysis

A Hierarchical Bayes meta-analytical model (Figure 3) formally synthesizes available trait data from multiple studies while accounting for various sources of uncertainty. This Hierarchical Bayes approach integrates prior information and provides a flexible approach to variance partitioning and parameter estimation.

The meta-analytical framework is useful for summarizing data sets that include summary statistics. The trait data queried by PEcAn consist of a trait name, sample mean, sample size, and a sample error statistic. PEcAn transforms error statistics to exact or conservative (i.e., erring toward inflating the variance) estimates of precision ($\tau = 1/SE^2$) (Appendix C).

The sample mean is drawn from a normal distribution:

$$Y_k \sim N(\Theta_k, \tau_k) \quad (1)$$

Where Y_k is the sample mean of the k^{th} unique site by treatment combination (sample unit), Θ_k is the unobserved 'true' value of the trait for the k^{th} sample unit.

The meta-analysis partitions trait variability into among site, among treatment, and within-unit variance. The unobserved 'true' trait mean Θ_k is a linear function of the global trait mean, β_0 plus random effects for study site (β_{site_j}) and treatment ($\beta_{tr|site_i,j}$) and a fixed effect for greenhouse (β_{gh}):

$$\Theta_k = \beta_0 + \beta_{site_i} + \beta_{tr|site_i,j} + \beta_{gh}I(i) \quad (2)$$

Where i indexes study site, j indexes each treatment within a study, and $I(i)$ is an indicator variable set to 0 for field studies and 1 for studies conducted in a greenhouse, growth chamber, or pot experiment. The parameter used in the ecosystem model is the posterior estimate of the global mean trait value, β_0 . β_0 , has an informed prior functional

239 form and parameter specification that varies by trait and species or PFT. Methods used
 240 to elicit priors for the present study are provided in the Application section under Priors.
 241 The “site” random effects (β_{site}), accounts for the spatial (among-site) heterogeneity of a
 242 parameter. The “treatment” random effect ($\beta_{\text{tr}|\text{site}}$) accommodates differences among
 243 experimental treatments. These random effects of treatment and site are assumed to be
 244 Normally distributed with zero mean and they have diffuse Gamma priors on precision
 245 τ_{site} and τ_{tr} . Control treatments and observational studies have $\beta_{\text{tr}|\text{site}} = 0$. PEcAn
 246 dynamically adjusts the meta-analysis model specification to include terms for each level
 247 of site and treatment, or greenhouse studies as required by available data. To ensure that
 248 the prior on precision remains sufficiently diffuse when the magnitude of a parameter is
 249 small, the scale parameters in the gamma priors on random effect precision terms (τ_{site}
 250 and $\tau_{\text{tr}|\text{site}}$) are scaled to $(\bar{\beta}_0^2/1000)$ when the prior on β_0 has a mean $\bar{\beta}_0 < \sqrt{10}$.
 251 A “greenhouse” fixed effect β_{gh} accounts for potential biases associated with plants grown
 252 in a greenhouse, growth chamber, pot, or other controlled environment. This “greenhouse”
 253 effect, β_{gh} , has a diffuse Normal prior with a mean of zero and a precision of 0.01.
 254 The observation precision (precision = $1/\text{variance}$) for the k^{th} sample mean, τ_k , is
 255 determined based on the within-unit precision, τ_Y , and the sample size, n , as $\tau_k = n \times \tau_Y$
 256 (since $SE = SD/\sqrt{n}$). A common within sample unit precision, τ_Y , is assumed in order
 257 to accommodate literature values with missing sample sizes or variance estimates. The
 258 sample standard error, se_k , is drawn from a Gamma distribution with parameters
 259 informed by the sample size, n , and within-site precision, τ_Y :

$$\frac{1}{n \times se_k^2} \sim \text{Gamma}\left(\frac{n}{2}, \frac{n}{2\tau_Y}\right) \quad (3)$$

260 τ_Y has a diffuse gamma prior. Unlike the mean and variance parameters, missing values
 261 of n cannot be estimated and are conservatively set either to 2 (when existence of a

262 variance estimate indicates $n \geq 2$) or to 1 (if no variance estimate is given).
 263 The random and fixed effects and the among study, among treatment, and within-unit
 264 precisions are used to evaluate the importance of different sources of uncertainty.
 265 The meta-analysis module in PEcAn is fit using JAGS software (version 2.2.0, (Plummer,
 266 2010)) called from within R code that handles data manipulations and meta-analysis
 267 model specification in JAGS. JAGS uses standard Markov Chain Monte Carlo (MCMC)
 268 methods (Gelman and Rubin, 1992) to approximate the posterior distribution of the
 269 terms in the meta-analysis. To overdispense the n MCMC chains, initial values of β_0 are
 270 set to the $\frac{1}{n+1}, \dots, \frac{n}{n+1}$ quantiles of the prior on β_0 ; for the default $n = 4$ chains, this
 271 would be the $\{0.2, 0.4, 0.6, 0.8\}$ quantiles. Following Gelman and Shirley (2011), PEcAn
 272 discards the first half of each chain, thins each chain to 5000 samples and then combines
 273 the chains into a single vector of samples for each term in the meta-analysis model. Trace
 274 plots and the Gelman-Rubin convergence diagnostic (Gelman and Rubin, 1992) are used
 275 to assess chain convergence. Density plots (Figure 4) are used to visually compare the β_0
 276 chain to data and priors. The significance of the greenhouse effect is evaluated by
 277 calculating a two-sided probability that $\beta_{\text{gh}} \neq 0$.
 278 When species-level data are unavailable, the posterior distributions are equivalent to the
 279 priors.
 280 Each term in the meta-analysis is represented as a vector of MCMC samples from the
 281 posterior distribution. Statistical summaries of the parameters can easily be calculated
 282 from these chains, and chains can also be directly sampled for use in ecosystem model
 283 parameterization. When the β_0 chains are sampled for the ecosystem model ensemble,
 284 the meta-analysis posteriors become the model ensemble priors.

Model Analysis

Ensemble Analysis

Typically, ecosystem models are run for a single model parameterization. For example, the model could be evaluated at the median value of each parameter. However, this approach only provides a point estimate with no accounting for parameter uncertainty. To propagate parameter uncertainty through the ecosystem model, PEcAn uses standard ensemble-based Monte Carlo approaches. An ensemble of model runs is a set (e.g. 500 or 1000) of model runs that are parameterized by sampling from the trait parameter distributions. For each ensemble member, parameter sets are sampled from the full joint parameter distribution of β_0 , the vector of all model parameters. As a result, the model ensemble approximates the posterior distribution of the ecosystem model output. The model ensemble produces a posterior distribution of the ecosystem model output that can be summarized with standard statistics (e.g. mean, standard error, and credible interval).

Sensitivity Analysis

Sensitivity analyses are used to understand how much a change in a model parameter affects model output; sensitivity is the derivative, $df/d\beta_{0t}$, of the model (f) with respect to the estimate of β_0 for trait t . PEcAn approximates the sensitivities based on univariate perturbations of model parameters. These approximations are necessary because analytical solutions for sensitivity are not tractable for most ecosystem models, and PEcAn is designed to be flexible and applicable to any such model. One disadvantage of traditional perturbation-based sensitivity analyses is that the perturbations are usually arbitrary, for example varying each parameter by a fixed percentage of its value (Larocque et al., 2008) rather than over a meaningful range of the parameter. These traditional approaches make interpretation of sensitivities difficult because they fail to

309 acknowledge the distribution or uncertainty of each parameter. In this regard, PEcAn
 310 offers a distinct advantage over traditional sensitivity analyses because parameters are
 311 varied based on the meta-analysis posterior parameter distributions.
 312 Based on initial exploratory analyses, we found a local perturbation to be inadequate for
 313 capturing the responses in most parameters so we instead estimate sensitivities using a
 314 global univariate sensitivity analysis. By default, PEcAn evaluates each parameter at the
 315 posterior median and at the six posterior quantiles equivalent to $\pm[1, 2, 3]\sigma$ in the
 316 standard normal while holding all other variables constant at their posterior median. The
 317 relationship between model output and each model parameter β_{0t} is then approximated
 318 by a natural cubic spline $g_t(\beta_{0t})$ that interpolates through the evaluation points. The
 319 model sensitivity to each parameter is approximated by the derivative of the spline
 320 $(dg_t/d\beta_{0t})$ at the parameter mean. In addition to the sensitivity analysis, this set of
 321 spline approximations is used in the variance decomposition, in partitioning residual
 322 variance, and in evaluating the effect of ensemble size on the estimate of model variance.
 323 To facilitate comparisons among the trait sensitivities, despite differences in the units on
 324 different traits, we tabulate the coefficient of variation (normalized parameter variance)
 325 and the elasticity (sensitivity with terms df and $d\beta_{0t}$ standardized by the mean model
 326 output and parameter mean respectively).

327 **Variance Decomposition**

328 Variance decomposition aims to explain how much each input parameter contributes to
 329 uncertainty in model output (Cariboni et al., 2007). Although the present analysis
 330 focuses on model parameters, these methods can be extended to address uncertainty in
 331 initial conditions or model drivers.

332 The Delta Method uses Taylor series expansion to approximate the probability
 333 distribution of a continuous function of random variables (Oehlert, 1992; pp. 240–245 in

334 Casella and Berger, 2001). In this study, the model output $f(\beta_0)$ is a function of a vector
 335 of the full set of parameters. After approximating the distribution of $f(\beta_0)$, it is possible
 336 to estimate the variance of the model output. The first step is to derive the Taylor series
 337 approximation of the variance of a function (Casella and Berger, 2001, equation 5.5.9 in):

$$Var(f(\beta_0)) \approx \sum_{t=1}^m Var \left(f(\overline{\beta_{0t}}) + \frac{df}{d\beta_{0t}}(\beta_{0t} - \overline{\beta_{0t}}) + \dots \right) \quad (4)$$

$$= \sum_{t=1}^m \left(\frac{df}{d\beta_{0t}} \right)^2 Var(\beta_{0t}) + \omega \quad (5)$$

338 where m is the number of parameters in the model, and the error term ω accounts for
 339 higher order terms in the Taylor series, and β_{0t} is the estimate of β_0 from the
 340 meta-analysis (equation 2) for each trait, t .

341 With this approximation, it is straightforward to estimate the variance contributed by
 342 each parameter. The terms in this form of the variance decomposition can be estimated
 343 directly from the preceding analyses: $Var(f(\beta_0))$ is the variance of the model ensemble;
 344 $Var(\beta_{0t})$ is the posterior variance of trait β_{0t} from the meta-analysis (equation 2); and
 345 $df/d\beta_{0t}$ is the model sensitivity at the parameter mean $\overline{\beta_{0t}}$. The resulting assertion is
 346 that the variance of model output is equal to the sum over the variance of each trait
 347 times its sensitivity squared plus a closure term, ω .

348 We found that the traditional Taylor polynomial approach to variance decomposition
 349 produced a poor closure of the total variance of the model output: for more sensitive
 350 parameters, a linear approximation of $f(\beta_0)$ provided unrealistic estimates of the
 351 sensitivity function that overestimated variance. Increasing the order of the Taylor series
 352 expansion actually exacerbated this problem (results not shown). One problem with the
 353 polynomial approximation is that, unlike polynomials, most response variables in
 354 ecosystems and ecosystem models tend to be asymptotic at both high and low values of a
 355 trait. For example, when assessing aboveground biomass there is a lower bound of zero

biomass and most parameters become progressively less sensitive, if not genuinely asymptotic, at their upper bound. This asymptotic behavior is poorly approximated by a polynomial because polynomials are unbounded at extreme parameter values. Therefore, we sought a better approximation for the variance decomposition. First, we formulated a more generalized form of the variance decomposition (equation 4):

$$Var(f(\boldsymbol{\beta}_0)) = \sum_{t=1}^m Var(g_t(\beta_{0t})) + \omega \quad (6)$$

The spline $g_t(\beta_{0t})$ is a statistical emulator of the model response to trait t that transforms β_{0t} from the parameter domain to the model domain. The univariate contribution of each parameter to variance of the model output is thus $Var(g_t(\beta_{0t}))$. Equation 6) only requires β_{0t} from the preceding meta-analysis, $g_t(\beta_{0t})$ from the sensitivity analysis, and $Var(f(\boldsymbol{\beta}_0))$ from the ensemble analysis. The final term, ω , is the closure between the right hand side and the left hand side of the variance decomposition; ω represents the effects of the higher order terms in the Taylor approximation and the covariance terms between parameters. This closure term is intended to represent parameter interactions that are excluded from the univariate variance decomposition (equation 6). Negative trade-offs among physiological traits would result in ω less than zero. However, our estimate of ω also includes errors associated with using finite sample sizes, the spline approximation in each $g_t(\beta_{0t})$, and biological range restrictions on model output that are not reflected in the variance decomposition (equation 6).

One approach to partition the error in the closure term is to use the univariate spline functions from the sensitivity analysis to estimate what the model output would be for each of the parameter sets used in the model ensemble; we call this estimate the “spline

ensemble”:

$$\mathbf{g}_\ell(\boldsymbol{\beta}_0) = \mathbf{g}(\hat{\boldsymbol{\beta}}_0) + \sum_{t=1}^m \left(g_t(\beta_{0t\ell}) - g_t(\hat{\beta}_{0t}) \right) \quad (7)$$

In this equation, $\mathbf{g}_\ell(\boldsymbol{\beta}_0)$ is the spline estimate of the model output for the ℓ^{th} ensemble member and $\hat{\beta}_{0t}$ is the posterior median parameter value.

Although the individual splines may respect range restrictions on output variables (e.g. biomass values cannot fall below zero), combinations of the splines evaluated for a set of unfavorable traits can fall outside these ranges. For parameter sets that give a biologically implausible estimate of negative biomass ($\mathbf{g}_\ell(\boldsymbol{\beta}_0) < 0$), the estimate is set to zero. The only difference between the variance of the spline ensemble (equation 7) and the variance decomposition (equation 6) is that range restrictions are not corrected for in the variance decomposition. Therefore, the spline ensemble allows us to estimate the effect of using combinations of spline estimates that do not respect the zero bound on biomass in the variance decomposition. The difference between the model ensemble and the spline ensemble provides an estimate of parameter interactions in the model because the spline ensemble does not include the parameter interactions that exist in the model. The precision of the estimate of model ensemble variance is affected by the number of runs in the ensemble. When the computational expense of the model itself limits the ensemble size, there can be significant uncertainty in the estimate of ensemble variance. The uncertainty in a sample variance is estimated as

$$Var(s^2) = \frac{1}{n} \left(\mu_4 - \frac{n-3}{n-1} \sigma^4 \right) \quad (8)$$

(Mood et al., 1974, , p 239) where μ_4 is the fourth central moment. $Var(s^2)$ scales inversely with sample size. The effect of the limited model ensemble size on uncertainty in the estimate of ensemble variance is measured in two ways. The first way is to calculate $Var(s^2)$ for the model ensemble ($n = 500$). The second way is to compare

400 $Var(s^2)$ of the spline ensemble with 500 and 10,000 runs. The 95%CI for s^2 is calculated
401 as $s^2 \pm 1.96s_{s^2}$ where $s_{s^2} = \sqrt{Var(s^2)}$.
402 The errors introduced from using a spline approximation of the model response can not
403 be estimated based on the existing output, but it is small in comparison to the other
404 effects given the range restrictions imposed by the spline interpolation.
405 The results of a model ensemble are posterior estimates of aboveground biomass.
406 However, we also distinguish between ensembles depending on the nature of model
407 parameters. First, we ran a “prior model ensemble” using an ensemble of parameter sets
408 drawn from prior distributions, and then a “posterior model ensemble” drawn from
409 meta-analysis posteriors.

410 **Application: Switchgrass Monoculture**

411 We demonstrate the application of PEcAn to estimate the aboveground yield of an
412 experimental switchgrass (*Panicum virgatum*) monoculture. The first step to applying
413 PEcAn was to construct an appropriate set of priors based on data syntheses and expert
414 knowledge. These priors were conservative estimates of the plant trait parameters based
415 on information other than species level data. Next, switchgrass trait data from both
416 previous studies and field measurements were summarized using meta-analysis to
417 constrain the prior parameter estimates. The Ecosystem Demography model version 2.1,
418 (Medvigy et al., 2009; Moorcroft et al., 2001) was used to simulate plant growth.
419 The model ensemble and sensitivity analysis were performed using both the prior and
420 posterior parameter estimates. By comparing the prior model ensemble to the posterior
421 model ensemble, we are able to evaluate the ability of species level data to reduce model
422 uncertainty.

423 To evaluate model performance, we compare the ensemble estimates of aboveground

biomass with observed yields (Heaton et al., 2008; Wang et al., 2010, Figure 5).

Site

Switchgrass (*Panicum virgatum*) is a perennial grass native to North America that has received attention as a potential cellulosic biofuel crop (McLaughlin and Kszos, 2005; Wang et al., 2010). We modeled the aboveground biomass production of a switchgrass monoculture and compared model estimates to a monoculture planted in 2002 at the University of Illinois Agricultural Research and Education Center in Urbana, IL (40.09 N, 88.2 W). The climate at this site is characterized by hot, humid summers and cold winters with a 50 year (1959-2009) mean annual temperature of 11 °C and mean annual precipitation of 1000 mm yr⁻¹ (Angel, 2010). Meteorological data used to drive the model were downloaded from the North American Regional Reanalysis (Mesinger et al., 2006). Soil is a silt loam from the Drummer-Flanagan soil series; texture data was obtained through the USDA NRCS web soil survey website (websoilsurvey.nrcs.usda.gov). The yield and other aspects of this ecosystem have previously been reported (Heaton et al., 2008).

Ecosystem Demography Model

We used the Ecosystem Demography Model, version 2 to model the productivity and soil carbon pools in this switchgrass agro-ecosystem. ED2 is a terrestrial biosphere model that couples age- and stage-structured plant community dynamics with ecophysiological and biogeochemical models. The biophysical land-surface model in ED2 allows plant uptake and growth to respond dynamically to changes in weather and soil hydrology (Medvigy et al., 2009). ED2 has the ability to link short-term, physiological responses to environmental conditions with realistic, long-term successional changes in ecosystem structure and composition (Moorcroft et al., 2001). While other models have both

448 succession and physiology, ED2 also has explicit spatial scaling, a sub-daily time-step,
 449 and the ability to couple with to a land surface model (Dietze and Latimer, 2011).
 450 ED2 incorporates a mechanistic description of plant growth that accounts for the fast
 451 temporal responses of plants to changes in environmental conditions. In this study, we
 452 vary fifteen model parameters based on observable plant traits that control carbon
 453 uptake, carbon allocation, turnover, and reproduction (Table 1, Figures 2, 4).
 454 ED2 calculates photosynthetic rates using the enzyme kinetic model developed for C3
 455 plants (Farquhar and Sharkey, 1982; Ball et al., 1987) and the modifications for C4
 456 (Collatz et al., 1992). $V_{c,max}$ sets the upper bound on the rate of Rubisco-limited
 457 photosynthesis, while light limited photosynthesis is constrained by the quantum
 458 efficiency parameter, and a threshold parameter controls the minimum temperature at
 459 which photosynthesis will occur. Stomatal conductance is calculated using the Leuning
 460 variant of the Ball-Berry model (Leuning, 1995) and is controlled by the stomatal slope
 461 parameter. Leaf boundary layer conductance depends on the leaf width parameter.
 462 Together, stomatal conductance and leaf boundary layer conductance affect carbon and
 463 moisture fluxes and the leaf energy balance. Specific leaf area (SLA) determines the
 464 amount of leaf area produced per unit leaf biomass investment.
 465 In addition to photosynthesis, ED2 also accounts for carbon allocation to growth,
 466 respiration, and for the turnover rate of carbon pools. These parameters include: one to
 467 partition between leaf and fine root growth; one for allocation to reproduction; two
 468 respiration parameters associated with growth respiration and root maintenance
 469 respiration; and two parameters to control the rates of leaf and root turnover.
 470 Finally, three demographic parameters control seed dispersal, seedling mortality, and
 471 adult mortality due to carbon limitation (Table 1).

Priors

Priors from data

Priors were estimated by finding the best fit distribution to raw data sets include SLA and leaf turnover rate from the GLOPNET database (Wright et al. (2004), $n = 125, 40$ respectively), root turnover rate (Gill and Jackson (2000), $n = 66$), and quantum yield (Skillman (2008), $n = 56$). Candidate distributions for these priors were Gamma, Weibull, log-Normal, and F because each of these traits is bound at zero. In all cases we are interested in using the full distribution of across-species data as our prior constraint on what one individual species is capable of doing, as opposed to using the estimate of the mean of this distribution as our prior.

Quantum yield data represent a survey of published values of quantum yield in C4 monocots (Skillman, 2008); original data were provided by the author and restricted to measurements made under photorespiratory conditions (ambient CO_2 and O_2) (J. Skillman, personal communication). Given the narrow range of data ($\text{CV} = 11\%$), the normal distribution was also considered but was not the best fit.

Priors from meta-analysis

We used meta-analysis to calculate a prior from data when summary statistics and sample sizes were available. The meta-analysis model used to calculate prior distributions is similar to the one used by PEcAn to summarize species-level data (equation 2), with three differences. First, there are no site, treatment, or greenhouse effects. Second, data from multiple species were used. Third, we generated a posterior predictive distribution to predict the distribution of trait values for an unobserved C4 plant species, unlike the species-level meta-analysis, which estimated the global mean parameter value. Thus, the model included plant functional type (PFT) as a random effect:

$$\Theta_{\text{species}} = \beta_0 + \beta_{\text{PFT}} \quad (9)$$

Stomatal slope is the empirical slope coefficient in the (Leuning, 1995) model of stomatal conductance. Estimates of this parameter are sparse, so we collected new data for this study. Stomatal slope was estimated using measurements of four leaves from each of five field-grown energy crop species during the 2010 growing season (Appendix A). The five species included two C4 grasses: Miscanthus (*Miscanthus x giganteus*) and Switchgrass (*P. virgatum*) planted in 2008 and three deciduous tree species: Red Maple (*Acer rubrum*), Eastern Cottonwood (*Populus deltoides*, and Sherburne Willow (*Salix x Sherburne*) planted in 2010 as 2 year old saplings. All plants were grown at the Energy Biosciences Institute Energy Farm (40°10'N, 88°03'W). We used the data from the three tree species and Miscanthus to calculate the posterior predictive distribution of an unobserved C4 grass species, and used this distribution as the prior estimate for Switchgrass stomatal slope.

Maximal carboxylation rate (V_{cmax}) data consists of ninety-four C3 species (Wullschleger, 1993) plus three C4 species (Kubien and Sage, 2004; Massad et al., 2007; Wang et al., 2011). To express V_{cmax} at a common temperature of 25°C for all species, we applied an Arrhenius temperature correction (Appendix C). The Wullschleger (1993) data set included a 95% CI and an asymptotic error calculated by the SAS nlin procedure (Stan Wullschleger, personal communication). We used the asymptotic error as an estimate of SE, the 95% CI to estimate SD ($\text{SD} = \frac{\frac{1}{2}\text{CI}}{1.96}$), and then estimated n as $\hat{n} = \left(\frac{\text{SE}}{\text{SD}}\right)^2$. Plant species were classified into five functional types: C3 grass, C4 grass, forb, woody non-gymnosperm, and gymnosperm based on species records in the USDA PLANTS Database (USDA and NRCS, 2011). Ambiguous species (those with both forb and woody growth forms, $n = 15$) were excluded.

Leaf width data represent 18 grass species grown in a common garden greenhouse experiment (Oyarzabal et al., 2008). *P. virgatum* was among the 18 species, and was excluded from the prior estimation but used as raw data in the meta-analysis. The remaining seventeen species were divided into C3 and C4 functional types. Relative to the small sample of C4 grasses, switchgrass leaf width was an outlier; inflating the variance four-fold reduced the prior information content to account for this discrepancy. Root respiration rate values were measured on thirty-six plants representing five functional types, including six C4 grass species (Tjoelker et al., 2005). As before, *P. virgatum* data was excluded from the prior estimation and used as raw data in the species-level meta-analysis.

Priors from limited data and expert knowledge

When available data were limited to a few observations, these were used to identify a central tendency such as the mean, median, or mode, while expert knowledge was used to estimate the range of a confidence interval. An optimization approach was used to fit a probability distribution to this combination of data and expert constraint.

In order to estimate the fine root to leaf ratio for grasses, we assume fine roots account for all belowground biomass (Jackson et al., 1997) and that leaves account for all above ground biomass. Roots account for approximately 2/3 of total biomass across temperate grassland biomes (Saugier et al., 2001, Table 23.1), so we constrained a beta prior on the root fraction to have a mean of 2/3 by setting $\alpha = \beta/2$ since the mean of a beta is defined as $\frac{\alpha}{\alpha+\beta}$. To convert from proportion to ratio, we used the identity: if

$X \sim \text{Beta}(\frac{\alpha}{2}, \frac{\beta}{2})$ then $\frac{X}{1-X} \sim F(\alpha, \beta) \times \frac{\alpha}{\beta}$. We constrained the 95%CI = $[1/3, 10/11]$, equivalent to a fine root to leaf ratio with a mean fixed at two and a 95%CI = $[1/2, 10]$.

We simulated the distribution of the fine root to leaf ratio by drawing 10000 samples from a $F(2\alpha, \alpha)$ distribution and multiplying these samples by two.

Seed dispersal in ED2 represents the proportion of seed dispersed outside of a 7.5m radius plot, which we approximate as a beta distribution. Although no direct measurements of seed dispersal were available, it was straightforward to parametrize a ballistic model of seed dispersal (Ernst et al. (1992), from Creemer 1977): $D = \frac{V_w H}{V_t}$. This model relates dispersal distance D to terminal velocity V_t , wind speed V_w , and seed height H . Although more sophisticated treatments of dispersal exist and are important for estimating low probability long distance dispersal events (Clark et al., 1999; Thompson and Katul, 2008), the Ernst et al. (1992) model is sufficient for relatively short dispersal distances required in the present context.

Values of terminal velocity, V_t , of grass seeds were taken from two references, (Ernst et al., 1992; Jongejans and Schippers, 1999) and these data were best described as $V_t \sim \text{Gamma}(2.93, 1.61)$.

Next the heights from which the seeds are dropped was estimated from calibrated photographs of reproductively mature switchgrass at a field site in Urbana, IL:

$H \sim N(2, 0.5)$. Finally, wind speed observed at the site were fit to a Weibull distribution (Justus et al., 1978). $V_w \sim \text{Weibull}(2.4, 0.712)$ (Marcelo Zeri, unpublished wind and height data). Given these distributions of V_w , H , and V_t , sets of 100 dispersal distances were simulated 10000 times to calculate the fraction of seeds in each simulation dispersed beyond 7.5m,

Priors informed by expert knowledge

When no data were available, expert knowledge was used to estimate the central tendency and confidence interval for a trait parameter. Again, optimization was used to fit a probability distribution to these constraints.

The minimum temperature of photosynthesis for C4 grasses was given a prior based on expert knowledge with a mean of 10°C and a 95%CI = [8, 12]°C that fits a normal

($\mu = 10$, $\sigma = 1.02$) distribution (Don Ort, UIUC, personal communication, 2010).

The growth respiration factor is the proportion of daily carbon gain lost to growth respiration. Because it is a proportion, the beta distribution was fit with a mean set equal to the ED2 default parameter value, 0.33 and a 95%CI = [0.05, 0.60], conservatively based on the range of construction costs reviewed by Amthor (2000).

Seedling mortality factor represents the proportion of carbon allocated to reproduction that goes directly to the litter pool. Given the default ED2 parameter is 0.95, we stated a beta prior with a median at 0.95, and a 95%CI = [$2/3$, 1].

The mortality factor in ED2 is the rate parameter in the negative exponential relationship between carbon balance and mortality (Medvigy et al., 2009). The default parameter for all plant functional types (PFT's) in ED2 is 20, and our weakly informed gamma prior sets this as the median and gives a 95%CI = [5, 80].

Reproductive allocation represents the proportion of carbon in the storage pool allocated to reproduction. This parameter is a proportion and has a default value of 0.33 in ED.

The Beta(2, 4) distribution has a mean of $1/3$ and a 95%CI = [0.05, 0.72] representing relatively high uncertainty. This distribution implies that the plant may allocate any fraction of the carbon pool to reproduction between but not equal to 0 and 1 and has an 80% probability that less than half of the carbon pool will be allocated to reproduction.

Switchgrass Trait Meta-analysis

Switchgrass trait data used to constrain model parameters are stored in the Biofuel Ecophysiological Trait and Yield database (BETYdb, www.betydb.org), a database designed to support research on biofuel crops. BETYdb includes both previously published and primary data (Appendix A). Prior to entry in the database, data were converted to standard units chosen for each variable (Table 1).

Trait data available for *Panicum virgatum* include V_{cmax} , SLA, leaf width, fine root to

leaf ratio, root respiration, stomatal slope, and root turnover rate (Figure 4, Table 2).
Methods used to collect these data and site descriptions are available in the source
references (Appendix A). Each meta-analysis was run with four 50,000 step MCMC
chains.

Model Analysis

We executed a ten-year, 500 run ensemble of ED2 using parameter values drawn from the
prior or posterior parameter distributions. The model was run for the years 1995-2006 to
simulate the field trials conducted by Heaton et al. (2008). The model output of interest
was the December mean aboveground biomass (AGB) during the years 2004–2006,
simulating the yields of the mature stand (Heaton et al., 2008). The ensemble estimate
was also compared to the larger set of all reported switchgrass yield data (Wang et al.,
2010).

Runs resulting in yields less than 2 Mg/ha were considered non-viable parameter
combinations. To test if prior and posterior parameter sets resulted in different fractions
of non-viable runs, we estimated the posterior probability of a non-viable run as a
binomial posterior from a beta-binomial model with a flat (Beta(1, 1)) prior. Then, we
calculated the two-tailed probability that the difference between these binomial posteriors
was $\neq 0$.

Results

Trait Meta-analysis

Switchgrass data were collected from the literature and field for seven of the model
parameters: specific leaf area (SLA) ($n = 24$), leaf width (39), V_{cmax} (4), fine root to leaf
allocation ratio (4), stomatal slope (4), root respiration rate (1), and root turnover rate

(1). Table 2 summarizes the meta-analysis for each of these parameters, including the posterior mean and 95% CI of the global mean, the fixed greenhouse effect, and each of the variance components (reported as standard deviations).

SLA and leaf width data were from multiple sites, but the meta-analysis provided no evidence for among site variability in excess of within site variability (σ_Y and σ_{site} , respectively, in Table 2). For the remaining traits, there was insufficient spatial sampling to assess site to site variability. Greenhouse growing conditions had a positive effect on both SLA ($P = 0.027$), and leaf width ($P = 0.052$).

Figure 4 compares parameters before and after incorporating data in the meta-analysis.

A reduction in parameter uncertainty is seen as the reduction in the spread of the posterior (black) compared to the prior (grey) parameter distributions. The influence of the prior information on the posterior distribution increased when the prior was more constrained or when less data were available for use in the meta-analysis. For example, data substantially constrained the uncertainty in the V_{cmax} and SLA posteriors relative to the priors. By contrast, there was little effect of additional data on the parameter estimates for fine root to leaf allocation and root respiration rate; these parameters had relatively well constrained priors and limited species-specific data.

Model Analysis

Ensemble

Within the model ensemble analysis (Figure 5), both the prior and posterior parameterizations produced yield estimates that were consistent with yields observed at the Urbana site for which the model was run (Heaton et al., 2008) and with 1902 previously reported yields of switchgrass (Wang et al., 2010). In both the prior and posterior ensembles, the predicted aboveground biomass was clearly bimodal. These two

modes had little overlap and a distinct break at two Mg/ha. We inferred that the first
 peak represents non-viable plants generated by unrealistic parameter sets so plants with
 aboveground biomass less than two Mg/ha were considered “non-viable”. When
 summarizing the model output, we consider viable and non-viable ensemble members
 separately; all runs are considered in the sensitivity analysis and variance decomposition.
 A greater percentage of runs in the prior ensemble fell below this threshold (53.4 vs 36.6,
 $P \simeq 0$).
 Compared to the prior ensemble prediction, the data-constrained posterior runs had lower
 median yields and a more constrained 95% credible interval (16.5[7.2, 37] Mg/ha vs
 25[7.7, 56] Mg/ha). This reflects the substantial shrinkage of the posterior relative to the
 prior SD estimates of model output uncertainty (from $\sigma = 19.7$ to $\sigma = 11.9$). In
 particular, the upper tail of the modeled yield was reduced toward the observed yields.
 Despite the reduction in ensemble uncertainty, the ensemble posterior yield was still
 relatively imprecise and had much greater uncertainty than the field trial (Heaton et al.,
 2008, $\sigma = 4.1$) or the meta-analysis of all observations (Wang et al., 2010, ($\sigma = 5.4$)).
 The spline ensemble viable plants had a median 18.8[2.9, 48] and $\sigma = 13$.

Sensitivity Analysis

Sensitivity analysis demonstrated that traits varied in their effect on on aboveground
 biomass (Figure 6), and many of these relationships are clearly non-polynomial. For
 example, parameters associated with photosynthesis and carbon allocation - including
 V_{cmax} , SLA, growth respiration, and root allocation - were particularly sensitive. For
 particularly sensitive parameters, the sensitivity functions had coverage of unrealistic
 yields greater than 30 Mg/ha. Constraining SLA and V_{cmax} parameters with data
 resulted in a more realistic range of yields. On the other hand, aboveground biomass was
 largely insensitive to leaf width, seed dispersal, and mortality rate.

Variance Decomposition

The variance decomposition showed that data-constrained parameters substantially reduced their contribution to overall model variance (Figure 7). Prior to including species-specific trait data, SLA contributed the most to model uncertainty, followed by growth respiration, fine root allocation, V_{cmax} , seedling mortality, and stomatal slope (right panel, grey bars Figure 7). Incorporating species level data substantially reduced the contributions of SLA, V_{cmax} , seedling mortality, and stomatal slope to model uncertainty. For example, SLA fell from first to fourth and stomatal slope fell from sixth to fourteenth in rank contribution to ensemble variance. Although the addition of data reduced parameter uncertainty for fine root to leaf allocation, aboveground biomass was more sensitive to this parameter at the posterior median. These changes cancelled each other out, and as a result the contribution of the fine root allocation parameter to ensemble variance remained constant.

The variance of the ensemble was less than the variance calculated in the variance decomposition, and this difference is the closure term, ω . The closure term accounted for approximately 28.455981108897% of the variance decomposition estimate (Table 3b).

There was no effect of increasing the sample size from 500 to 10000 on the variance estimates (Table 3a).

Discussion

Switchgrass Trait Meta-analysis

When species-level data were available, the meta-analysis constrained estimates of the trait mean parameter (Figure 4) and provided insight into the sources of parameter uncertainty (Table 2). In the context of constraining model parameters, we were

interested in accounting for but not directly investigating the specific effects of site,
 treatment, or greenhouse effects. However, we can use the meta analysis results to
 identify sources and scales of parameter variability. This insight into parameter
 variability helps inform future sampling designs, development of the ecosystem model,
 and improvement of methods used to parametrize the ecosystem model.
 Where data from multiple sites were available, we could evaluate the relative importance
 of within versus among-site variance for the range of sites represented in the data
 (Table 2). Recent studies demonstrate important effects of intraspecific trait variability
 on ecosystem functioning (Breza et al., 2012; Albert et al., 2011; Violle et al., 2012).
 Therefore, for traits that do exhibit greater variability across than within sites, explicit
 inclusion of spatial, environmental, and even genetic information into the meta-analytical
 model would be justified. This approach would enable the estimation of site-specific
 parameters for use in the ecosystem model and will be investigated in future development
 of the meta-analysis module.
 For the other parameters (V_{cmax} , fine root allocation, root respiration rate, and root
 turnover rate) data came from one site, so it is not possible to estimate the across-site
 variability. For these traits, obtaining data from additional sites would better constrain
 both the global mean and the across-site variance. This additional data collection is
 particularly justified for traits that contribute most to the uncertainty in the model
 ensemble.

Model Ensemble

Despite the large reduction in model uncertainty from the prior to the posterior model
 ensemble, the uncertainty in projected yield is substantial (Figure 5) and further
 constraint would increase the utility of this model output. However, the explicit
 accounting of parameter uncertainty is an important first step toward producing more

informative model output. If model parameters had been treated as fixed constants, we would have no estimate of model uncertainty; without an estimate of uncertainty, the similarity between the modeled (16.5 Mg/ha) and observed (12.0 Mg/ha) median yields would be difficult to interpret; a naive interpretation could create false confidence in the model. Including the non-viable plants would have made the model mean more accurate (Figure 5), but the 90%CI would have been less accurate, containing the possibility that switchgrass would not grow in Champaign County, Illinois, even though extensive research (Heaton et al., 2008; VanLoocke et al., 2012, personal observation) demonstrates that it does grow very well in this area.

The reduction in median modeled yield in the posterior relative to the prior model ensemble 5 is consistent with the reduced probability of high SLA and V_{cmax} values in the posterior relative to the prior distributions. As expected, the use of switchgrass trait data to inform model parameters succeeded in both reducing total uncertainty and bringing modeled yield in line with observations of switchgrass yields both at this site (Heaton et al., 2008) and globally (Wang et al., 2010). Reducing uncertainty in model outputs, in this case yield, is key to increasing the value of ecological forecasts (Clark et al., 2001).

While reducing uncertainty does not necessarily increase model accuracy, an estimate of model uncertainty is a first step toward generating meaningful statistical inference from the model itself. Without an estimate of model uncertainty, it is not possible to make such a basic inference as the probability that the model predictions overlap with observed data; this limits the capacity of models to inform research and applied problems (Clark et al., 2001). Instead, comparisons of ecosystem models with observations have focused on differences and correlations between model output and data (Bellocchi et al., 2010; Schwalm et al., 2010; Dietze et al., 2011) without providing a confidence interval around the model output itself. The ability to identify, with confidence, the conditions under which a model produces valid output helps determine appropriate applications of the

740 model and it helps to identify targets for further model development (Williams et al.,
 741 2009). While parameter uncertainty is clearly just one of many sources of uncertainty in
 742 models (McMahon et al., 2009), and constraining model parameters by no means
 743 guarantees that a model will match reality, is difficult to assess the accuracy of a model if
 744 it has low precision. In deterministic models, such as most ecosystem models, parameter
 745 uncertainty is a major driver of the precision of a model.

746 In this study, we can state with 90% Confidence that the mean Switchgrass yield during
 747 the Heaton et al. (2008) study should have been between 7.2 and 37, and if we had made
 748 this prediction in advance, we could have said that we were correct because the mean did
 749 fall within this range. We can also see that the model uncertainty contains the 90% CI
 750 for observed switchgrass yields globally (Wang et al., 2010), even though we know that
 751 important drivers of variability in the global meta-analysis (e.g., climate, soil) are
 752 different from the source of uncertainty in our model predictions (e.g., parameters). The
 753 model ensemble left open the possibility that the yields could have been much more or
 754 much less than was actually observed, and we conclude that much of this variability could
 755 be constrained with additional trait level data or data assimilation. Wang et al 201x (in
 756 review, Ecological Applications #12-0854) provides an example of combining the PEcAn
 757 meta-analysis and variance decomposition with data assimilation of biomass to constrain
 758 uncertainty in parameter estimates and improve the accuracy and precision of model
 759 output. Once the model can make more precise predictions, it will be possible to begin
 760 investigation of other sources of uncertainty, such as model structure and state variables
 761 (e.g. climate, soil).

762 Although the present analysis focuses on modeled aboveground biomass, PEcAn can
 763 analyze any model output, including pools and fluxes of water, energy, and carbon.
 764 Indeed, PEcAn's approach to the synthesis of data and mechanistic models is
 765 independent of the system being modeled, and thus has potential applications far beyond

the scope of its current development to support ecosystem modeling.

Variance Decomposition

Variance decomposition quantified the contribution of each parameter to model uncertainty, helping to identify a subset of parameters for further constraint. SLA, V_{cmax} , fine root to leaf ratio, and leaf turnover rate dominated uncertainty in yield prior to incorporating species level data. Therefore, SLA, which can be measured quickly and at low cost, would make a good first target for reducing uncertainty when a new species is evaluated. SLA also correlates strongly with other important model parameters, such as photosynthetic rate, leaf lifespan, and nitrogen content (Wright et al., 2004). The ranking of parameters based on variance contribution depends on the response variable of choice (in this case, aboveground biomass) as well as the conditions of the run (duration, soil, climate), and the species or PFT being evaluated. In general, for a given species and model output, overall patterns of parameter importance are consistent across a broad range of climates (Wang et al., 201x, in review, Ecological Applications #12-0854). Variance decomposition (equation 6) demonstrates that it is not parameter uncertainty or model sensitivity alone, but the combination of the two, that determines the importance of a variable. For example, despite the high uncertainty in seed dispersal, switchgrass yield is insensitive to this parameter (Figures 6, 7), therefore a better understanding of switchgrass seed dispersal would not reduce model uncertainty. By contrast, although uncertainty in the growth respiration is not particularly large, switchgrass yield was very sensitive to growth respiration, and for this reason growth respiration is the greatest contributor to model uncertainty. In addition, although no seedling mortality data were available, model sensitivity to this parameter was much lower in the posterior compared to prior runs. Using the sensitivity analysis or parameter uncertainties alone would thus lead to incorrect conclusions about what parameters are most important and an

inefficient approach to reducing predictive uncertainties.

This analysis only represents the first step toward more comprehensive accounting of known sources of uncertainty. The next step in reducing uncertainty would be to use the results of the variance decomposition to target the most influential model parameters for further constraint through data collection. We have demonstrated how the use of species-level data to constrain parameter uncertainty reduced ensemble variance, resulting in a new set of targets for additional field observations and refined literature surveys.

Traits that are difficult to measure, such as belowground carbon cycling, can be indirectly constrained with ecosystem-level observations using data assimilation (Luo et al., 2009, 2011). Integrating data assimilation into PEcAn will allow ecosystem-level observations to further constrain parameters for which trait level observations are difficult or impossible to obtain. To date most Bayesian data assimilation approaches applied by ecologists have employed flat, uninformative priors (assigning equal probability to values over many orders of magnitude) , which has lead to the problems of parameter identifiability and the criticism that model parameters are allowed to take on biologically unrealistic values. The use of the meta-analysis posteriors as priors in the data assimilation step ensures that any parameter estimates are consistent with what is known about different plant traits. In this way Bayesian methods are, in effect, updating the literature-derived estimates with new data and providing a coherent and rigorous framework for combining multiple different types of data.

In addition, by effectively restricting parameter space based on observed values, the use of informed priors in data assimilation reduces problems of equifinality and identifiability. This is consistent with the argument by Beven and Freer (2001) that only the feasible parameter range should be sampled.

To a first order the spline interpolations of the univariate relationships between parameters and aboveground biomass (Figure 6) provide a good estimate of the total

model variance. The closure term (Table 3b) accounted for approximately 5.2Mg ha⁻¹ or 28% of the variance decomposition estimate (18.1 Mg ha⁻¹, Table 3a), suggesting that while parameter interactions are important, univariate parameter uncertainty drives overall model variance. One key concern of parameter interactions is that the combination of the variance decomposition terms would result in the prediction of negative yields, which is clearly biologically impossible. By comparing the spline ensemble, where this term is truncated, to the spline-based variance decomposition we can conclude that this truncation effect accounts for 4.1 or 80% of the closure term in the variance decomposition.

By contrast, evaluating the spline ensemble for different ensemble sizes shows that ensemble size had negligible effect on the mean variance estimate although it does improve the precision of this estimate (Table 3a). Finally, the difference between the model and spline ensembles (Table 3b) suggests that the absence of parameter interactions in the variance decomposition account for the remaining 20% of the closure term (< 6% of the variance decomposition calculation), which could be improved by a multivariate meta-analysis and sensitivity analysis, both of which are planned for future development of PEcAn. Overall, the closure term is relatively well constrained even when the parameter interactions are assumed to be linear.

Model-field work feedback

Variance decomposition can be used to inform data collection by identifying candidate parameters for further refinement based on their contributions to model variance. Recall that this variance contribution is a function of parameter sensitivity and the parameters' probability density (equation 6, Figure 7). Sensitivity is a function of the model and so there is no direct way to reduce sensitivity. However, because $Var(f) \propto Var(\beta_0)$, it is possible to reduce the model uncertainty by reducing parameter variances.

842 Through simple power analyses one can explicitly estimate the relationship between an
843 increase in sample size and the reduction in posterior variance. Not only can we calculate
844 the reduction in parameter uncertainty that would be expected for a given sample size,
845 but using equation 6 we can also express this in terms of reductions in the variance of the
846 model output. This then allows us to directly compare the value of different data types in
847 a common currency.

848 Quantitatively evaluating the relationship between data and model uncertainty provides a
849 path of communication between field research and modeling, opening the door for a new
850 framework in which modeling and field work can be mutually informative. Given the
851 current data and model uncertainties, it is possible to identify effective data acquisition
852 strategies based on this analysis. For example, data could be ranked by the ratio of
853 reduction in model uncertainty to the cost of acquiring each sample in terms of dollars
854 and/or man hours. In this way, data collection could be optimized in terms of the
855 efficiency at which new information is gained.

856 These approaches close the model-data loop by enabling models to inform data collection,
857 and data to inform models. Such a shift has the potential to put field ecologists and
858 modelers in closer connection. It also gives us the tools to answer the long standing
859 question among many field ecologists about what exactly modelers need to know. Indeed,
860 this shift highlights a need for greater accessibility to models by the general research
861 community so that field ecologists can drive this loop directly. This is exactly the
862 objective of PEcAn – to encapsulate these tasks in a way that makes the analysis of
863 models transparent, repeatable, and accessible.

864 In addition to informing sample size, the parameter meta-analysis can inform
865 experimental design by providing valuable information on the scales of variability. For
866 example, when data from multiple sites is available, the meta-analysis partitions among
867 site and within site variance. This information can be used to construct optimal sampling

designs which balance intensive vs extensive sampling, and may help identify environmental covariates that should be measured in order to explain parameter variability.

Based on our switchgrass example, variance partitioning also highlights broad data needs and the discrepancy between the relative ease of parameterizing aboveground processes compared to below ground processes. Indeed, one of the greatest challenges in ecosystem ecology is the ability to constrain below ground processes such as root allocation, respiration, and turnover. These parameters are uncertain precisely because measurement is difficult, often indirect, and data may reflect the diverse methods used to indirectly estimate a pool or flux. Many parameters in the ED2 model correspond to processes that are not directly observable. For example, the root respiration parameter in ED2 is not total root respiration but just maintenance respiration, while measurements typically can not separate growth, maintenance, and rhizosphere respiration. Whole-plant growth respiration, which is currently the most important model parameter, is also difficult to estimate directly from measurements (Amthor, 2000). In this case, data assimilation is likely the most efficient route to constrain this parameter; data assimilation would effectively use mass balance of ecosystem carbon exchange to estimate this respiration parameter once other parameters are more directly constrained by data.

Future Directions

The analyses presented here represent the first phase in the development of the PEcAn project. In the near term we intend to expand the existing analyses to include a multivariate meta-analysis and sensitivity analysis to reduce model uncertainty by accounting for parameter covariances. In addition, we plan to implement the power analyses discussed above to more quantitatively inform data collection. A data assimilation module is in progress for PEcAn that will allow the use of ecosystem level

data including plot-level inventory data, eddy covariance fluxes, and remote sensing imagery to enter the analysis and provide additional constraint on uncertainty in both parameters and output. The basic concept of variance decomposition will also be expanded to investigate other sources of variability, such as uncertainty in initial conditions or in driver data. We are implementing ecosystem models other than ED2 within the PEcAn workflow. This will provide opportunities for multi-model ensemble forecasting and assessing data requirements across models. Integrating modeling into a workflow system has distinct advantages not just for model analysis but also for managing the flows of information coming in and out of the model. In this sense we also envision PEcAn as an informatics and data management tool. Finally, it is our hope that other researchers will find PEcAn useful and contribute modules that extend the functionality of the system in creative and exciting ways.

Conclusion

In this paper, we demonstrate an approach to the parametrization of a terrestrial biosphere model that synthesizes available data while providing a robust accounting of parameter uncertainty. We also present a scientific workflow that enables more efficient constraint of this uncertainty by identifying the key areas requiring data collection and model refinement. By quantifying the effect that each parameter has on model output uncertainty, this analysis identifies additional data that, once obtained, would improve model precision. In addition, the analysis calculates probabilities of alternate potential outcomes, resulting in more useful assessments.

Acknowledgments

Author contributions: DSL and MCD developed the concept, statistical analyses, and writing; DSL implemented analyses and visualization; MCD provided scientific and statistical guidance; DW provided feedback on design and implementation of BETY, PEcAn, and this manuscript; KTR collected stomatal slope data; CDD contributed code. Other contributions: Michael Sterling Burns assisted the writing process. Patrick Mulrooney assisted with database design and implementation. John Skillman provided raw data; John Skillman, Peter Reich, Ian Wright, and Stan Wullschleger provided assistance with interpretation of their published data. DSL, DW, David Bettinardi, Andy Tu, Xiaohui Feng, and others entered previously published data into the database and contributed to the design of the web-based data entry workflow interface. Andrew Leakey, Evan Delucia, and Don Ort assisted with prior estimation. Statistical and programming advice was obtained from contributors to crossvalidated.com, stackoverflow.com, and other stackexchange websites. We appreciate all individuals, including but not limited to the authors of publications cited in this paper, whose time, efforts, and intellect designed and produced the data used here, and whose recognition of the importance of archiving and sharing large data sets contributes to the advancement of science.

Funding for this research was provided by the Energy Biosciences Institute and an NSF Advances in Bioinformatics grant #1062547 to MCD.

References

Albert, C. H., F. Grassein, F. M. Schurr, G. Vieilledent, and C. Violle. 2011. When and how should intraspecific variability be considered in trait-based plant ecology? *Perspectives in Plant Ecology, Evolution and Systematics* **13**:217–225. URL <http://linkinghub.elsevier.com/retrieve/pii/S143383191100028X>.

938 Amthor, J. 2000. The McCree-de Wit-Penning de Vries-Thornley Respiration Paradigms:
 939 30 Years Later. *Annals of Botany* **86**:1–20. URL
 940 <http://linkinghub.elsevier.com/retrieve/doi/10.1006/anbo.2000.1175>.

941 Angel, J., 2010. Illinois State Climatologist Data for Station 118749 (Urbana). URL
 942 <http://www.isws.illinois.edu/atmos/statecli/Summary/118740.htm>.

943 Ball, J., I. Woodrow, and J. Berry, 1987. A model predicting stomatal conductance and
 944 its contribution to the control of photosynthesis under different environmental
 945 conditions. Pages 221—224 *in* J. Biggins, editor. *Progress in Photosynthesis Research*.
 946 Martinus Nijhoff Publishers, Netherlands.

947 Bellocchi, G., M. Rivington, M. Donatelli, and K. Matthews. 2010. Validation of
 948 biophysical models: issues and methodologies. A review. *Agronomy for Sustainable*
 949 *Development* **30**:109–130. URL
 950 <http://www.springerlink.com/index/10.1051/agro/2009001>.

951 Beven, K. 2006. A manifesto for the equifinality thesis. *Journal of Hydrology* **320**:18–36.
 952 URL <http://linkinghub.elsevier.com/retrieve/pii/S002216940500332X>.

953 Beven, K., and J. Freer. 2001. Equifinality, data assimilation, and uncertainty estimation
 954 in mechanistic modelling of complex environmental systems using the GLUE
 955 methodology. *Journal of Hydrology* **249**:11–29. URL
 956 <http://linkinghub.elsevier.com/retrieve/pii/S0022169401004218>.

957 Breza, L. C., L. Souza, N. J. Sanders, and A. T. Classen. 2012. Within and between
 958 population variation in plant traits predicts ecosystem
 959 functions associated with a dominant plant species. *Ecology and Evolution* **1**:no–no. URL
 960 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3287307&tool=pmcentrez&ren](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3287307&tool=pmcentrez&rendition=full)
 961 <http://doi.wiley.com/10.1002/ece3.223>.

962 Cariboni, J., D. Gatelli, R. Liska, and A. Saltelli. 2007. The role of sensitivity analysis in
 963 ecological modelling. *Ecological Modelling* **203**:167–182. URL
 964 <http://linkinghub.elsevier.com/retrieve/pii/S0304380006005734>.

965 Casella, G., and R. L. Berger. 2001. *Statistical Inference*. second edition. Duxbury Press,
 966 Pacific Grove, CA.

967 Chapin III, F. S., P. A. Matson, and H. A. Mooney. 2002. *Principles of Terrestrial*
 968 *Ecosystem Ecology*. Springer.

969 Clark, J. S., S. R. Carpenter, M. Barber, S. Collins, A. Dobson, J. a. Foley, D. M. Lodge,
 970 M. Pascual, R. Pielke, W. Pizer, C. Pringle, W. V. Reid, K. a. Rose, O. Sala, W. H.
 971 Schlesinger, D. H. Wall, and D. Wear. 2001. Ecological forecasts: an emerging
 972 imperative. *Science* (New York, N.Y.) **293**:657–60. URL
 973 <http://www.ncbi.nlm.nih.gov/pubmed/11474103>.

974 Clark, J. S., M. Silman, R. Kern, E. Macklin, and J. HilleRisLambers. 1999. Seed
 975 dispersal near and far: patterns across temperate and tropical forests. *Ecology*
 976 **80**:1475–1494. URL [http://www.esajournals.org/doi/full/10.1890/0012-](http://www.esajournals.org/doi/full/10.1890/0012-9658(1999)080[1475:SDNAFP]2.0.CO;2)
 977 [9658\(1999\)080\[1475:SDNAFP\]2.0.CO;2](http://www.esajournals.org/doi/full/10.1890/0012-9658(1999)080[1475:SDNAFP]2.0.CO;2).

978 Collatz, G., M. Ribas-Carbo, and J. Berry. 1992. Coupled Photosynthesis-Stomatal
 979 Conductance Model for Leaves of C4 Plants. *Functional Plant Biology* **19**:519—538.
 980 URL <http://www.publish.csiro.au/paper/PP9920519>.

981 Denman, K., G. Brasseur, A. Chidthaisong, P. Ciais, P. Cox, R. Dickinson,
 982 D. Hauglustaine, C. Heinze, E. Holland, D. Jacob, U. Lohmann, S. Ramachandran,
 983 P. D. S. Dias, S. W. Lohmann, S. Ramachandran, P.L. Da Silva Dias, X. Zhang,
 984 X. Zhang, S. R. Lohmann, P. da Silva Dias, S. Wofsy, and X. Zhang, 2007. Couplings
 985 between changes in the climate system and biogeochemistry. 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. Chapter 7, pages 499–587 . Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA. URL http://www.osti.gov/bridge/product.biblio.jsp?osti_id=934721.

Dietze, M. C., and A. M. Latimer, 2011. Forest simulators. Pages 307–316 *in* A. Hastings and L. Gross, editors. Sourcebook in Theoretical Ecology. Univ California Press, Berkeley.

Dietze, M. C., R. Vargas, A. D. Richardson, P. C. Stoy, A. G. Barr, R. S. Anderson, M. A. Arain, I. T. Baker, T. A. Black, J. M. Chen, P. Ciais, L. B. Flanagan, C. M. Gough, R. F. Grant, D. Hollinger, R. C. Izaurralde, C. J. Kucharik, P. Laflleur, S. Liu, E. Lokupitiya, Y. Luo, J. W. Munger, C. Peng, B. Poulter, D. T. Price, D. M. Ricciuto, W. J. Riley, A. K. Sahoo, K. Schaefer, A. E. Suyker, H. Tian, C. Tonitto, H. Verbeeck, S. B. Verma, W. Wang, and E. Weng. 2011. Characterizing the performance of ecosystem models across time scales: A spectral analysis of the North American Carbon Program site-level synthesis. *Journal of Geophysical Research* **116**. URL <http://www.agu.org/pubs/crossref/2011/2011JG001661.shtml>.

Ellison, A. M. 2010. Repeatability and transparency in ecological research. *Ecology* **91**:2536–2539. URL <http://www.esajournals.org/doi/abs/10.1890/09-0032.1>.

Ernst, W. H. O., E. M. Veenendaal, and M. M. Kebakile. 1992. Possibilities for dispersal in annual and perennial grasses in a savanna in Botswana. *Vegetatio* **102**:1–11. URL <http://www.springerlink.com/index/10.1007/BF00031700>.

Farquhar, G. D., and T. D. Sharkey. 1982. Stomatal Conductance and Photosynthesis. *Annual Review of Plant Physiology* **33**:317–345. URL <http://www.annualreviews.org/doi/abs/10.1146/annurev.pp.33.060182.001533>.

1010 Gelman, A., and D. B. Rubin. 1992. Inference from Iterative Simulation Using Multiple
 1011 Sequences. *Statistical Science* **7**:457–472. URL
 1012 <http://projecteuclid.org/euclid.ss/1177011136>.

1013 Gelman, A., and K. Shirley, 2011. Inference from simulations and monitoring
 1014 convergence. Pages 163—174 *in* S. Brooks, A. Gelman, G. Jones, and X.-L. Meng,
 1015 editors. *Handbook of Markov Chain Monte Carlo*. CRC Press LLC.

1016 Gill, R. A., and R. B. Jackson. 2000. Global Patterns of Root Turnover for Terrestrial
 1017 Ecosystems. *New Phytologist* **147**:13–31. URL
 1018 <http://www.jstor.org/stable/2588686>.

1019 Heaton, E. A., F. G. Dohleman, and S. P. Long. 2008. Meeting US biofuel goals with less
 1020 land: the potential of *Miscanthus*. *Global Change Biology* **14**:2000–2014. URL
 1021 <http://www3.interscience.wiley.com/journal/120119109/abstract>.

1022 Jackson, R. B., H. A. Mooney, and E. D. Schulze. 1997. A global budget for fine root
 1023 biomass, surface area, and nutrient contents. *Proceedings of the National Academy of*
 1024 *Sciences of the United States of America* **94**:7362–6. URL
 1025 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=23826>.

1026 Jones, M. B., M. P. Schildhauer, O. Reichman, and S. Bowers. 2006. The New
 1027 Bioinformatics: Integrating Ecological Data from the Gene to the Biosphere. *Annual*
 1028 *Review of Ecology, Evolution, and Systematics* **37**:519–544. URL
 1029 <http://www.annualreviews.org/doi/abs/10.1146/annurev.ecolsys.37.091305.110031>.

1030 Jongejans, E., and P. Schippers. 1999. Modeling seed dispersal by wind in herbaceous
 1031 species. *Oikos* **87**:362–372. URL <http://www.jstor.org/stable/3546752>.

1032 Justus, C., W. Hargraves, A. Mikhail, and D. Graber. 1978. Methods for estimating wind

1033 speed frequency distributions. *J. Appl. Meteorol.* **17**:350–353. URL
1034 http://www.osti.gov/energycitations/product.biblio.jsp?osti_id=5127748.

1035 Kubien, D. S., and R. F. Sage. 2004. Low-temperature photosynthetic performance of a
1036 C4 grass and a co-occurring C3 grass native to high latitudes. *Plant, Cell and*
1037 *Environment* **27**:907–916. URL
1038 <http://doi.wiley.com/10.1111/j.1365-3040.2004.01196.x>.

1039 Kucharik, C. J., J. A. Foley, C. Delire, V. A. Fisher, M. T. Coe, J. D. Lenters,
1040 C. Young-Molling, N. Ramankutty, J. M. Norman, and S. T. Gower. 2000. Testing the
1041 performance of a dynamic global ecosystem model: Water balance, carbon balance, and
1042 vegetation structure. *Global Biogeochemical Cycles* **14**:795. URL
1043 <http://www.agu.org/pubs/crossref/2000/1999GB001138.shtml>.

1044 Larocque, G. R., J. S. Bhatti, R. Boutin, and O. Chertov. 2008. Uncertainty analysis in
1045 carbon cycle models of forest ecosystems: Research needs and development of a
1046 theoretical framework to estimate error propagation. *Ecological Modelling*
1047 **219**:400–412. URL
1048 <http://linkinghub.elsevier.com/retrieve/pii/S0304380008003542>.

1049 Leuning, R. 1995. A critical appraisal of a combined stomatal-photosynthesis model for
1050 C3 plants. *Plant, Cell and Environment* **18**:339–355. URL
1051 <http://doi.wiley.com/10.1111/j.1365-3040.1995.tb00370.x>.

1052 Luo, Y., K. Ogle, C. Tucker, S. Fei, C. Gao, S. LaDeau, J. S. Clark, and D. S. Schimel.
1053 2011. Ecological forecasting and data assimilation in a data-rich era. *Ecological*
1054 *applications* : a publication of the Ecological Society of America **21**:1429–42. URL
1055 <http://www.ncbi.nlm.nih.gov/pubmed/21830693>.

1056 Luo, Y., E. Weng, X. Wu, C. Gao, X. Zhou, and L. Zhang. 2009. Parameter

1057 identifiability, constraint, and equifinality in data assimilation with ecosystem models.
1058 Ecological applications : a publication of the Ecological Society of America **19**:571–4.
1059 URL <http://www.ncbi.nlm.nih.gov/pubmed/19425417>.

1060 Massad, R.-S., A. Tuzet, and O. Bethenod. 2007. The effect of temperature on C(4)-type
1061 leaf photosynthesis parameters. *Plant, cell & environment* **30**:1191–204. URL
1062 <http://www.ncbi.nlm.nih.gov/pubmed/17661755>.

1063 McLaughlin, S., and L. Kszos. 2005. Development of switchgrass (*Panicum virgatum*) as
1064 a bioenergy feedstock in the United States. *Biomass and Bioenergy* **28**:515–535. URL
1065 <http://linkinghub.elsevier.com/retrieve/pii/S0961953404001904>.

1066 McMahon, S. M., M. C. Dietze, M. H. Hersh, E. V. Moran, and J. S. Clark. 2009. A
1067 predictive framework to understand forest responses to global change. *Annals of the*
1068 *New York Academy of Sciences* **1162**:221–36. URL
1069 <http://www.ncbi.nlm.nih.gov/pubmed/19432650>.

1070 Medvigy, D., S. C. Wofsy, J. W. Munger, D. Y. Hollinger, and P. R. Moorcroft. 2009.
1071 Mechanistic scaling of ecosystem function and dynamics in space and time: Ecosystem
1072 Demography model version 2. *Journal of Geophysical Research* **114**:1–21. URL
1073 <http://www.agu.org/pubs/crossref/2009/2008JG000812.shtml>.

1074 Mesinger, F., G. DiMego, E. Kalnay, K. Mitchell, P. C. Shafran, W. Ebisuzaki, D. Jović,
1075 J. Woollen, E. Rogers, E. H. Berbery, M. B. Ek, Y. Fan, R. Grumbine, W. Higgins,
1076 H. Li, Y. Lin, G. Manikin, D. Parrish, and W. Shi. 2006. North American Regional
1077 Reanalysis. *Bulletin of the American Meteorological Society* **87**:343. URL
1078 <http://journals.ametsoc.org/doi/abs/10.1175/BAMS-87-3-343>.

1079 Michener, W. K., and M. B. Jones. 2012. Ecoinformatics: supporting ecology as a

data-intensive science. *Trends in Ecology & Evolution* **27**:85–93. URL
<http://linkinghub.elsevier.com/retrieve/pii/S0169534711003399>.

Miguez, F., X. Zhu, S. Humphries, G. Bollero, and S. Long. 2009. A semimechanistic
model predicting the growth and production of the bioenergy crop *Miscanthus*
giganteus: description, parameterization and validation. *GCB Bioenergy* **1**:282–296.
URL <http://doi.wiley.com/10.1111/j.1757-1707.2009.01019.x>.

Mood, A. M., F. A. Graybill, and D. C. Boes. 1974. *Introduction to the Theory of*
Statistics. 3rd edition. McGraw-Hill College.

Moorcroft, P. R., G. C. Hurtt, and S. W. Pacala. 2001. A Method for Scaling Vegetation
Dynamics: the Ecosystem Demography Model (ED). *Ecological Monographs*
71:557–586. URL <http://www.esajournals.org/doi/abs/10.1890/0012-9615%282001%29071%5B0557%3AAMFSVD%5D2.0.CO%3B2>.

Moore, D. J., J. Hu, W. J. Sacks, D. S. Schimel, and R. K. Monson. 2008. Estimating
transpiration and the sensitivity of carbon uptake to water availability in a subalpine
forest using a simple ecosystem process model informed by measured net CO₂ and
H₂O fluxes. *Agricultural and Forest Meteorology* **148**:1467–1477. URL
<http://linkinghub.elsevier.com/retrieve/pii/S0168192308001263>.

Oehlert, G. 1992. A note on the delta method. *American Statistician* **46**:27–29. URL
<http://www.jstor.org/stable/10.2307/2684406>.

Oyarzabal, M., J. M. Paruelo, F. Pino, M. Oesterheld, and W. K. Lauenroth. 2008. Trait
differences between grass species along a climatic gradient in South and North
America. *Journal of Vegetation Science* **19**:183–192. URL
<http://blackwell-synergy.com/doi/abs/10.3170/2007-8-18349>.

1103 Parton, W. J., M. Hartman, D. Ojima, and D. Schimel. 1998. DAYCENT and its land
 1104 surface submodel: description and testing. *Global and Planetary Change* **19**:35–48.
 1105 URL <http://linkinghub.elsevier.com/retrieve/pii/S092181819800040X>.

1106 Plummer, M., 2010. JAGS Version 2.2.0 user manual.

1107 Reich, P. B., and J. Oleksyn. 2004. Global patterns of plant leaf N and P in relation to
 1108 temperature and latitude. *Proceedings of the National Academy of Sciences of the*
 1109 *United States of America* **101**:11001–6. URL
 1110 <http://www.pnas.org/cgi/content/abstract/101/30/11001>.

1111 Reichstein, M., J. Tenhunen, O. Roupsard, J.-M. Ourcival, S. Rambal, F. Miglietta,
 1112 A. Peressotti, M. Pecchiari, G. Tirone, and R. Valentini. 2003. Inverse modeling of
 1113 seasonal drought effects on canopy CO₂ /H₂O exchange in three Mediterranean
 1114 ecosystems. *Journal of Geophysical Research* **108**. URL
 1115 <http://www.agu.org/pubs/crossref/2003/2003JD003430.shtml>.

1116 Richardson, A., and D. Hollinger. 2005. Statistical modeling of ecosystem respiration
 1117 using eddy covariance data: Maximum likelihood parameter estimation, and Monte
 1118 Carlo simulation of model and parameter uncertainty, applied to three simple models.
 1119 *Agricultural and Forest Meteorology* **131**:191–208. URL
 1120 <http://linkinghub.elsevier.com/retrieve/pii/S0168192305001139>.

1121 Saltelli, A., M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana,
 1122 and S. Tarantola. 2008. *Global Sensitivity Analysis*. John Wiley & Sons, West Sussex,
 1123 England.

1124 Saugier, B., J. Roy, and H. Mooney, 2001. Estimations of global terrestrial productivity:
 1125 Converging toward a single number? Chapter 23, pages 543—557 *in* B. Saugier and
 1126 H. Mooney, editors. *Terrestrial Global Productivity*. Academic Press, San Diego, CA.

1127 Schwalm, C. R., C. a. Williams, K. Schaefer, R. Anderson, M. A. Arain, I. Baker,
 1128 A. Barr, T. A. Black, G. Chen, J. M. Chen, P. Ciais, K. J. Davis, A. Desai, M. Dietze,
 1129 D. Dragoni, M. L. Fischer, L. B. Flanagan, R. Grant, L. Gu, D. Hollinger, R. C.
 1130 Izaurrealde, C. Kucharik, P. Lafleur, B. E. Law, L. Li, Z. Li, S. Liu, E. Lokupitiya,
 1131 Y. Luo, S. Ma, H. Margolis, R. Matamala, H. McCaughey, R. K. Monson, W. C.
 1132 Oechel, C. Peng, B. Poulter, D. T. Price, D. M. Riciutto, W. Riley, A. K. Sahoo,
 1133 M. Sprintsin, J. Sun, H. Tian, C. Tonitto, H. Verbeeck, and S. B. Verma. 2010. A
 1134 model-data intercomparison of CO₂ exchange across North America: Results from the
 1135 North American Carbon Program site synthesis. *Journal of Geophysical Research* **115**.
 1136 URL <http://www.agu.org/pubs/crossref/2010/2009JG001229.shtml>.

 1137 Skillman, J. B. 2008. Quantum yield variation across the three pathways of
 1138 photosynthesis: not yet out of the dark. *Journal of experimental botany* **59**:1647–61.
 1139 URL <http://www.ncbi.nlm.nih.gov/pubmed/18359752>.

 1140 Stodden, V., D. Donoho, S. Fomel, M. Friedlander, M. Gerstein, R. LeVeque, I. Mitchell,
 1141 L. Larrimore, C. Wiggins, N. W. Bramble, P. O. Brown, V. J. Carey, L. DeNardis,
 1142 R. Gentleman, J. D. Gezelter, A. Goodman, M. G. Knepley, J. E. Moore, F. A.
 1143 Pasquale, J. Rolnick, M. Seringhaus, and R. Subramanian, 2010. Reproducible
 1144 Research. URL <http://doi.ieeecomputersociety.org/10.1109/MCSE.2010.113>
 1145 <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5562471>.

 1146 Thompson, S., and G. Katul. 2008. Plant propagation fronts and wind dispersal: an
 1147 analytical model to upscale from seconds to decades using superstatistics. *The*
 1148 *American naturalist* **171**:468–79. URL
 1149 <http://www.ncbi.nlm.nih.gov/pubmed/18248297>.

 1150 Tjoelker, M. G., J. M. Craine, D. Wedin, P. B. Reich, and D. Tilman. 2005. Linking leaf
 1151 and root trait syndromes among 39 grassland and savannah species. *The New*

1152 phytologist **167**:493–508. URL
 1153 <http://onlinelibrary.wiley.com/doi/10.1111/j.1469-8137.2005.01428.x/full>.

1154 USDA, and NRCS, 2011. The PLANTS Database. URL <http://plants.usda.gov/>.

1155 VanLoocke, A., T. E. Twine, M. Zeri, and C. J. Bernacchi. 2012. A regional comparison
 1156 of water use efficiency for miscanthus, switchgrass and maize. Agricultural and Forest
 1157 Meteorology **164**:82–95. URL
 1158 <http://linkinghub.elsevier.com/retrieve/pii/S0168192312001931>.

1159 Violle, C., B. J. Enquist, B. J. McGill, L. Jiang, C. H. Albert, C. Hulshof, V. Jung, and
 1160 J. Messier. 2012. The return of the variance: intraspecific variability in community
 1161 ecology. Trends in Ecology & Evolution **27**:244–52. URL
 1162 <http://www.ncbi.nlm.nih.gov/pubmed/22244797>.

1163 Wang, D., D. S. LeBauer, and M. C. Dietze. 2010. A quantitative review comparing the
 1164 yield of switchgrass in monocultures and mixtures in relation to climate and
 1165 management factors. GCB Bioenergy **2**:16–25. URL
 1166 <http://blackwell-synergy.com/doi/abs/10.1111/j.1757-1707.2010.01035.x>.

1167 Wang, D., M. W. Maughan, J. Sun, X. Feng, F. E. Miguez, D. K. Lee, and M. C. Dietze.
 1168 2011. Impacts of canopy position and nitrogen on nitrogen allocation and
 1169 photosynthesis of switchgrass (*Panicum virgatum* L.). Aspects of Applied Biology 112
 1170 **112**:341—351.

1171 Williams, M., A. D. Richardson, M. Reichstein, P. C. Stoy, P. Peylin, H. Verbeeck,
 1172 N. Carvalhais, M. Jung, D. Y. Hollinger, J. Kattge, R. Leuning, Y. Luo, E. Tomelleri,
 1173 C. Trudinger, and Y.-P. Wang. 2009. Improving land surface models with FLUXNET
 1174 data. Biogeosciences Discussions **6**:2785–2835. URL
 1175 <http://www.biogeosciences-discuss.net/6/2785/2009/>.

1176 Wright, I. J., P. B. Reich, M. Westoby, D. D. Ackerly, Z. Baruch, F. Bongers,
1177 J. Cavender-Bares, T. Chapin, J. H. C. Cornelissen, M. Diemer, J. Flexas, E. Garnier,
1178 P. K. Groom, J. Gulias, K. Hikosaka, B. B. Lamont, T. Lee, W. Lee, C. Lusk, J. J.
1179 Midgley, M.-L. Navas, U. Niinemets, J. Oleksyn, N. Osada, H. Poorter, P. Poot,
1180 L. Prior, V. I. Pyankov, C. Roumet, S. C. Thomas, M. G. Tjoelker, E. J. Veneklaas,
1181 and R. Villar. 2004. The worldwide leaf economics spectrum. *Nature* **428**:821–7. URL
1182 <http://www.ncbi.nlm.nih.gov/pubmed/15103368>.

1183 Wullschleger, S. D. 1993. Biochemical Limitations to Carbon Assimilation in C3 Plants -
1184 Retrospective Analysis of the A/Ci Curves from 109 Species. *Journal of experimental*
1185 *botany* **44**:907–920.

List of Tables

1	Prior Distributions Prior distributions used in meta-analysis and model parameterization. Prior sources come from citations as indicated except * by authors or † based on default ED2 parameterizations, as described in text. The 'Clade' column indicates the group of plants for which the priors were derived.	54
2	Meta-analysis Results Results of meta-analysis of Switchgrass data for six physiological traits. The global mean parameter, β_0 , is used to parametrize the Ecosystem Demography model and is described in more detail by Figure 4. The variance components are transformed from precision to the standard deviation scale for ease of interpretation. Values are reported as the parameter median with the 95% CI in parentheses. Units are the same as in Table 1.	55

1199	3	a) Variance Estimates	Comparison of sample variances (s , on stan-	
1200			dard deviation scale) for the aboveground biomass estimated from data-	
1201			constrained parameters calculated from model ensemble, spline-emulated	
1202			model ensembles, and variance decomposition. Values in parentheses are	
1203			estimates of uncertainty in the sample estimate of variance. Sample size, n ,	
1204			refers to the size of the sample from the posterior parameter distribution.	
1205		b) Components of closure term, ω	The closure term ω (equation 6) is	
1206			5.2, the difference between the variance decomposition and model ensemble	
1207			estimates of σ . The closure due to parameter interactions is estimated as the	
1208			difference between the spline ensemble and the model ensemble; the closure	
1209			due to the absence of a lower bound of zero on the spline functions is esti-	
1210			mated as the difference between the variance decomposition and the spline	
1211			ensemble estimates. * Analysis of the closure term is based on estimates	
1212			with $n = 10000$ parameter sets, except in the case of the model ensemble	
1213			because evaluation of the model ensemble at $n = 10000$ is computationally	
1214			prohibitive.	56

Parameter	Units	Clade	Distribution	a	b	N	mean	LCL	UCL	Citation
Specific Leaf Area	$\text{m}^2 \text{kg}^{-1}$	Grass	Gamma	2.06	19.00	125	17	3.2	36	(Wright et al., 2004)
Leaf Turnover Rate	1/yr	Grass	Weibull	2.90	0.63	40	4.6	0.91	11	(Wright et al., 2004)
Root Turnover Rate	1/yr	Grass	Gamma	1.67	0.66	66	0.59	0.073	1.4	(Gill and Jackson, 2000)
Quantum Efficiency	percent	C4 grasses	Weibull	90.90	1580.00	56	0.058	0.046	0.07	(Skillman, 2008)
Stomatal Slope	ratio	C4 Grass	Gamma	3.63	3.81	4	3.4	1.4	5.5	*
Vcmax	$\text{umol CO}_2 \text{m}^{-2} \text{s}^{-1}$	graminoid	Gamma	3.49	24.70	97	22	8.6	36	(Wullschleger, 1993)
Leaf Width	mm	C4 Grass	Weibull	26.10	5.94	18	4.4	2.9	6.2	(Oyarzabal et al., 2008)
Root Respiration Rate	$\text{umol CO}_2 \text{kg}^{-1} \text{s}^{-1}$	C4 Grass	F	5.61	2.33	35	5.6	1	10	(Tjoelker et al., 2005)
Fine Root Allocation	ratio	Grass	Beta	0.80	0.81	0	3.1	0.46	11	(Chapin III et al., 2002)
Seed Dispersal	percent	Grass	log-Normal	20.10	74.90	30	0.21	0.14	0.3	(Jongejans and Schippers, 1993)
Photosynthesis min temp	Celsius	C4 Grass	F	10.00	1.02	0	10	8	12	*
Growth Respiration	percent	Grass	log-Normal	2.63	6.52	0	0.29	0.062	0.6	*
Seedling Mortality	percent	monocots	log-Normal	3.61	0.43	0	0.89	0.5	1	*
Mortality Coefficient	1/yr	plants	Weibull	1.47	0.06	0	25	1.8	80	*
Reproductive Allocation	percent	Plants	log-Normal	2.00	4.00	0	0.33	0.053	0.72	*

Table 1

Variable	n	β_0	σ_Y	σ_{site}	$\sigma_{\text{treatment site}}$	$\beta_{\text{greenhouse}}$
Specific Leaf Area	24	16(12, 20)	2.8(2.5, 3.2)	3.2(1.6, 7.3)	2.4(1.1, 6)	6.5(1, 12)
Leaf Width	39	6(4.7, 6.6)	0.46(0.44, 0.48)	0.47(0.2, 2.1)	6.4(1.9, 130)	1.6(-0.033, 3.5)
Vcmax	4	24(18, 30)	12(8.1, 17)		1.2(0.098, 47)	
Fine Root Allocation	4	1.3(0.5, 2.6)	2.2(1.2, 6.2)			
Root Respiration Rate	1	5.1(3.7, 6.6)	1.2(0.39, 2.3)			
Root Turnover Rate	1	0.67(0.2, 1.1)	0.45(0.14, 0.88)			
Stomatal Slope	4	4.1(3.9, 4.3)	0.33(0.23, 0.45)			

Table 2

	model	spline	variance
	ensemble	ensemble	decomposition
n	$s_f(\beta_0)$	$s_g(\beta_0)$	$\sum s_{g_i}(\beta_{0i})$
500	13(14)	13.8(13)	18.2(6)
10000	*	14.1(2.8)	18.1(1.2)

(a)

	calculation	Mg/ha
ω_{total}	$\sum s_{g_i} - s_f$	5.2
$\omega_{\text{covariance}}$	$s_g - s_f$	1.1
$\omega_{\text{truncation}}$	$\sum s_{g_i} - s_g$	4.1

(b)

Table 3

List of Figures

1	Overview of the PEcAn workflow. The synthesis of plant trait data begins by querying a database of plant trait data for data on a single species or a plant functional type, and then mapping these data to the model parameters that they inform. The database also provides probability distributions that describe our prior information about the range of values that a model parameter can take. Next, this information is synthesized in a Bayesian meta-analysis, resulting in a posterior trait distribution that summarizes the uncertainty in each parameter. The ensemble of model runs produces the posterior distribution of model outputs, representing a probabilistic assessment or forecast that accounts for input parameter uncertainty. The final steps in the workflow are the sensitivity analysis and variance decomposition; these steps gives insight into the relative contribution of each parameter to the uncertainty in the model output, and can be used to guide additional data collection that will most efficiently reduce model uncertainty.	62
---	--	----

2 **Prior distributions** PDFs of priors with data constraints. Parameter value
is on the x-axis and probability density is on the y-axis, and the area under
each curve equals one. Three points on each line, from left to right, indicate
the 2.5th, 50th, and 97.5th quantiles. (From top left) Four priors fit to data
(data points shown as rug plot) using maximum likelihood: specific leaf
area and leaf turnover rate (Wright et al., 2004), root turnover rate (Gill
and Jackson, 2000), and quantum yield (Skillman, 2008). Four priors fit to
the posterior predictive distribution of an unobserved C4 grass species using
Bayesian meta-analysis of data from multiple plant functional types (C4
data shown in black, other functional types in grey): stomatal slope (present
study data provided in Appendix A), V_{cmax} of C3 plants (Wullschlegel, 1993)
and C4 grasses (Kubien and Sage, 2004; Massad et al., 2007; Wang et al.,
2011), leaf width (Oyarzabal et al., 2008), and root respiration (Tjoelker
et al., 2005). Priors fit to 95% CI (dashed grey line) and median (solid grey
line) based on ED2 defaults and expert opinion as described in the text:
fine root to leaf ratio (Chapin III et al., 2002), seed dispersal (Ernst et al.
(1992) model parameterized with site level data), minimum temperature
of photosynthesis (Don Ort, personal communication), growth respiration,
seedling mortality factor, mortality factor, and reproductive allocation. . . 63

3 **Overview of the Hierarchical Bayesian meta-analysis model.** For each trait, the posterior estimate of the global trait mean (β_0) is used as an input parameter in the sensitivity analysis and model ensemble (Figures 6 and 5). Results from the meta-analysis of specific leaf area are as an illustrative example; x-axes have units of m^2kg^{-1} and all plots are on the same scale. Each of the k sample means (Y_k) are taken from published articles and unpublished field measurements, and may be associated with a sample standard error and sample size. When sufficient data were available, site, treatment, and greenhouse effects were estimated. The within-unit standard deviation, σ_Y , is estimated from se and n . Site and treatment random effects, represented by β_{site} and $\beta_{\text{tr}|\text{site}}$, are estimated for each site and treatment within site with from normal distributions with mean zero and standard deviations σ_{site} and $\sigma_{\text{tr}|\text{site}}$, respectively. Greenhouse is a fixed effect. Table 2 summarizes the global mean, variance terms, and greenhouse effect for the seven model parameters informed by species-level data. . . . 64

4 **Prior (gray) and posterior (black) densities of trait parameters used in the analysis.**

Priors distributions are based on the traits of plants within broad taxonomic or functional type groupings (e.g. all grasses). When species-level data were available, they are used in a hierarchical Bayesian meta-analysis, and the posterior estimate of the mean parameter value is shown. Data used in the meta-analysis come from both published and direct measurements of the trait on the perennial C4 grass Switchgrass (*Panicum virgatum*). These data are represented as mean \pm SE. Mismatch between data and the posterior estimate of the global trait mean results from site, treatment, and greenhouse effects. Data from plants grown under an experimental treatment or in a controlled environment (e.g. in a pot or greenhouse) are presented in grey; data from field-grown plants under control treatments are in black. Site-level effects account for disparity between raw data and parameter distribution in the SLA and leaf width plots. 65

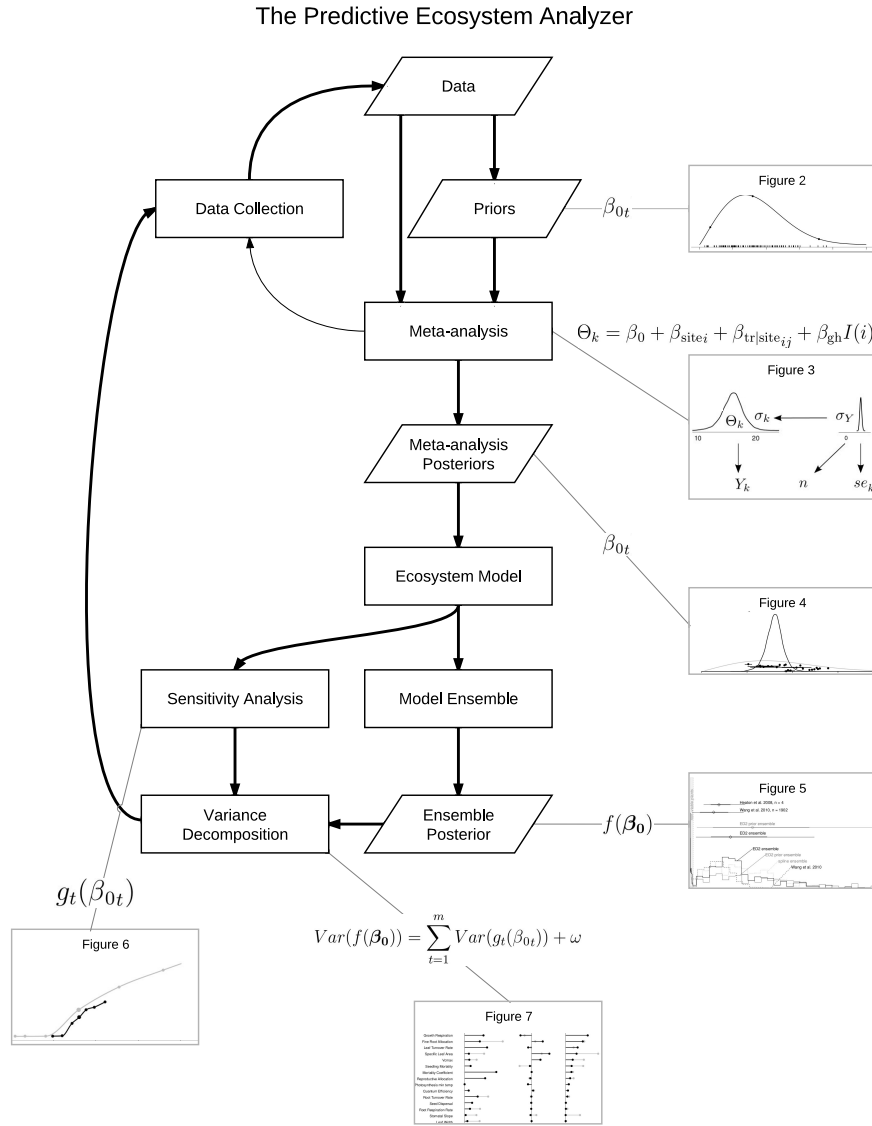
5 **Ensemble average 2004-2006 post-senescence yield.** Histogram of results from prior ensemble runs (dashed), posterior ensemble runs (solid line), and the spline ensemble (gray line). The gray box on the left represents non-viable ensemble members ($\leq 2\text{Mg/ha}$, see text). Horizontal bars provide a summary of yields, from top: a three year trial at the modeled site (Heaton et al., 2008), all 1902 observations included in a recent meta-analysis (Wang et al., 2010), viable runs from the ED2 ensemble based on prior and posterior parameterizations. Diamonds indicate the median; thick and thin lines indicating the 68% and the 95% CI, respectively. Histogram-style plots provide comparison of the distributions of observations and model runs. For clarity, non-viable and viable runs are plotted with different width bins. . 66

6 **Univariate relationships between parameters and 2004-2006 average modeled yield.** Parameter values are on the x-axis and biomass is on

the y-axis while runs centered around the prior median are in gray and those centered around the posterior median are in black. The univariate responses were estimated using a cubic spline to fit model output at the median and $\pm[1, 2, 3]\sigma$ quantiles of each parameter while holding other parameters to the median value. 67

7 **Partitioning of variance by parameter** results from variance decomposition conducted before (grey) and after (black) updating parameter estimates with species-level data in the meta-analysis. From left to right, panels present: a) the uncertainty associated with each parameter (coefficient of variation, $CV = \sigma/\mu$). The degree to which some parameters have been constrained by data is indicated by the reduction in CV in the black compared to the grey bars; sample sizes are indicated in Table 2. b) the sensitivity of modeled aboveground biomass to each parameter presented as elasticity (elasticity is normalized sensitivity, and an elasticity of one indicates that model output will double when the parameter value doubles). Sensitivity is the slope of the line at the median in Figure 6). Parameters with larger bars have greater influence on model output. c) Partial variance is the contribution of each parameter to explained variance. This is a function of both the parameter variance and sensitivity. Parameters with both large CV and elasticity contribute the most to uncertainty in model output. 68

Figure 1



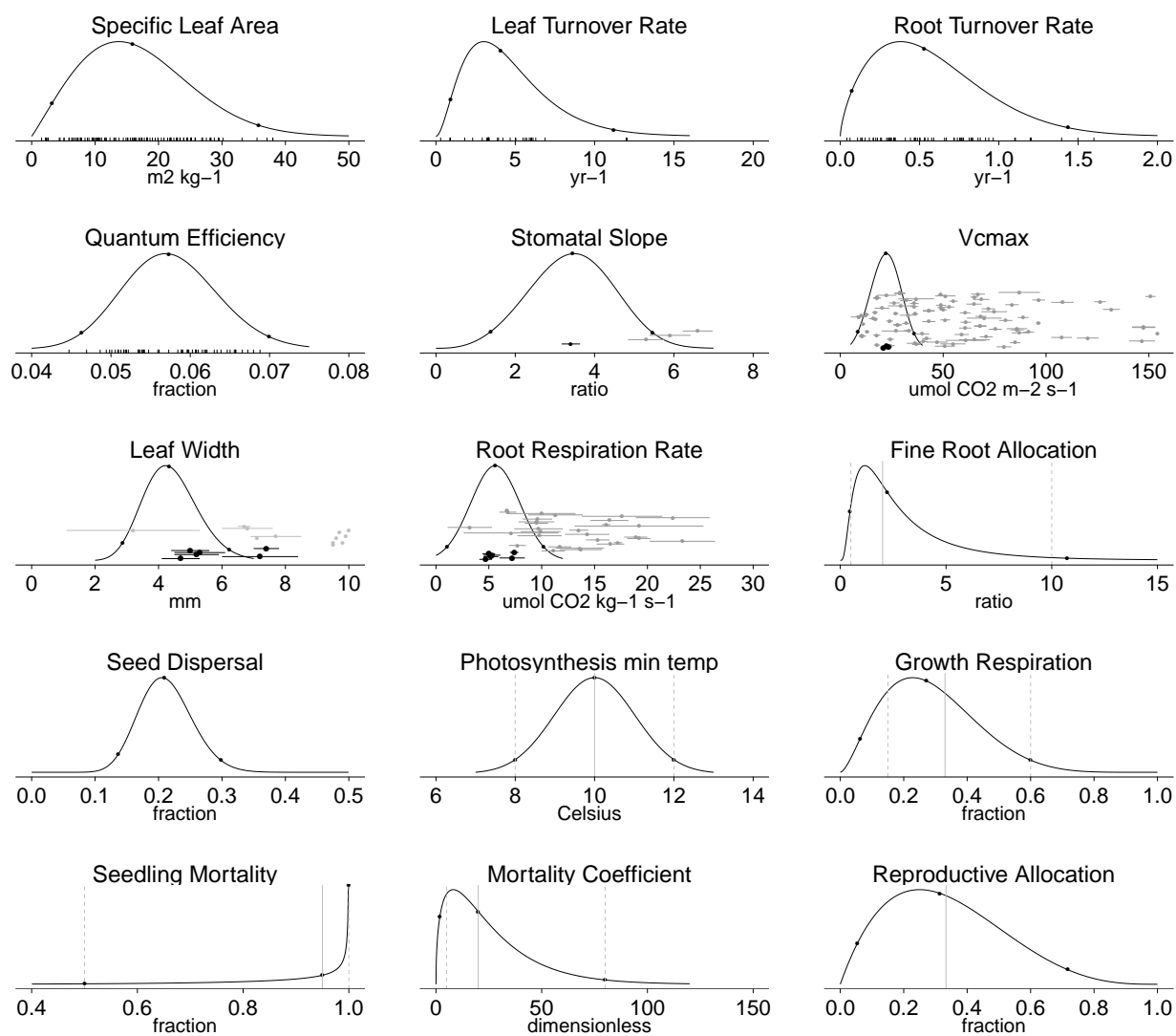


Figure 2

Figure 3

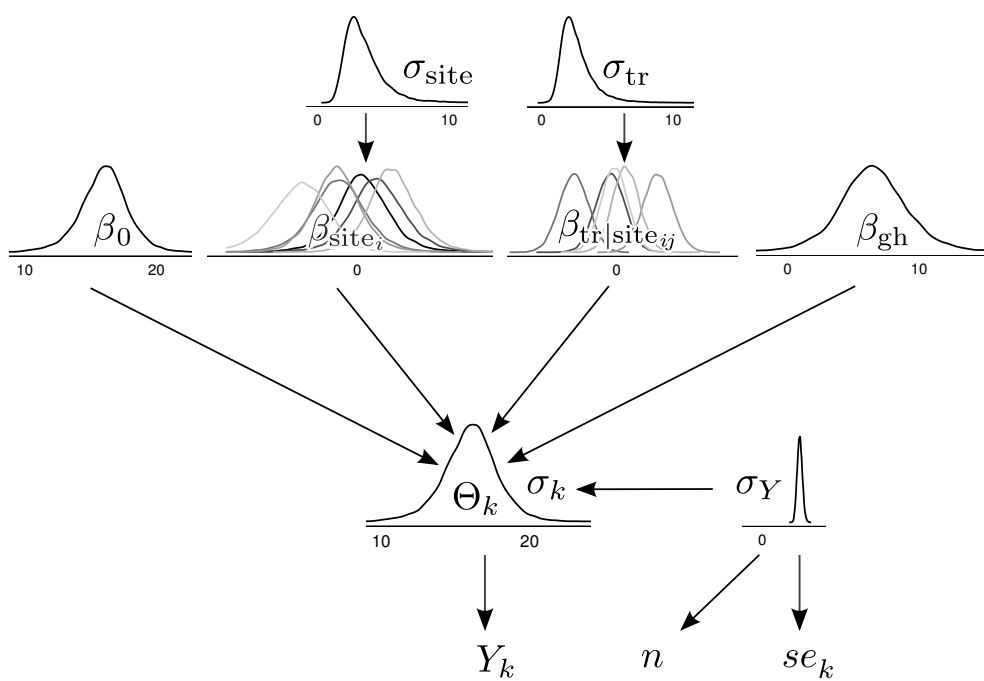


Figure 4

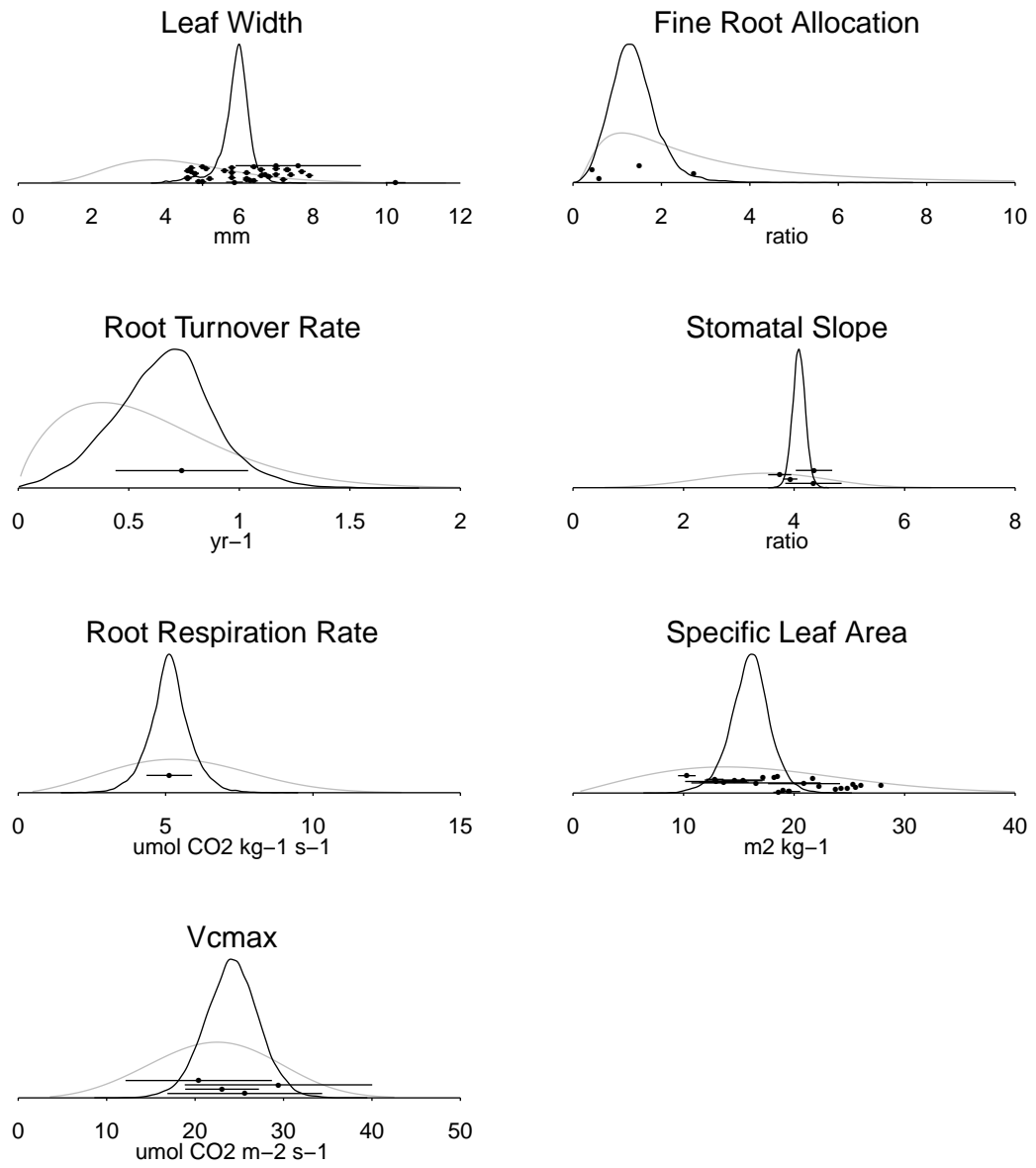


Figure 5

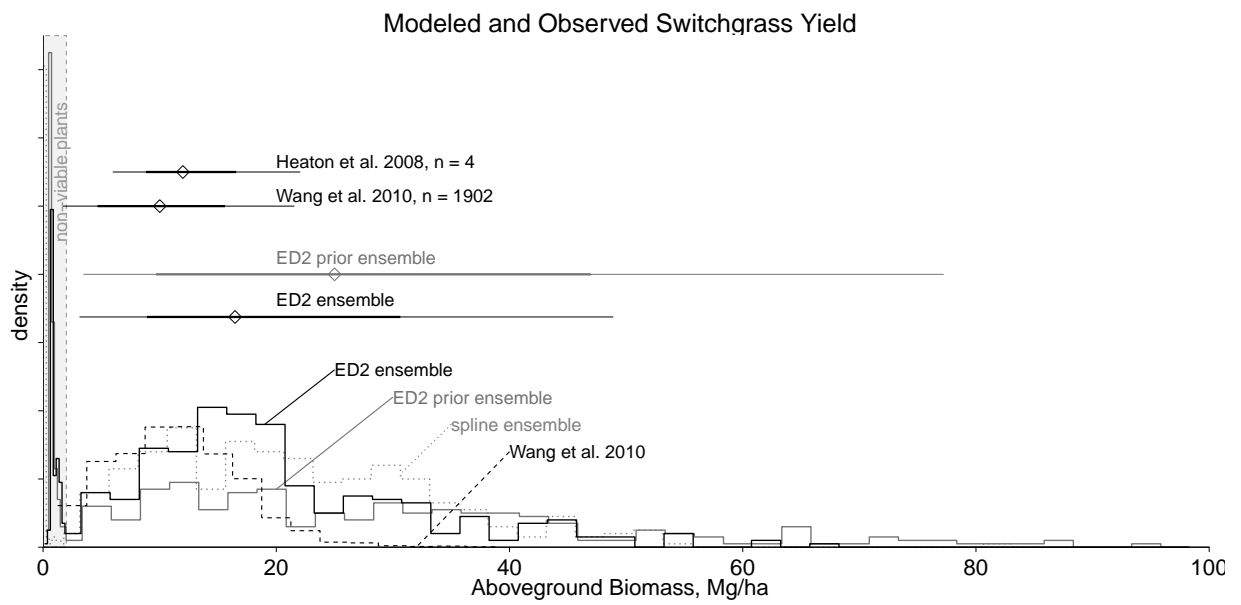


Figure 6

Sensitivity of Aboveground Biomass (Mg/ha) to Fifteen Plant Traits

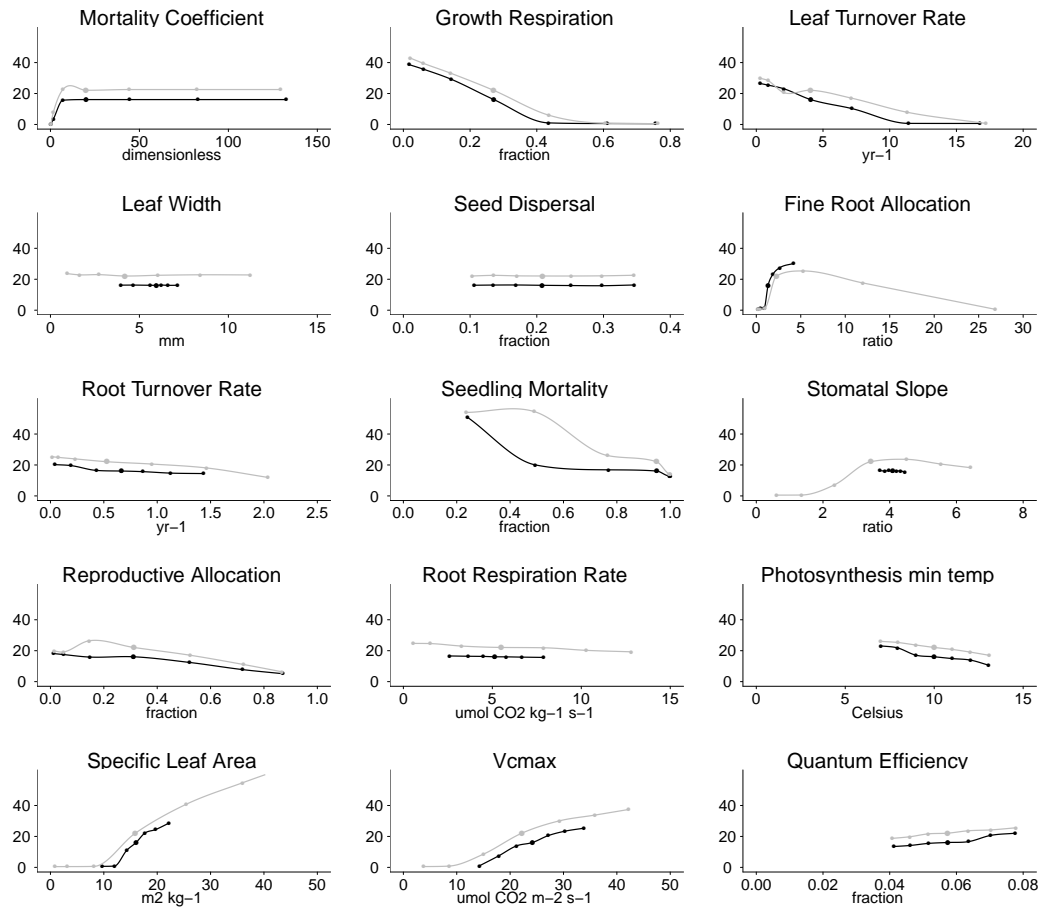


Figure 7

