

**Capturing Uncertainty in GHG Biogeochemical and Process Models:
A Path Forward for Ag Offset Protocols**

Draft Authors: William Salas (Applied Geosolutions, LLC), Steven DeGryze (Terra Global Capital), Mark Ducey (University of New Hampshire) and Johan Six (University of California at Davis). Please address questions and comments to wsalas@agsemail.com.

Goal of the white paper: Initiate discussion on how to assess uncertainty in biogeochemical process models for quantifying GHG offset projects. The objectives of the white paper are to: (1) initiate discussions on how to assess uncertainty in applications of biogeochemical process models for agricultural GHG offset projects; (2) identify sources of model uncertainty; (3) present some statistical approaches for model evaluation, since there is no single approach or right set of criteria; and (4) create a living document on this issue that improves over time.

A preliminary draft of the white paper was presented at the July 2011 C-AGG meeting in Chicago. Key points from the questions and discussion at that meeting included:

- The white paper touches on both technical questions and management and policy questions. C-AGG is a good group to tackle the latter. When beginning to evaluate and weigh in on these issues, it is important to determine whether these questions are important from a decision support perspective. Also, with regards to determining certainty/uncertainty, there is no need to determine a threshold level of either; what is important is that the level of uncertainty is or will be known. The level of uncertainty and confidence that is acceptable varies among the currently available carbon standards and protocols. It can be expected that, once the carbon market matures, a common acceptable level of uncertainty will emerge, as well as mechanisms to compensate for this uncertainty.
- Field data for independent model validation still poses substantial challenges. There are still many gaps, though different participants identified different priority gaps to fill. While there is little data on specialty crops, research on row crops, as a result of their prevalence, could be a higher impact use of limited resources. In either case, when setting data collection priorities, it is important to examine data already published in the peer reviewed literature.
- Variability in field data collection and quality is a continuing challenge for independent validation of process models. Even though there are some statistical work-arounds, comparing historical, current, and ongoing field data sets is a problem. Fortunately, procedures for collecting field data are becoming more standardized as time goes on and models are improving as they are calibrated with a growing set of high quality field data.
- C-AGG could keep a list of other issues and questions, beyond the scope of this white paper, for future development.

- Identifying and forming consensus around key questions could drive research in the future.

Background:

Agriculture represents an important near-term option for GHG offsets. Currently, the most widely accepted low-cost approaches to quantify N₂O and CH₄ emissions are based on emission factors. Emission factors relate simple input data such as the total amount of fertilizer N applied directly to N₂O or CH₄ emissions while largely disregarding factors related to weather or soil characteristics. However, given that N₂O and CH₄ emissions from agricultural practices exhibit high spatial and temporal variability, emission factors are not very sensitive to estimate this variability in emissions. It is clear that if agricultural offset projects are going to include N₂O and CH₄ reductions, then process-based biogeochemical models (PBM) are potentially important tools to quantify emission reductions within offset protocols. Success of offset projects hinges on the development of transparent tools for quantifying the full GHG impact (soil carbon, N₂O and CH₄ reductions) of changes in management practices as well as the uncertainty of the estimated GHG impact. The PBMs must be accompanied with unambiguous procedures detailing how to parameterize and calibrate them. While there have been many PBM validation efforts in the past, they have been done in an *ad hoc* fashion and with a goal of validating the model performance for a single site. Each of these PBM validation efforts employ a unique approach to parameterization and calibration that is challenging to apply beyond the conditions and circumstances of the field the PBM was validated for. In addition, most validation efforts have focused on evaluating the model's performance to predict absolute emissions accurately. The question remains how good a model's performance is with respect to emission reductions. If DNDC, or other PBMs, are to be integrated into protocols for agricultural GHG offsets, then a comprehensive and systematic validation exercise is needed to statistically quantify uncertainties in model-based estimates of GHG emission reductions that are obtained by standardized approach to parameterization and calibration that can be applied across a whole region. In addition to providing critical documentation of model performance and statistically valid uncertainty metrics, this validation work can help guide development of offset protocols and the science for improving the tools by highlighting conditions where the models does not perform well.

The total uncertainty of PBMs is usually split into two sources of uncertainty: (1) structural uncertainty and (2) uncertainty of PBM input data. The structural uncertainty is related to the inherent uncertainty of PBMs that remains even if all input data were error-free; the uncertainty of PBM input data is related to the impact of errors in the input data on simulated results. The distinction is an important one since the structural uncertainty is inherent to the model and cannot be reduced unless the model is improved, while the uncertainty in input data can be controlled by users of a PBM by e.g., expanding the number of samples on which input data is based.

This paper aims at investigating four key questions related to the use of PBMs for GHG emission reductions.

1. How do we demonstrate that a PBM is sufficiently calibrated and validated for a given set of conditions such as cropping system, region, etc.?
2. How do we develop a quantitative estimate of PBM structural uncertainty?
3. How do we capture the impact on uncertainties of PBM input data on modeled emission reductions?
4. How do we account for uncertainties in agricultural GHG offset projects?

In addition, three concrete examples of attempts to quantify different elements of the uncertainty of PBMs are provided.

Question 1: Process Model Calibration and Validation: What is sufficient?

While empirical models can be easy to use, they are limited in agricultural applications to those conditions under which they were developed. Given the variability in environmental drivers (soils biophysical conditions, climate, soil redox potential, etc), crop types (including broad range of specialty crops), cropping systems (tillage, irrigation, fertilizer use, etc) and animal production systems, it is difficult to develop empirical models that can capture the diverse conditions of agricultural systems. Modeling tools that simulate the biological, chemical and physical process based on these environmental and management drivers have a distinct advantage in that they are designed to be applicable across the broad conditions of agroecosystems. The mechanistic, or process-based, models require parameterization and calibration to simulate agricultural systems. Parameterization of a PBM is the step of selecting variables (parameters) that the model will use for simulation. Agricultural PBM typically include parameters for soil conditions (organic matter, texture, pH, porosity, wilting point, bulk density, etc), weather (temperature, precipitation, wind speed, solar radiation, etc), crop characteristics (crop type, maximum photosynthesis, crop growth rates, water requirements, nitrogen requirements, maximum yields, etc), agricultural management practices (planting and harvest dates, tillage, fertilizer use, irrigation, etc), and internal reaction rates (e.g. mass transfer coefficients, potential growth rates of microbial biomass, etc). PBM internal parameters (e.g. mass transfer coefficients) are those parameters that are set during the development of the model and cannot be modified by the user. The external parameters are those parameters that are input parameters (e.g. crop yield potential) that a user must provide to run the model. Calibration of a PBM is the process of tuning the coefficients of parameters to observations. For example, setting the maximum yield or C/N values of roots, leaves and stems of a particular crop is calibration. The calibration process can be applied to both internal and external parameters. However, calibration of the internal parameters is done only in model development by the developer. Validation of PBM requires independent measurements (measurements that were not used in calibration of internal parameters) for comparison with model estimates. One can validate internal parameters or model outputs.

In the context of using PBM for agricultural offset quantification, how do we demonstrate that models have been sufficiently calibrated and validated? Answering this question first requires an understanding of the sensitivity of model outputs to variability of input parameters. This is

typically done through sensitivity tests. The sensitivity tests provide a measure of impact of variability of individual or combinations of input parameters on model outputs (e.g. how does crop C/N ratio impact modeled GHG emissions). Once the sensitivity to a specific parameter is known, then one can decide if a model is sufficiently calibrated for that parameter by comparing model results with measurements from a controlled experiment in which this parameter was varied. In addition, a sensitivity analysis will elucidate the parameters that are most important to reduce uncertainty. Determination if a model has been sufficiently validated depends on its desired application. Ideally one would want to validate the model across the expected range of conditions for which the users wants to apply the model. The range of conditions may include various crop types, soil condition, climate conditions, management regimes, etc. Thus, this can require extensive field data that are often not available. The lack of field data is particularly challenging in case second- or higher-order interactions among input variables are calibrated. However, an advantage of models that are more process oriented over empirical models is that one does not need re-validate the model once it has been validated. It remains a challenging exercise to define the conditions and circumstances within which a model can be considered well validated.

Most validation of PBM focuses on comparing measurements and model estimates of absolute emissions rather than changes in emissions across management regimes. However, the application of PBM for agricultural offset projects focuses on estimating changes in emissions. So, if possible, model validation should assess model uncertainty for both absolute emissions and changes in emissions across the management systems under consideration for offset projects. Again validation of changes in emissions is limited by availability of independent field data.

Questions to address: When is re-calibration needed? Differences in crop varieties, differences in cultural practices, soil biophysical and climate drivers... What statistical approaches should be used to assess the need for re-calibration? (See mixed effects discussion). How can one delineate the conditions within a model can be considered to be well validated?

Question 2: Quantifying model structural uncertainty

Model structural uncertainty is defined as model accuracy when all inputs are known and free of error. Quantifying model structural uncertainty requires the use of independent validation data., The basic approach entails statistical analysis of the residuals of modeled versus measured estimates of greenhouse gas emission. As part of the statistical analysis one can examine heteroskedasticity to asses if model uncertainty increases with higher emissions. Approaches for quantifying uncertainty are separated into parametric approaches (where one must assume distribution of the residuals) and non-parametric (no distribution assumed) approaches. There are several key questions to address in quantifying model structural uncertainty, including the following:

- ❖ What is the appropriate uncertainty metric? RMSE or % deviation. Dividing the error in different components: lack of correlation, non-unity slope, and bias
- ❖ How to define confidence intervals? As a factor of the average or absolute amount (related to assumption on distribution)
- ❖ Minimum number of validation points for deriving structural uncertainty metrics?

- ❖ Using local/regional/global datasets: (relates to the question regarding recalibration)

One approach for quantifying model structural uncertainty is provided in Example 1 at the end of this document.

Question 3: Assessing uncertainty in Input Data

As discussed above, PBM require input data on soils, climate, crops and agricultural management. However, for many applications collecting these data can be expensive. If collecting these data are too expensive for agricultural offset project developers, then is it possible to account for uncertainties in model inputs on modeled offsets? There are several key questions to be addressed, including:

- ❖ What are major sources of uncertainty in input data for process models: soils (using soil surveys versus measured data), crop parameters (relates to the previous section discussion regarding crop varieties), spatial variability in precipitation (climate stations location relative to project fields).
- ❖ Often, a Monte Carlo approach is used to simulate the impact of variability in a PBM's input data on the simulated results. In a Monte Carlo approach, a large number of model runs are executed using input parameters that are slightly different in each run. While a Monte Carlo approach is an elegant solution to analyze uncertainty in emissions from variability in input data, several challenges remain.
 - How can one estimate the parameter probability density functions, i.e., the range in which parameters are expected to vary?
 - How can one account for correlation in input parameters, in which certain combinations of input parameters are much more likely than others (e.g. high soil SOM contents is much more likely on soils with high clay content, compared to soils with small clay content)?
- ❖ Discuss tradeoffs between measurements and other input sources Examine approaches for integrating a sensitivity analysis of parameters with costs to measure parameters and thus create a graph of uncertainty vs. cost

One approach for assessing impact of uncertainty of input data on modeled agricultural offsets is described in Example 2 at the end of the document.

Question 4: Accounting for uncertainties in Agricultural GHG Offset projects

How can we discount modeled reductions based on uncertainty statistics derived from analyses of model structural uncertainty and the impact of input uncertainties on modeled reduction? Carbon projects can have drastically different circumstances that affect how accurate calculations of emission reductions are. In addition, even with great care and effort, there remains a small probability that the calculated emission reductions are not conservative, i.e., smaller than the true emission reductions. To ensure consistency across projects and across project types, one must standardize calculations according to the certainty of the calculations for a specific project. The most logical way of doing this is by using a confidence interval and calculating the minimum amount of credits that is generated with a certain confidence. In this

context, the value of the confidence can be interpreted as the percentage of projects with similar uncertainties as the one in question that do not over-estimate credits when only this minimum amount of credits is claimed. The exact value of the confidence limit is a political decision that lays with the carbon standard.

By using a deduction based on uncertainty, one does not only achieve consistency across projects, but also incentivizes better monitoring and provides flexibility for projects with a small field monitoring budget.

Agricultural GHG projects calculate GHGs as the difference between project emissions and baseline emissions. The uncertainty around the difference of two emissions will be different than the uncertainty around absolute emissions. More specifically, since the baseline scenario is so similar to the project scenario (except for the variables that are affected by project activities), it is to be expected that there will be a strong correlation between baseline emissions and project emissions. The extent of this correlation affects how the uncertainty around absolute emissions is related to the uncertainty of emission reductions. A Monte Carlo analysis can be used to calculate the variation in emission reductions as input parameters are changed.

There are two key issues to be addressed in accounting for uncertainties in model quantification of offsets:

1. Impact of input variables on structural uncertainty

Accounting for interactions between model structural uncertainties and uncertainty from input parameters. As was discussed before, there are two components to the uncertainty of model simulations. Under the simplest (and most conservative) assumptions, these components can be combined as if there was no correlation between the two uncertainty terms. However, it is known that a model's performance is strongly dependent on the values of the input parameters. For example, it is known that the accuracy of modeled nitrous oxide emissions decreases with increasing soil organic matter content. The stronger the correlation between these two error terms, the smaller the combined uncertainty. If enough samples in a dataset are available, one can theoretically calculate the structural uncertainty given a set of input parameters.

2. The impact of aggregation in reducing model structural uncertainty

Carbon credits are always calculated and managed for an individual project that can consist of different individual fields. In fact, due to the complexity of GHG offset protocols and the transaction costs associated with offset projects, it is likely that several agricultural fields will be combined within a discrete agricultural GHG offset project. The aggregation of individual fields into one project will affect the calculation of uncertainty and the uncertainty deduction. The uncertainty will likely decrease as more fields are included within a project since the uncertainty can be averaged over individual fields. The uncertainty deduction mechanism incentivizes project aggregation since the (relative) deduction will be smaller if more fields are combined within a carbon project package. This is discussed in Example 1 below.

In Example 1 we discussed an approach to reduce model structural uncertainties through aggregation of multiple fields. One approach for accounting for uncertainties in agricultural offset projects is to apply a deduction for model structural uncertainty (μ_{struct}) on the project aggregate and the deduction for uncertainty due to input uncertainties ($\mu_{inputs,i}$) to each field in the aggregate. Using this approach one can quantify total primary effect greenhouse gas emission reductions (tCO₂eq) for a project area as follows:

$$M - Offsets = \mu_{struct} * \sum_{i=1}^m \mu_{inputs,i} * (GHG_{BSL,i} - GHG_{P,i})$$

Where,

$M - Offsets$	=	Modeled GHG emissions reductions over the project area
m	=	Number of individual fields included in the project area
$\mu_{inputs,i}$	=	Accuracy deduction factor for individual field i due to input uncertainties (% reduction for each field)
$GHG_{P,i}$	=	Project emissions in year y for individual field i
$GHG_{BSL,i}$	=	Baseline emissions in year y for individual field i
μ_{struct}	=	Accuracy deduction from model structural uncertainty (% reduction)

Example 1: Quantifying Model Structural Uncertainty

The following illustrates an approach for quantifying process model (DNDC model in this case) structural uncertainty for quantifying methane emissions for rice systems in California.

Basic Assumptions

- The structural error induced by a biogeochemical model such as DNDC is multiplicative, not additional:

$$Y_{field,i} = Y_{model,i} \cdot \varepsilon_i$$

The multiplicative nature of the deviation between modeled and measured results originates from increasing deviations with increasing modeled values.

- No bias exists between measured and modeled results, so that $\langle Y_{field} \rangle = \langle Y_{model} \rangle$. In addition, the error ε is log-normally distributed, $\varepsilon \sim \ln \mathcal{N}(0, \sigma)$.
- Model results of an alternative treatment are 100% correlated with the model results of the baseline treatment:

$$Y_{field,project} = k \cdot Y_{field,baseline}$$

Where k is dependent on all factors that were not impacted by the project. In other words, changes in emissions due to weather or other non-critical variables are similar between project and baseline scenarios, apart from a linear constant.

Theoretical Foundation

Since the structural error is multiplicative, the residual of the log-transformed field and measured results is normally distributed:

$$\ln(Y_{field} - Y_{model}) \sim \mathcal{N}(0, \sigma)$$

Assume that n is the number of $(Y_{field,i}, Y_{model,i})$ pairs, σ can be estimated as:

$$s = stdev(\ln(Y_{field,i} - Y_{model,i}))$$

Since σ is not known, traditional statistical theory dictates that confidence and prediction intervals need to be estimated based on the student t-distribution with n degrees of freedom. We are interested in the effect of taking averages of individual fields on the decrease in the uncertainty. However, since the sum of different student t-distribution does not have an easy analytical form, we will assume that the error σ is normally distributed. In this case, the 95%-confidence prediction interval becomes:

$$[-s \cdot \phi(0.025); +s \cdot \phi(0.975)]$$

In case one is looking at the average of m field measurements, the 95% -confidence prediction interval around the m measurements becomes:

$$\left[\frac{-s}{\sqrt{m}} \cdot \phi(0.025); \frac{+s}{\sqrt{m}} \cdot \phi(0.975) \right]$$

The discounting factor u_{struct} must be set so that, with 95% confidence:

$$u_{struct} \cdot (Y_{model,alternative} - Y_{model,baseline}) < Y_{field,alternative} - Y_{field,baseline}$$

Using assumption 2, this comparison can be simplified as following

$$u_{struct} \cdot Y_{model,baseline} \cdot (1 - k) < Y_{field,baseline} \cdot (1 - k)$$

$$u_{struct} \cdot Y_{model,baseline} < Y_{field,baseline}$$

After taking a logarithm and rearranging:

$$\ln(u_{struct}) < Y_{field,baseline} - Y_{model,baseline}$$

The discounting factor for structural uncertainty is therefore:

$$u_{struct} = e^{\frac{-s}{\sqrt{m}} \phi(0.025)}$$

Application to DNDC model runs in California

Table 1 Data from Fitzgerald et al. 2000¹ and Horwath/Assov² unpublished, preliminary data

Seeding	Tillage	Winter Flooding	Residue	Modeled kg CH ₄ - C ha ⁻¹	Measured
Water ¹	Conventional	Yes	incorporation	121	130
Water ¹	Conventional	Yes	burn	56	52
Water ¹	Conventional	No	incorporation	68	75
Water ¹	Conventional	Yes	incorporation	166	273
Water ¹	Conventional	Yes	burn	56	57
Water ¹	Conventional	Yes	incorporation	71	165
Water ²	Conventional	N/A (GHG measurements)	N/A, no GHG	465	354

C-AGG Whitepaper: Uncertainty in Process Models and Ag Offset Protocols

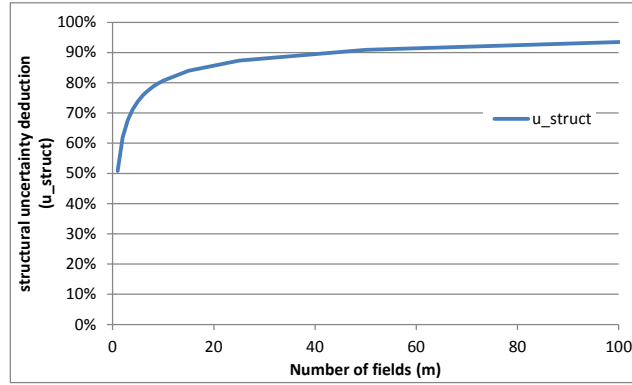
		were only for in-season emissions)	measurements after residue management		
Water ²	Stale seedbed (essentially no-till prior to plant)	N/A (GHG measurements were only for in-season emissions)	N/A, no GHG measurements after residue management	417	390
Dry ²	Conventional.	N/A (GHG measurements were only for in-season emissions)	N/A, no GHG measurements after residue management	254	229

The average of the ln of the deviations is 0.112; the standard deviation is 0.346. Using the equations above, the appropriate discounting factors are:

$$u_{struct}(m) = e^{\frac{-0.346}{\sqrt{m}} - 0.196}$$

Number of fields (<i>m</i>)	<i>u_{struct}</i>
1	51%
2	62%
3	68%
4	71%
5	74%
6	76%
7	77%
8	79%
9	80%
10	81%
15	84%
25	87%
50	91%
100	93%
1000	98%

C-AGG Whitepaper: Uncertainty in Process Models and Ag Offset Protocols



A minimum of 5 required fields would correspond to a maximal uncertainty deduction of 26% (100% - 76%).

Example 2: Quantifying impact of input uncertainty: NRCS Soils

One approach for assessing the impact of model input uncertainty on modeled GHG emissions and assessment of GHG reductions from the implementation of an offset project is to use Monte Carlo methods. To apply this method for assessing the impact of uncertainty of soil conditions, the first step entails defining a possible range and probability distribution of the soil conditions. To assess a range of results for greenhouse gas emissions at a single rice field, this example demonstrates a Monte Carlo modeling approach using DNDC. Our general approach was to assume some variability in site soil attributes (clay fraction, organic matter fraction, bulk density, and pH) as modeled in the USDA NRCS SSURGO soil model. Using a Monte Carlo simulation, we modeled identical crop management practices and meteorological conditions while varying soil conditions through 1,000 iterations. Rather than using the uncertainty tools included with DNDC (which assume an even distribution of the tested parameters), our approach assumed a log-normal distribution of each of the soil attributes as well as some amount of correlation between them. Our approach involved three aspects:

1. an analysis of correlation between the four soil attributes
2. programmatic generation of DNDC inputs based on the Monte Carlo method
3. running the DNDC model in site mode and synthesizing the results

We applied this approach in two ways, the first assumes no correlation between soil parameters, which is conservative since we know that there is significant correlation. The second set of Monte Carlo runs utilized correlation statistics as part of the sampling procedure.

Soil attributes are stored within the SSURGO database according to the following relationships:

Horizon	Contains soil attribute data (low, representative, and high values) based on an assessment of soil field conditions
↓	[one to many]
Component	The basic soil type (roughly equivalent to soil series) – soil components have many horizons and have no explicit spatial location
↓	[one to many]
Map Unit	The smallest mapped polygon in the SSURGO model – soil map units have many components of varying fractions

To assess correlation among soils in rice growing areas of California, we selected all map units intersecting rice fields as mapped in the CA Department of Water Resources land use database. From this selection, we identified all soil components contained within the map units. Soil attribute data came from the top horizon for each component. Thus, our final database represents all soil horizons intersecting rice fields.

We calculated Pearson correlation coefficients for each set of pairs for representative values of the four soil attributes:

	Clay fraction	OM fraction	Bulk Density	pH
Clay fraction	1	-	-	-
OM fraction	0.139	1	-	-
Bulk Density	-0.526	-0.685	1	-
pH	0.263	0.098	-0.126	1

Generation of DNDC inputs

Our Monte Carlo simulation randomly generated 1,000 numbers for each of the four soil properties with the correlation matrix and with each following a log-normal distribution. This was done by using the Cholesky decomposition of the correlation matrix to transform a set of standard-normal random numbers in the logarithm space. The representative value was used as the mean, while the low and high values were transformed into log space and treated as a range of +/-3 standard deviations. This resulted in four sets of 1,000 correlated random numbers, normally distributed. The soil properties, other than pH, were finally calculated by taking the exponent of the numbers.

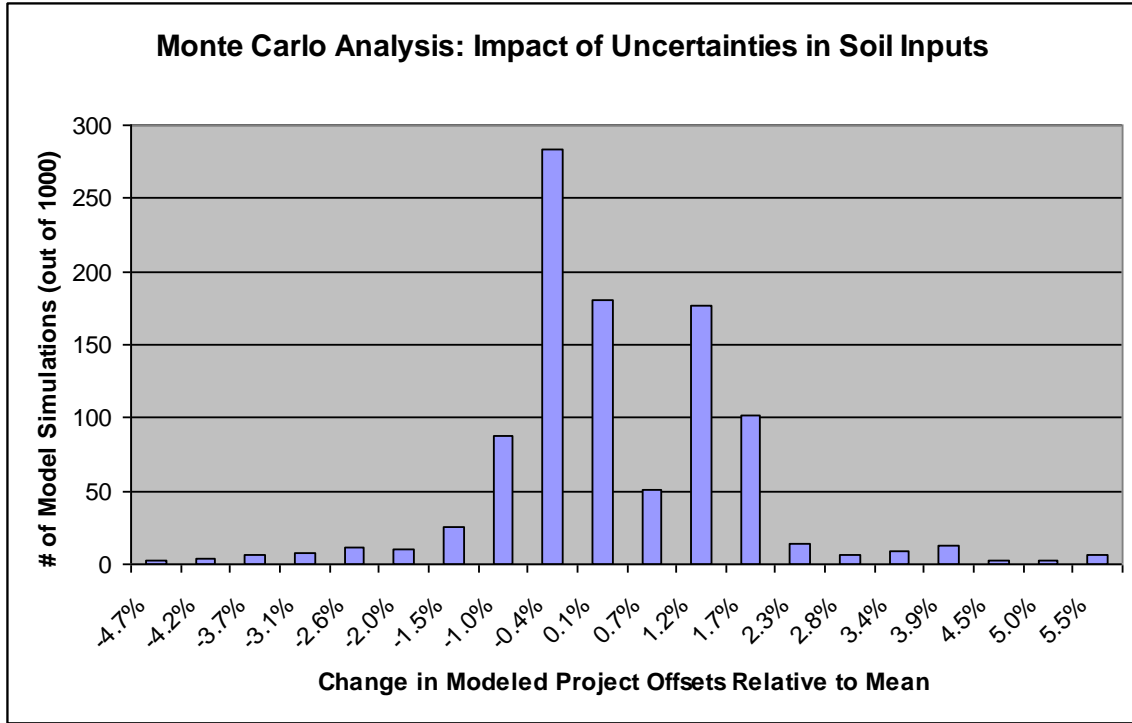
DNDC modeling

The model was run as a batch using the DNDC site mode. We ran two scenarios (one with a winter flood, one without a winter flood) for a single field as follows:

- ♦ Rice, planted 5/1, harvested 9/11
- ♦ Tillage on 4/23, 4/26, 4/27, 4/29, and 9/15
- ♦ Fertilizer on 4/30 (injected anhydrous ammonia), 5/1 (surface application of $(\text{NH}_4)_2\text{HPO}_4$), 5/26 (surface application of $(\text{NH}_4)_2\text{SO}_4$)
- ♦ Flooded from 5/1 to 9/1
- ♦ Winter flood from 11/15 to 1/31 (only for the winter flood scenario)
- ♦ Rice straw burned once every 8 years

We also ran a 16-year initialization period (using the baseline crop management) to stabilize soil organic matter. Only the last eight years (one year with residue removal, seven years with residue incorporation) were averaged in the final summary.

Our results indicate the modeled methane emissions and net GHG emissions are quite sensitive to soil conditions. At 90% confidence interval, the modeled CH₄ and net GHG emissions were significant (over 14% in both baseline and project simulations). However, the impact of soil uncertainties on modeled changes in emissions from baseline to project conditions were quite small (<3%). The figure below shows the histogram of the Monte Carlo simulation results for the case assuming no correlation between soil input parameters. It is clear for this baseline and project scenario, that uncertainty in soil input parameters impacted both baseline and project modeled emissions in a similar degree. Accounting for correlation between soil input parameters reduced uncertainties. The table below summarizes these results.



	Assuming no correlation in soil input parameters		Accounting for correlation of input soil parameters	
	CH4 GWP (90% CI / Mean)	Global Warming Potential (90% CI / Mean)	CH4 GWP (90% CI / Mean)	Global Warming Potential (90% CI / Mean)
Baseline	14.7%	14.4%	14.0%	13.7%
Project	18.5%	20.0%	17.5%	19.1%
Baseline-Project	1.0%	2.2%	0.2%	1.4%

Appendix: Statistical evaluation of model performance: A Philosophy

What we assume here is that there is a model, with a documented procedure for parameterization and calibration. This procedure has been followed and now we have an independent test data set (i.e. the model has *not* been calibrated to the particular outputs we are now testing). We would like to test whether the model is performing adequately. It is very easy to get hung up here on the statistical procedures and ignore the all-important steps of asking the right questions, and defining what “perform adequately” means. Let’s focus on the definition part first, because that determines the statistical measures and approaches we should use (and not vice versa!).

There are three basic levels of performance we may wish to consider. Not all are equally important for every context. As a result, there is no single “right” criterion for deciding whether a model is performing well, there is no single “right” statistical test, and there is no single “right” threshold. We will focus this discussion, on the right criterion and threshold within the context of GHG offset projects.

1. Closeness of individual predictions to individual measurements. On a case-by-case basis, do model predictions match up well to observations? This is the most relevant performance test within the context of GHG emission protocols. Of course, unless we have a perfect model, the match will be less than perfect for many or all observations, and there will likely be some for which large errors occur no matter how good the model is (because of input data issues, or because processes that rarely matter and are not represented well in the model do occasionally influence the outcome in the field). What constitutes “closeness of predictions”? Is it a relative matter (e.g. within 10%) or an absolute one (e.g. within 10 kg C/ha/yr)? For what fraction of observations would it be acceptable to have out-of-tolerance results? 10% of observations, or 1% of observations? (Demanding that all predicted and observed values have very close tolerances is likely to lead to frustration with real models and real data for complex phenomena.) If one observation corresponds with one carbon project, an alternative way of posing the same question is which percentage of all projects within a carbon standard’s project portfolio will have overestimated its carbon offsets? Most likely, this percentage will not be tolerable, and the percentage of projects that do not overestimate credits will have to be increased by applying an uncertainty deduction. The question then becomes, if a deduction is applied, which percentage of carbon projects does the standard allow to be overestimating its credits.

2. Correlation of modeled and measured values. Over a large number of specific cases, is it generally true that high predictions match high observations, and low predictions match low observations? If not, then the model will do a poor job of representing within-region variability and maps developed using the model will not provide good spatial representations. Ideally, we would like the modeled and measured values to fall close to a 1:1 line with an intercept near zero. Of course, unless we have the perfect model the slope of the “best fit” line between predictions and measurements will not be exactly 1. But how much error can we tolerate? How should utilize information on model fit when the intercept is not near zero? Defining reasonable tolerance here is not a statistical

question, but a management and policy question that depends on the eventual use of the model.

3. Average model performance. Over a large number of specific cases, does the average prediction given by the model match up well with the average value that is observed? If not, then in the most basic statistical sense we would call the model biased. If a model is to be used within the context of GHG protocols, a model will have to be unbiased, or, at least, not biased so that credits are systematically overestimated. Of course all models will be biased to some degree (unless we are modeling a trivially simple phenomenon associated with well-defined physical constants). How much bias can we tolerate?

4. Extending the use of models to new regions. Often we wish to implement a model in a new region. We may have a “global” validation or test data set, or at least a test set that comes from several other regions. Some reasonable questions to ask include: do we need to perform a new validation? If so, how many data points would we need, either to provide an evaluation or better yet to recalibrate our expectations about patterns of bias or variability?

Statistical Approaches

Above, we said there is no single “right” criterion for deciding whether a model is performing well, there is no single “right” statistical test, and there is no single “right” threshold. However, there are some approaches that are sensible and others that are less so.

Traditionally, statistical assessment of model performance has applied standard Fisherian hypothesis testing approaches, with varying levels of sophistication. Although new approaches may be better suited to the task, the traditional approaches are still widely used and are not completely without value, so let’s review them. We’ll assume we have pairs of predicted and measured values to work with, and that these pairs are in some way a representative sample of some domain of interest (such as a region or project).

Over a large number of specific cases, does the average prediction given by the model match up well with the average value that is observed? This is really a question about two means, the mean of the observed values and the mean of the predicted values. Since we have pairs of values, we can calculate differences

$$d_i = obs_i - pred_i$$

and then perform an appropriate test to see if the mean of the differences is significantly different from zero. For example, if we are willing to assume that the differences are approximately normally distributed, we could perform a simple t test, rejecting the null hypothesis that the average difference is zero if $p \leq 0.05$. If we do not believe the differences are normally distributed, then a nonparametric (bootstrap) test could be used instead.

Over a large number of specific cases, is it generally true that high predictions match high observations, and low predictions match low observations? Most conventional approaches to this question involve regression. At the most elementary level, one could fit a regression and examine R^2 to see if it is close to 1. However, it is difficult to define what an acceptable value of R^2 should be (what proportion of variance explained is good enough? and how does that differ between populations that are highly variable and those where the variability is relatively small?) Moreover, one can have a high value of R^2 even if the slope is very far from 1 (or even negative, though we might hope that model predictions would at least be positively correlated with observed values!) A more sophisticated approach would be to use the estimated slope and its standard error to construct a t test of the null hypothesis that the slope equals 1 (which is not the same as the null in the usual test reported by statistical software, that the slope equals 0). Where the residuals of the regression are not normally distributed, a bootstrap t test would be a logical alternative. Where outliers are present, robust regression procedures could be used. And so on. Of course one must also answer the question of which variable – the observation or the prediction – should go on the y axis and which will go on the x axis, as this will influence the estimate of the slope and the outcome of the test. Although it is the model prediction that “has the error,” and on that basis we might think predicted values should be y and observed x , the eventual use of the results suggests the opposite: down the line we will have a large number of predictions and we are trying to draw inferences about unknown true values that have not been observed. So the case is stronger for putting predicted values on the x axis and observed values on the y axis.

On a case-by-case basis, do model predictions match up well to observations? This question is both the simplest and the hardest to deal with. If the data set being used for validation is representative of the population we are eventually interested in, then for a great many criteria the estimation process is straightforward. For example, if we wish 90% of the predictions to be within 10% of the actual observations, and we have a representative sample, we can simply calculate the proportion of the data that meet the criterion, and even use familiar confidence limit calculations to see if our particular results might have met (or failed to meet) the specification by virtue of good (or bad) luck. Likewise, quantities such as variances can be estimated unbiasedly from the sample (and, for a sample of reasonable size, the estimates of quantities like relative root mean squared error will be nearly unbiased). The difficulty, of course, is in getting a “representative sample,” especially when quality studies that can be used for validation are sparse or come from a broader, narrower, or just plain different situation than the one we wish to know about. The results of such tests and calculations should typically be taken with some humility.

Equivalence Testing. There is, unfortunately, a fundamental challenge to the traditional approaches, especially to the first and second questions. This challenge can be illustrated by a very simple example. Suppose we have only a handful of data for testing a model in a particular context, and the data are somewhat noisy. We perform a t -test of the difference between the predicted and observed values, or fit a regression of observed vs. predicted values, and test whether the slope of the regression equals 1. Of course with

very little data, our confidence limits on either the mean or the slope will be very wide. So even if our estimates (or the true values) are far from 0 for the mean, or far from 1 for the slope, we will be unable to reject the null hypothesis. We have too little power. Therefore, we are likely to “accept” a model not because its performance is good but because our evidence is weak. Conversely, if we have a great deal of data we may be able to reject the null that the difference is 0, even if it is so small as to be of no practical consequence, and we may do the same for the null that the slope is 1 even if it is very close. In either case, from a statistical perspective we can never prove that the model is right. The take-home message is simple: collect as little data as possible so that your model will not be rejected!

The reason for this perverse situation is that in conventional statistical practice, the roles of the null and alternative hypotheses are essentially backwards. What we want is to demand that the evidence substantiate that the model is good enough. (In conventional practice, we are asking the evidence to show that the model isn’t perfect – something we already knew, but we treat perfection as the null hypothesis that must be disproven.) A lucid but somewhat theoretical review can be found in the article by Robinson et al. (2005). One solution is the set of statistical techniques known as equivalence testing (Berger and Hsu 1996, Wellek 2003). Although these approaches have been used for some time in biomedical work, they are only beginning to be explored in agricultural and forestry contexts.

Extending the use of models to new regions: These kinds of problems and questions are not unique to agricultural and environmental modeling. In the United States, for example, they arise commonly in the analysis of data from the U.S. Census Bureau. Recall that every 10 years, the U.S. Census is supposed to identify and enumerate all residents of the U.S., and this is done with a “short form.” A small sample of people receives a “long form” that includes a great deal of additional identifying information. But at many levels of aggregation, for example for an individual town, there may not be enough “long form” data to estimate certain quantities or relationships accurately. The challenge is to make optimal use of the data at a regional level to identify general patterns and trends, then use the local data to calibrate or modify the regional expectation. In the statistical literature these techniques fall under the general name of “small area estimation” techniques (Rao 2003).

There are many approaches for small-area estimation problems, but a simple one can be illustrated through a mixed-effect model (Pinheiro and Bates 2000). Let’s write out the usual regression equation for observed vs. predicted values:

$$Obs_i = a + b Pred_i + e_i$$

Here a is the intercept (which we hope is very close to 0), b is the slope (which we hope is very close to 1), and e_i is the error term for an individual observation (which we most often model using the normal distribution, and which we hope has very small variance).

But suppose the relationship between observed and predicted might change from region to region. To account for this we should add one or two terms to our equation:

$$Obs_i = (a + c_{region}) + (b + d_{region})Pred_i + e_i$$

Here c_{region} and d_{region} represent the departure of the region where observation i falls, from the “average” values of the slope and intercept. Under some mild assumptions, we can incorporate c_{region} and d_{region} into the model – treating them as random variables akin to e_i – and develop not only predictions for the regions we have in our data, but also some understanding of how c_{region} and d_{region} vary among possible new regions (for example, what their standard deviation might be, indicating how the slope and intercept vary between regions). Armed with that knowledge, we could in principle predict how much data we would need to recalibrate our expectations for a model in a new region.

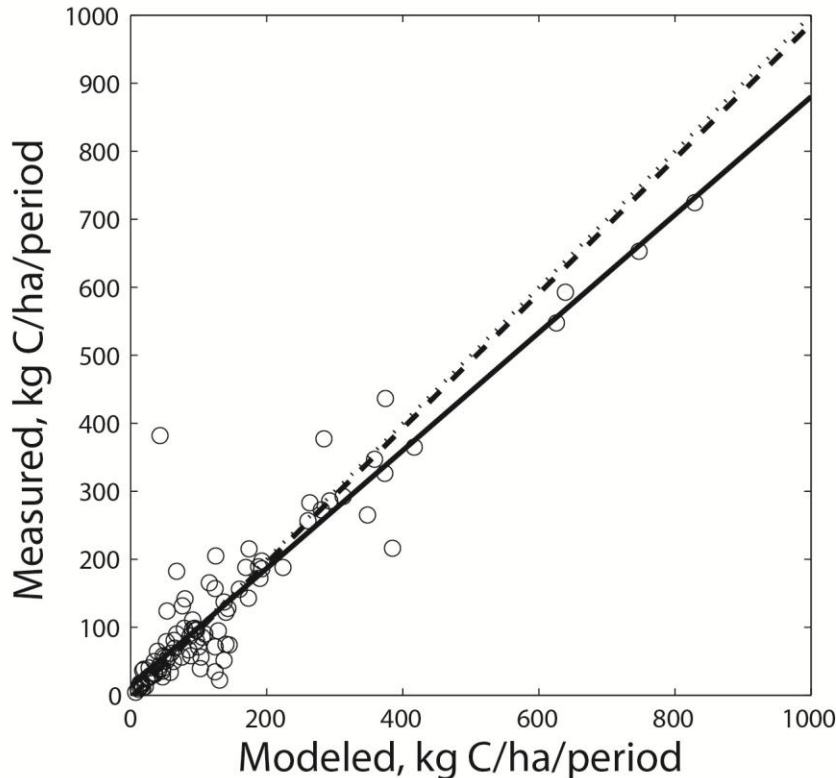
Applying the DNDC Model for Rice Methane Offset projects: An Example

To illustrate these concepts, we compiled a dataset of 99 independent validation data for rice methane. These data were compiled from 9 peer reviewed papers and a recent validation exercise in California using data from one peer reviewed publication and a recent dataset collected at Rice Experiment Station in Biggs, CA (data courtesy of Dr. W. Horwath, UC Davis). For each validation, the DNDC model was calibrated to the local rice variety using methods similar to calibration approaches discussed in draft and existing protocols for ACR, VCS and CAR. However, it is important to point out that this dataset was compiled from several versions of the DNDC model. We are in the process of trying to compile all of the input data and field data to redo the validation of the 99 datapoints using a single version of the model. Nevertheless, we feel this dataset is useful for demonstrating the concepts outlined in this whitepaper.

On a case-by-case basis, do model predictions match up well to observations? The answer here is “it depends” – depending, of course, on what “match up well” means and how many observations are allowed to match poorly. For example, half of the individual data points have predictions within 18.0 kg C/ha/period of the observed values, and fully 75% are within 43.5 kg/ha/period. In relative terms, half of the predictions are within 19% of their corresponding observations, and 75% are within 38%. Is this good enough *at an individual site scale* – where we can probably tolerate much larger errors than we would at a regional or project-level scale?

Over a large number of specific cases, is it generally true that high predictions match high observations, and low predictions match low observations? Our answer to this question will depend on some statistical choices. If we make the observed values the x -axis and the predicted values the y -axis, then $R^2=0.88$, which we might conclude is a pretty good fit for a process model but the number itself carries little practical meaning. More importantly, the estimate of the intercept (3.41 ± 7.48) is not significantly different

from zero, and with this choice of axes we would conclude that the slope (1.01 ± 0.04) is not significantly different from 1. However, as suggested above there is a very strong case that we have the axes backward here: observed values should be on the y-axis and predicted values on the x-axis. In that case the R^2 doesn't change, but the regression coefficients do. The intercept (13.39 ± 6.80) is almost significantly different from zero ($p=0.0518$), and the slope (0.87 ± 0.03) shows a strong significant difference from 1 ($p<0.0001$). See figure below. Now we might well be concerned about the influence of possible outliers on this relationship, and could easily launch ourselves down a whole series of alternative regression analyses to deal with them, but those would change the interpretation of our results. (For example, the simplest solution – deletion of outliers – might move the slope closer to 1, but then our conclusion would be that the model corresponds well *except for some outlying cases that were deleted because they disagreed with the model* – hardly a satisfying conclusion).



Over a large number of specific cases, does the average prediction given by the model match up well with the average value that is observed? The average value of the observed flux (130.7 kg/ha/period) is very close to the average value of the predictions (135.3 kg/ha/yr). The standard error of the paired differences is only 5.5 kg/ha/period (or about 3.5%), so the typical *t*-test approach would lead us to conclude that the difference is not statistically significant ($t=0.85$, $p>0.05$). This does **not** mean there is no

difference! It only means our data are too small to resolve it, and it is probably small (within about 11 kg/ha/period, or 7% of the mean, with 95% confidence.)

Equivalence Testing. Without going into gory statistical detail, it's fair to say that the equivalence testing approach – using a nonparametric bootstrap test to avoid the impact of assuming normality – gives very similar results to the traditional approaches, in part because this is a large data set. For example, if we demand that the evidence show the modeled average is within 10% of the observed average, and use a 95% probability basis for our inference, the equivalence test rejects the null hypothesis that the model performance is worse than our specification – in other words, we can accept the average performance of the model. But if we demand that the evidence show that the slope of our observed vs. modeled relationship is between 0.9 and 1.1, we fail to reject that null – in other words, the data don't substantiate that the model meets that specification. The result is similar to that of the corresponding traditional regression test, but the interpretation is somewhat different. Again, the equivalence testing approach is sensitive to the presence of outliers, so that is an issue we would probably explore in a detailed investigation.

Mixed Effects Example: Can our global validation data set tell us anything about likelihood of DNDC applicability for estimating methane emission from rice in new regions (like Arkansas and Louisiana, for example)? Unfortunately, with our example data set, this approach is rather anticlimactic! The test data used here suggest that c_{region} and d_{region} actually have zero variability – in other words, the slope and intercept appear to be consistent across all regions. In a sense, that is good news: it is reasonable to expect that the kinds of variability around the central trend that we see in this data, would play out similarly when implementing the model in a new region. While it would certainly enhance our confidence in the model to have some region-specific test data to confirm that, it also does not appear critical for modeling CH₄ fluxes with this particular model. On the other hand if either c_{region} and d_{region} did have an appreciable standard deviation, that would imply we really would want to acquire some data in a new region just to understand how the model performance might differ with previously studied regions. The results of estimating a , b , c_{region} , and d_{region} could be used in a variety of modeling approaches to simulate the recalibration process and identify an appropriate sample size.

References

Berger, R.L. and Hsu, J.C. 1996. Bioequivalence trials, intersection tests and equivalence confidence sets. *Stat. Sci.* 11: 283–319.

Pinheiro, J.C. and Bates, D.M. 2000. *Mixed-Effects Models in S and S-Plus*. Springer, New York.

Rao, J.N.K. 2003. *Small Area Estimation*. Wiley: New York.

Robinson, A.P., Duursma, R.A. and Marshall, J.D. 2005. A regression-based equivalence test for model validation: shifting the burden of proof. *Tree Physiology* 25: 903-913.

Wellek, S. 2003. Testing statistical hypotheses of equivalence. Chapman and Hall, London, 284 p.