

## Article (refereed)

---

**Van Oijen, M.; Cameron, D.R.;** Butterbach-Bahl, K.; Farahbakhshazad, N.; Jansson, P.-E.; Kiese, R.; Rahn, K.-H.; Werner, C.; Yeluripati, J.B..  
2011 A Bayesian framework for model calibration, comparison and analysis: application to four models for the biogeochemistry of a Norway spruce forest. *Agricultural and Forest Meteorology*, 151 (12). 1609-1621.  
[10.1016/j.agrformet.2011.06.017](https://doi.org/10.1016/j.agrformet.2011.06.017)

Copyright © 2011 Elsevier Ltd.

This version available <http://nora.nerc.ac.uk/15938/>

NERC has developed NORA to enable users to access research outputs wholly or partially funded by NERC. Copyright and other rights for material on this site are retained by the authors and/or other rights owners. Users should read the terms and conditions of use of this material at <http://nora.nerc.ac.uk/policies.html#access>

**This document is the author's final manuscript version of the journal article prior to the peer review process. Some differences between this and the publisher's version may remain. You are advised to consult the publisher's version if you wish to cite from this article.**

[www.elsevier.com/](http://www.elsevier.com/)

Contact CEH NORA team at  
[noraceh@ceh.ac.uk](mailto:noraceh@ceh.ac.uk)



27 comparison, (3) Analysis of model-data mismatch.

28 Data were available for four output variables common to the models: soil water  
29 content and emissions of N<sub>2</sub>O, NO and CO<sub>2</sub>. All datasets consisted of time series of daily  
30 measurements. Monthly averages and quantiles of the annual frequency distributions of daily  
31 emission rates were calculated for comparison with equivalent model outputs. This use of the  
32 data at model-appropriate temporal scale, together with the choice of heavy-tailed likelihood  
33 functions that accounted for data uncertainty through random and systematic errors, helped  
34 prevent asymptotic collapse of the parameter distributions in the calibration.

35 Model behaviour and how it was affected by calibration was analysed by quantifying  
36 the normalised RMSE and  $r^2$  for the different output variables, and by decomposition of the  
37 MSE into contributions from bias, phase shift and variance error. The simplest model,  
38 BASFOR, seemed to underestimate the temporal variance of nitrogenous emissions even  
39 after calibration. The model of intermediate complexity, DAYCENT, simulated the time  
40 series well but with large phase shift. COUP and MoBiLE-DNDC were able to remove most  
41 bias through calibration.

42 The Bayesian framework was shown to be effective in improving the parameterisation  
43 of the models, quantifying the uncertainties in parameters and outputs, and evaluating the  
44 different models. The analysis showed that there remain patterns in the data - in particular  
45 infrequent events of very high nitrogenous emission rate – that are unexplained by any of the  
46 selected forest models and that this is unlikely to be due to incorrect model parameterisation.

47

48

49 **Keywords: Bayesian calibration; carbon cycle; likelihood of data; nitrogen cycle; NO;**  
50 **N<sub>2</sub>O; prior and posterior probability distributions for parameters; uncertainty analysis;**  
51 **water cycle**

52

53

## 54 **1. Introduction**

55

### 56 1.1 *Rationale*

57 Various recent reviews have assessed the evidence for impacts of environmental change on  
58 European forests (Hyvönen et al., 2007; Kahle et al., 2008; Luysaert et al., 2010). Most  
59 studies have focused on changes in growth and carbon balance, but the importance of the  
60 interaction with the nitrogen cycle is increasingly recognised (de Vries et al., 2009; Sutton et  
61 al., 2008; Van Oijen et al., 2008a; Van Oijen et al., 2004). Research programmes to measure  
62 and model emissions of nitrogenous greenhouse gases from European forests and other  
63 ecosystems have been set up (Sutton et al., 2007).

64 The measurement of nitrous oxide (N<sub>2</sub>O) and nitric oxide (NO) emissions from forest  
65 soils is hampered by the large spatial and temporal heterogeneity in the fluxes, and modelling  
66 these processes is still limited by availability of data (Kesik et al., 2005). Moreover, the  
67 relevant underlying mechanisms have not yet been clarified fully, and large uncertainties are  
68 present in both data and models. Available data sets not only suffer from random  
69 measurement error, but also from systematic errors associated with the positioning of  
70 measurement chambers in the field and their functioning (Butterbach-Bahl et al., 2002; Kroon  
71 et al., 2010). When modelling the systems, there is uncertainty about how to represent  
72 processes, i.e. model structural uncertainty (de Bruijn et al., 2009). Furthermore, there is  
73 uncertainty about environmental drivers and parameter values.

74 To improve the applicability of models to the analysis of the greenhouse gas balance  
75 of forests, these uncertainties need to be quantified and reduced. Probabilistic methods of  
76 model-data fusion or data-assimilation have come to the fore in recent years, and offer the

77 prospect of improved data use and uncertainty quantification (Fox et al., 2009; Wang et al.,  
78 2009). Because these methods are applications of probability theory, they require all  
79 uncertainties – in data, model inputs and model structure - to be expressed in the form of  
80 probability distributions. Bayes' Theorem can then be employed to update the distributions  
81 when new information becomes available.

82 In biogeochemical modelling, most Bayesian applications have focused on  
83 parameterisation of individual models, with little attention for systematic errors in data and  
84 model structure. Wang et al. (2009) thus concluded, in a recent review on model-data fusion  
85 studies for terrestrial ecosystems, that there is a need for “developing an integrated Bayesian  
86 framework to study both model and measurement errors systematically”. The work presented  
87 here is intended to contribute to that goal.

88

## 89 1.2 *Towards a Bayesian framework for dynamic modelling in forest biogeochemistry*

90 We propose a framework which requires that multiple models are used in any given study,  
91 and which consists of three operations: (1) Bayesian calibration, (2) Bayesian model  
92 comparison, (3) Analysis of model-data mismatch.

93 The overarching objective of this paper is to demonstrate that this three-stage  
94 framework is an effective tool for the analysis of models in forest biogeochemistry. For that  
95 purpose, we used four different published models and one rich data set from the Norway  
96 spruce forest in Höglwald, Germany (Kreutzer et al., 2009). Most of the data were on the  
97 nitrogen cycle, with long time series of measurements of emissions of N<sub>2</sub>O and NO, but we  
98 also used time series of the carbon and water cycles in the form of soil respiration and soil  
99 water content.

100 Bayesian calibration, i.e. the first operation in the framework, consists of defining a  
101 prior probability distribution for a model's parameters and updating that distribution using the

102 data. The method has not often been applied to parameter-rich nonlinear process-based  
103 ecosystem models (Luo et al., 2009). One reason is the high computational demand  
104 associated with the technique, which is exacerbated by the long running time of the models.  
105 A second issue is the difficulty of quantifying uncertainties about random and systematic  
106 measurement errors. We show in this paper how both types of error can be accommodated in  
107 a Markov Chain Monte Carlo algorithm for Bayesian calibration.

108 Bayesian model comparison, the second operation used in the framework, aims to  
109 determine the extent to which the data support the different models. This is done by providing  
110 a probability distribution over models rather than parameter values. The attempt in this paper  
111 to assess whether Bayesian model comparison as a method can be useful for model selection  
112 purposes is, as far as we are aware, new for parameter-rich process-based ecosystem models.

113 Detailed analysis of model-data mismatch, the third operation in our framework, is  
114 not a common step in Bayesian model studies, which tend to focus on the probabilistic  
115 aspects of model behaviour rather than the internal structure of the models (Gelman and  
116 Shalizi, 2010). Bayesian calibration and model comparison effectively treat models as black  
117 boxes that convert parameter values into outputs, so this further analysis is needed to  
118 facilitate model improvement.

119 In summary, this paper aims to show the strengths and weaknesses of this three-  
120 operation Bayesian framework using a case-study with four models simulating the  
121 biogeochemistry of a Central European spruce forest.

122

## 123 **2. Materials and Methods**

124

### 125 2.1 *Data*

126 All data were taken from the Norway spruce (*Picea abies* L.) site at Höglwald, Germany,

127 latitude 48°30'N, longitude 11°10'E, altitude 540 m (Papen and Butterbach-Bahl, 1999).  
128 Trees were planted in 1907. Soil C and N were around 90,000 and 5000 kg ha<sup>-1</sup> (Kreutzer et  
129 al., 2009; Rothe, 1997). For the years 1985-1995, mean annual temperature was 7.9 °C,  
130 precipitation 888 mm, and atmospheric N-deposition as measured in the throughfall 39.4 kg  
131 N ha<sup>-1</sup> (Rothe, 1997). For 1975-1990, average global radiation was 11.3 MJ m<sup>-2</sup> d<sup>-1</sup> and wind  
132 speed 2.8 m s<sup>-1</sup> (data from JRC-Ispra, as cited by (Van Oijen et al., 2008b)).

133         The primary data from the site were in the form of time series of daily measurements  
134 of soil water content and soil emissions of N<sub>2</sub>O, NO and carbon dioxide (CO<sub>2</sub>) (see e.g. Wu et  
135 al., 2010). Measurements at the Höglwald Forest are continuous throughout the year with  
136 fluxes being available in sub-daily resolution. However, daily mean values were used here for  
137 various years between 1994 and 2003 (1994-1996 and 2001-2003). For use in the calibration,  
138 the data were aggregated to monthly averages (Fig. 1). For N<sub>2</sub>O and NO, we also calculated  
139 intra-annual quantiles of the frequency distribution of daily emission magnitudes (10, 50 and  
140 90%). Monthly averages and annual statistics were only calculated for months and calendar  
141 years with more than 75% coverage, no gap-filling being applied. The data transformations  
142 led to ten different time series of data being available for use in the Bayesian analysis, four  
143 with monthly averages and six with annual quantiles, and with a combined number of data  
144 points of  $n = 225$  (Table 1).

145

146 [Fig. 1 HERE]

147

## 148 2.2 Models

149 Four different deterministic process-based models of forest biogeochemistry were used in this  
150 study: BASFOR, COUP, DAYCENT and MoBiLE-DNDC (Table 2).

151         BASFOR is the simplest model in the group. It was designed to simulate the

152 interactive effects of changes in N-deposition, atmospheric [CO<sub>2</sub>] and climate on the carbon  
153 balance of forests (Van Oijen et al., 2010a; Van Oijen and Thomson, 2010). The model has  
154 been subjected to Bayesian calibration before, using data from the United Kingdom (Van  
155 Oijen et al., 2005).

156 COUP and MoBiLE-DNDC are the two most complex models in the group. Both  
157 models were originally designed with special focus on soil processes, but recent versions of  
158 the models simulate the whole ecosystem. MoBiLE-DNDC calculates soil microclimate and  
159 hydrology, plant growth and plant-soil interactions, biogeochemical processes of the C and N  
160 cycle in soils, microbial growth and subsequent trace gas emissions. The core functionality  
161 follows the concepts developed in the DNDC suite of models (Li et al., 2000; Li et al., 1992;  
162 Werner et al., 2007). COUP was subjected to uncertainty quantification by Klemedtsson et al.  
163 (2008) and Svensson et al. (2008). The version of the model used in this paper, referred to as  
164 CoupModel, includes an N-flux submodel taken from the PnET-N-DNDC model (Norman et  
165 al., 2008). A preliminary uncertainty assessment was also carried out for MoBiLE-DNDC (de  
166 Bruijn et al., 2009).

167 DAYCENT, a model of intermediate complexity, traces its origins to the grassland  
168 soil model CENTURY (Parton et al., 1993), but like the previous two models it has  
169 developed into a full model for various ecosystems (Del Grosso et al., 2001), of which a  
170 small part was subjected to Bayesian calibration before (Yeluripati et al., 2009). The model  
171 version used in this paper, referred to as DailyDAYCENT, uses a daily time step for all  
172 processes.

173

### 174 2.3 *Parameter screening*

175 In the case of the simplest model, BASFOR, no parameter screening was applied, so all its 48  
176 model parameters and initial constants for state variables were included in the Bayesian

177 calibration. Initial constants were not included in the calibrations of the other three models.

178 The COUP model was subjected to informal screening, based on previous work  
179 involving the same model and experimental site (Norman et al., 2008). The fraction of  
180 COUP's output variability that was caused by uncertainty regarding the selected parameters  
181 was not quantified.

182 Calibration parameters for DAYCENT were selected using Morris screening (Morris,  
183 1991) on average model outputs for soil water content and emissions of N<sub>2</sub>O, NO and CO<sub>2</sub>.  
184 DAYCENT has over 300 parameters but only 214 were subjected to Morris screening  
185 because the majority of the about hundred parameters of the soil water dynamics module  
186 were known to contribute about equally to the overall uncertainty, precluding identification  
187 of a subset of essential parameters. Morris screening is a global parameter sensitivity  
188 analysis, i.e. it explores combinations of parameter values across parameter space rather than  
189 just in the neighbourhood of a default parameterisation. Compared to other global sensitivity  
190 analyses, the method requires relatively few model runs (proportional to the number of  
191 parameters) which permits its use even for models with long runtimes. A subset of 17  
192 DAYCENT parameters was selected in this way. The  $r^2$  of the relations between model  
193 output variability (from runs where all parameters were varied) and the selected parameters  
194 was 0.88-0.98 for the three emission rates and 0.20 for water content, the latter value  
195 reflecting the difficulty in selecting key parameters for soil water dynamics in this model.

196 In the case of MoBiLE-DNDC, uncertainty in the 67 parameters of the soil chemistry  
197 submodel was considered. These were mainly parameters used to adapt or scale physical or  
198 chemical processes observed by lab studies to real world conditions, semi-empirical ratios,  
199 and k-values (decay rates for the various litter and microbial pools of varying  
200 decomposability). A substantial number of those parameters also describe Michaelis-Menten  
201 kinetics for the microbial turnover processes. The Morris screening method (Morris, 1991)

202 was applied to MoBiLE-DNDC using the same selection criteria as for DAYCENT, and 26  
203 parameters were selected which together accounted for more than 60% (NO-emission) and  
204 >90% (N<sub>2</sub>O- and CO<sub>2</sub>-emission, soil water content) of model behaviour.

205

#### 206 2.4 *Bayesian calibration*

207 Bayesian calibration is the application of probability theory to parameter estimation (Jaynes,  
208 2003; Sivia, 2006). The method finds increasing use in ecological modelling (Ogle, 2009;  
209 Ogle and Barber, 2008; Van Oijen et al., 2005). Uncertainty about parameters is represented  
210 as a joint probability distribution for the possible parameter values. Bayes' Theorem is used  
211 to determine how this distribution changes in the light of new data:

212

$$213 \quad P(\theta|D) \propto P(\theta) P(D|\theta) \quad (1)$$

214

215 Where  $P(\theta)$  and  $P(\theta|D)$  are the prior and posterior distributions for the parameters  $\theta$ , i.e.  
216 before and after conditioning on the data. The factor that modifies the prior,  $P(D|\theta)$ , is the  
217 likelihood function, which is the probability of the data for a given  $\theta$ . A formal likelihood  
218 function, integrating to unity in data space, needs to be used to be consistent with the  
219 probability calculus, allowing Bayes' Theorem to be applied (Rougier, (in press)).

220 The prior probability distribution for the parameters of a model,  $P(\theta)$ , reflects a  
221 modeller's uncertainty about parameter values before using the data. This uncertainty is  
222 subjective, and there was no effort in this study to impose any harmonisation on the priors for  
223 the four different models (except for data-scaling factors, discussed later). All modellers  
224 assigned prior distributions that could be written as the product of independent marginal  
225 distributions for individual parameters, but different types of marginal distribution were used:  
226 beta for BASFOR and uniform for the other three models.

227 The likelihood function  $P(D|\theta)$  is the probability of the data  $D$  (the 225 data points)  
 228 given model output generated by parameter vector  $\theta$ . It accounts for possible measurement  
 229 error. The same likelihood function was used for all models to allow formal Bayesian model  
 230 comparison. As described above, the calibration data were in the form of time series of ten  
 231 different variables, six for annual quantiles and four for monthly averages of daily  
 232 measurements (Table 2). Estimates for the uncertainty of these variables, both random and  
 233 systematic, were elicited from the data-providers: co-authors Butterbach-Bahl, Kiese and  
 234 Werner. Uncertainty about random measurement error is often represented by the use of  
 235 independent Gaussian distributions for the data points. However, the squared exponential in  
 236 the Gaussian tends to cause asymptotic collapse of the parameter distribution even with  
 237 moderate amounts of data (Clark, 2005), and may represent an overestimate of their  
 238 information content. We therefore used the more heavy-tailed function proposed by Sivia  
 239 (2006):

240

$$241 \quad P(D | \theta) = \prod_{i=1}^{225} \frac{1}{\sigma_i \sqrt{2\pi}} \frac{1 - \text{Exp}(-R_i^2 / 2)}{R_i^2} \quad (2)$$

242

243 Where  $\sigma_i$  is a measure of the uncertainty about random error of the  $i$ -th data point, and  $R_i$  is  
 244 the difference between model output and  $i$ -th data point, divided by  $\sigma_i$ . The values of the  $\sigma_i$   
 245 were considered to be specific to the type and magnitude of the data points, with relative  
 246 values of 0.2 for the medians (Q50), 0.3 for the tail-quantiles (Q10 and Q90), and 1.0 for the  
 247 monthly averages. Besides random measurement error, the data-points were considered to be  
 248 subject to possible systematic error, which could result from unrepresentative positioning of  
 249 the soil measurement chambers or errors in instrument calibration. This was implemented by  
 250 means of four multiplicative data-scaling factors  $\gamma_j$ , one each for  $N_2O$ ,  $NO$ ,  $CO_2$  and water.

251 As in other recent studies (Raupach et al., 2005), we treated the  $\gamma_j$  as additional parameters to  
252 be calibrated. We considered that errors larger than a factor 2 would be very unlikely, but  
253 otherwise no assumptions about the systematic data-scaling factors were made. We therefore  
254 used an uninformative Jeffrey's prior (Jaynes, 2003; Sivia, 2006) for the marginal prior  
255 distribution for each of the four factors:

256

$$257 \quad P(\gamma_j) \propto 1 / \gamma_j \quad (3)$$

258

259 The Jeffrey's priors for the uncertainty about multiplicative error thus are log-uniform  
260 distributions on the interval  $[\frac{1}{2}, 2]$ . Keats et al. (2009) also used multiplicative errors, for air  
261 pollution data in Bayesian calibration of an atmospheric transport model, but used log-normal  
262 distributions instead. As Keats et al. (2009) argued, multiplicative error is often the natural  
263 choice for measurements of non-negative quantities.

264 For each of the four models, the Bayesian calibration was carried out by means of  
265 Markov Chain Monte Carlo sampling (MCMC), using the Metropolis algorithm (Metropolis  
266 et al., 1953; Van Oijen et al., 2005), but model-specific choices were made of proposal  
267 distribution and method of testing chain convergence.

268 Before and after each model's calibration, a preliminary parameter sensitivity analysis  
269 was carried out by calculating the partial correlation coefficients (PCC) for the relationships  
270 between individual parameters and the average simulated values of N<sub>2</sub>O, NO, CO<sub>2</sub> and water.  
271 In contrast to the ordinary correlation coefficient ( $r$ ), the PCC calculates the association  
272 between parameter and output after correcting for the linear effects of the other parameters.

273

## 274 2.5 *Bayesian model comparison*

275 The formal Bayesian model comparison consisted of quantifying the relative probabilities of

276 correctness of the four models, under the assumption that at least one of them was a correct  
277 model for the data. The comparison was, like the calibration described above, based on  
278 application of Bayes' Theorem (Kass and Raftery, 1995):

279

$$280 \quad P(M/D) \propto P(M) P(D|M) = P(M) \int P(D|M, \theta) P(\theta) d\theta \quad (4)$$

281

282 Where  $P(M)$  and  $P(M|D)$  are the prior and posterior distributions for the models  $M$ . In  
283 contrast to the parameter distributions, these are discrete distributions, over the four models  
284 in our comparison. As shown in the right-hand term, the factor that modifies the prior,  
285  $P(D|M)$ , is the integral of the likelihood function over the space spanned by the prior  
286 parameter distribution of each model. We refer to this integral as the model's 'integrated  
287 likelihood'. Another common name for this quantity is 'marginal likelihood' which expresses  
288 the fact that it is found by marginalising out the parameters  $\theta$ .

289 We *a priori* assigned equal probabilities to the different models of being correct, so  
290  $P(M)$  is uniform and the integrated likelihoods represented the relative probability for each  
291 model of being correct given the information in the data (Kass and Raftery, 1995). We  
292 quantified the integrated likelihoods as follows. For each model, 1000 parameter vectors were  
293 drawn from its prior parameter distribution. Comparison of the model outputs for these  
294 parameter vectors with the data yielded a sample of 1000 values of the likelihood, and the  
295 sample mean was taken as the estimate for the integrated likelihood of the model. Because  
296 any sampling-based method is subject to sampling error (McCulloch and Rossi, 1992), we  
297 additionally calculated the integrated likelihoods using the method suggested by Kass and  
298 Raftery (1995), as the harmonic mean of the sample generated by the MCMC.

299

## 300 2.6 Analysis of model-data mismatch

301 Besides the calculation of the likelihood function, the mismatch between model outputs and  
302 measurements was also quantified using more classical means. This was done separately for  
303 each of the ten output variables (Table 2) by calculating the Normalised Root Mean Square  
304 Error (NRMSE = square root of the average squared difference between model output and  
305 data, divided by the average of the data) and the squared correlation coefficient ( $r^2$ ). NRMSE  
306 and  $r^2$  are distributed quantities, because they depend on the parameterisation, so they were  
307 calculated across the range of prior and posterior parameter distributions.

308 Additional analysis was carried out for just the modes of the prior and posterior  
309 parameter distributions. We ran each model with the two modal parameter vectors and  
310 calculated the Mean Squared Error (MSE) for each of the four time series of monthly  
311 averages. The MSE-values were then decomposed as suggested by Kobayashi and Salam  
312 (2000):

313

$$314 \quad + (\sigma_M - \sigma_D)^2 + 2(\sigma_M \sigma_D) (1-r) \quad (5)$$

315

316 Where M is a simulated time series consisting of monthly averages of N<sub>2</sub>O, NO, CO<sub>2</sub> or  
317 water content, D is the matching data,  $\sigma_M$  and  $\sigma_D$  are their standard deviations, and  $r$  is the  
318 correlation between the two. The decomposition consists of three terms, which can be  
319 interpreted as measures for model-data mismatch due to bias, variance error and phase shift  
320 (Kobayashi and Salam, 2000).

321

322

### 323 **3. Results**

324

#### 325 **3.1 Bayesian calibration**

326 All four models were calibrated using the same MCMC-algorithm, i.e. Metropolis sampling.  
327 Burn-in and convergence were determined visually, by each modelling group separately, but  
328 an additional analysis of the Markov chains was carried out to confirm that parameter  
329 distributions had properly stabilised. The analysis was based on the fact that, after a chain  
330 reaches convergence, subsequent distinct and sufficiently long sub-chains should have similar  
331 sample means and variances. We compared the first and second halves of the chains after  
332 burn-in. The results showed that convergence was adequate for BASFOR and DAYCENT,  
333 with all parameters having similar means and variances in the two halves. However, for  
334 COUP and MoBiLe-DNDC, some parameters had not stabilised to the same extent, so the  
335 posterior parameter distributions for these two models were likely less accurate.

336 The calibration modified the means and reduced the variances of most marginal  
337 distributions. The average variance reduction for process-parameters was small in BASFOR  
338 and DAYCENT (3%, 4%), but larger in COUP and MoBiLE (26%, 29%). The data-scaling  
339 factors  $\gamma_j$  showed greater variance reductions except for soil water content (Fig. 2).

340

341 [Fig. 2 HERE]

342

343 Parameter uncertainty induced output uncertainty. The degree of prior output  
344 uncertainty was assessed by determining the quantiles of the output distributions (Table 3).  
345 For all models except DAYCENT, prior Q95 was one or two orders of magnitude larger than  
346 Q5 for the eight nitrogenous emission variables (Table 3). DAYCENT was already strongly  
347 constrained by its prior parameter distribution. For all models, soil respiration was *a priori*  
348 slightly more constrained than the N-emissions, whereas soil water content, which was the  
349 only state variable in the set of ten output variables, was most narrowly delimited. Overall,  
350 prior ranges were widest for BASFOR. The calibration had only little impact on the output

351 distributions for soil water content and almost no effect on the soil respiration distributions  
352 (Table 3). However, the posterior outputs for the nitrogenous emission variables were much  
353 more narrowly constrained than the prior distributions, with posterior Q95/Q5 ratio's ranging  
354 from 2-5 (BASFOR, see also Fig. 1), 2-3 (COUP) and ~1.5 (MoBiLE-DNDC) (Table 3).  
355 DAYCENT was the exception with posterior ratio's that were similar to the prior.

356 Using the samples from both the prior and posterior parameter distributions, we  
357 calculated the partial correlations between individual parameters and outputs. Posterior PCC-  
358 values tended to be higher than prior values for BASFOR and COUP, whereas the calibration  
359 decreased PCC-values for DAYCENT and MoBiLE-DNDC. However, for all models and  
360 output variables, the PCC-based ranking of the parameters changed little, so we restrict  
361 ourselves to reporting on the posterior values.

362 In the case of BASFOR, only one parameter was strongly correlated ( $|PCC| > 0.5$ )  
363 with N<sub>2</sub>O- and NO-emission: the soil water content at which both emissions are equal. Soil  
364 respiration was strongly correlated with the parameters that govern decomposition rate of  
365 organic matter, and also with the light-use efficiency. Soil water content was mainly  
366 correlated with specific leaf area and leaf longevity, both of which affect the active surface  
367 area for transpiration. The results for COUP were similar. N-emissions were also mainly  
368 correlated with the N<sub>2</sub>O-NO balance, soil respiration with decomposition rates and light-use  
369 efficiency, and water content with specific leaf area. In the case of DAYCENT, no individual  
370 parameters were *a posteriori* strongly correlated with model outputs, although there had been  
371 some strong correlations with the prior parameter distribution (i.e. the NO<sub>3</sub>-N<sub>2</sub>O conversion  
372 efficiency for N-emissions and the leaf area ratio for soil water content). In the case of  
373 MoBiLE-DNDC, no parameters were strongly correlated with N<sub>2</sub>O-emission, but the K<sub>m</sub>-  
374 value for NO<sub>2</sub> did have a high PCC with NO-emission. Soil respiration and water content  
375 were both mainly correlated with the parameter that scales decomposition of active organic

376 substance as a function of soil porosity.

377

### 378 3.2 *Bayesian model comparison*

379 The log-transformed integrated likelihood values, calculated from samples from the prior,  
380 were as follows (between brackets the values from the alternative calculation using the  
381 harmonic mean of the sample from the posterior): BASFOR: -661.7 (-654.7), COUP: -663.5  
382 (-651.2), DAYCENT: -738.5 (-761.2), MoBiLE-DNDC: -657.0 (-758.9). For comparison: a  
383 parameter vector whose model outputs would have gone exactly through the 225 data points,  
384 would have had a log-likelihood of -581.2. Both methods of calculating the integrated  
385 likelihoods showed that the data provided greater support for BASFOR and COUP than for  
386 DAYCENT. The two estimates of the integrated likelihood for MoBiLE-DNDC differed  
387 strongly, so it is less clear how plausible this model is.

388

### 389 3.3 *Analysis of model-data mismatch*

390 First, the data were compared with the ranges of model output uncertainty induced by the  
391 parameters. All time series averages of measurements were in the central intervals of the  
392 prior output distributions, between the 5% and 95% quantiles, except for soil respiration as  
393 predicted by BASFOR and MoBiLE-DNDC, and the lower quantiles of daily N<sub>2</sub>O emission  
394 rates as predicted by COUP (Table 3).

395 Although the distributions of simulated time series averages were found to cover the  
396 data fairly well, inspection of the time series themselves revealed considerable differences  
397 between models and measurements (Fig. 3). For example, none of the models was able to  
398 reproduce the large peak in N<sub>2</sub>O-emissions in early 1996 after a strong freeze-thaw event  
399 (Papen and Butterbach-Bahl, 1999), neither with the mode of the prior, not with the posterior  
400 mode (Fig. 3). This is likely to be a consequence of incomplete process representation in the

401 models, although a more recent version of MoBiLE-DNDC showed a possible way to  
402 account for the freeze-thaw effect (de Bruijn et al., 2009). Despite these remaining  
403 differences between model outputs and data, overall the calibration was able to remove much  
404 of the mismatch for the N-emissions and, to lesser extent, for CO<sub>2</sub>-emission (compare the  
405 lower ‘posterior’ panels of Fig. 3 to the upper ‘prior’ ones, and see also the increased  
406 likelihoods depicted in Fig. 4). Note that some of the reduced bias in the posterior results was  
407 due to the impact of calibration on data-scaling factors (Fig. 2) rather than on model  
408 parameters. There was no apparent improvement in the simulation of soil water content,  
409 which reflected the fact that for this variable the least amount of information was available  
410 (Table 2). In one specific case, simulation of soil water by model DAYCENT using the mode  
411 of its posterior distribution, bias was increased relative to the prior mode (Fig. 3), but this  
412 result was not representative of the full posterior distribution for the water data-scaling factor  
413 ( $\gamma_{\text{WATER}}$ ) of this model (Fig. 2). However, it does suggest that a useful – but for this model  
414 computationally demanding - additional step in the procedure would have been to determine  
415 the mode of the posterior distribution by targeted optimisation, rather than relying on the  
416 parameter sample generated by the MCMC.

417

418 [Fig. 3 HERE]

419

420 Whereas for the prior output distributions of the models most time-series averages  
421 were in the central Q5-Q95, the same did not apply to the posterior distributions. Most time-  
422 series averages were to be found in the upper tails (>Q95) of the posterior output distributions  
423 (Table 3, and compare also the examples for BASFOR in Fig. 1). There were differences,  
424 however, in how the likelihood distributions responded to the calibration, as can be seen from  
425 the posterior distributions of likelihoods (Fig. 4).

426

427 [Fig. 4 HERE]

428

429 For each model and each parameter vector sampled from the prior and posterior  
430 distributions, we compared the simulated time series of the ten output variables with the  
431 corresponding data, by calculating the correlation coefficient ( $r$ ) and the normalised root  
432 mean square error (NRMSE). Each parameter distribution thus induced ten different  
433 distributions of  $r$  and NRMSE, of which we show the 5, 50 and 95% quantiles (Table 4). The  
434 posterior values of the quantiles for  $r$  are often not improvements over the prior, except for  
435 Q5. So calibration tended to remove only the parameter vectors with the poorest output-data  
436 correlation. In contrast, NRMSE was improved for almost every quantile of every variable in  
437 each model (Table 4).

438 The MSE-decompositions for time series with monthly averages, both for the prior  
439 and posterior parameter modes, are shown in Fig. 5. Phase shift, variance error and bias were  
440 reduced to different extent for the different models.

441

442 [Fig. 5 HERE]

443

## 444 **4. Discussion**

445

### 446 4.1 *Bayesian calibration: methodological issues*

447 Bayesian calibration uses data to update the joint probability distribution for a model's  
448 parameters. The Bayesian approach allows for non-Gaussian distributions for both parameter  
449 uncertainty and measurement error. Our calibration was therefore based on sampling by  
450 means of MCMC rather than on matrix inversion methods. This in turn allowed us to include

451 systematic data error in the calibration, rather than having to estimate error terms in a first  
452 separate step, as was done for example by Michalak et al. (2005), using maximum-likelihood  
453 estimation.

454         Although the theory is straightforward, it is easy to overestimate the information  
455 content of any dataset, and this may lead to unsupported changes in the parameter  
456 distributions. For example, when the common assumption is made that each new data point  
457 adds independent new information to the calibration, the parameter distributions will  
458 asymptotically collapse with sample size (Clark, 2005). Modellers thus need to elicit realistic  
459 assessments of measurement uncertainty from the data-providers (Moala and O'Hagan,  
460 2010). This issue was important in the case-study presented here because the dataset was  
461 fairly large ( $n=225$ ), covering times series of four variables. We applied four techniques to  
462 ensure a realistic, albeit subjective assessment of the information content of the data: (1)  
463 using the monthly temporal scale as the one at which the models were supposed to be  
464 applicable, together with the frequency distribution of daily emission events, (2) allowing for  
465 random errors in the data, (3) allowing for systematic errors in the data by the use of the four  
466 scaling factors, (4) using a heavy-tailed likelihood function (Sivia, 2006). The adjustment of  
467 temporal scale is a common technique in atmospheric physics, applied whenever models  
468 produce more smooth results than measurements and thereby induce apparent correlations  
469 between measurement errors (Prinn, 2000). We considered the implementation of these four  
470 techniques to be partly the responsibility of the data-providers (in the case of random and  
471 systematic errors) and partly that of the modellers and data-providers together (temporal scale  
472 and likelihood function). Using this approach, parameter uncertainties were reduced  
473 markedly but the distributions did not collapse. The techniques applied are generic and may  
474 be widely applicable to calibration of complex dynamic models using long time series.

475

476 4.2 *Bayesian calibration: impact on parameter uncertainty of the forest models*

477 Parameter uncertainties were reduced strongly compared to the prior, and the likelihood  
478 distributions were shifted toward higher values for all four models (Fig. 4). The degree of  
479 uncertainty reduction varied between the models, as did the balance between changing data-  
480 scaling parameters and process parameters. Although our dataset included measurements of  
481 the three major biogeochemical cycles (nitrogen, carbon and water), there was still some lack  
482 of balance because about 75% of the data points were for emissions of N<sub>2</sub>O and NO (Table  
483 1). Therefore, most of the improvement of model behaviour (reduction of likelihood and  
484 NRMSE, increase in  $r$ ) was for these variables. For all models, the parameters that were  
485 changed the most were related to the soil nitrogen dynamics.

486 The results of our preliminary parameter sensitivity analysis, consisting of calculating  
487 partial correlations with the different outputs, need to be interpreted with care. A high value  
488 of the PCC for a specific parameter-output combination suggests that the parameter – within  
489 its range of uncertainty - strongly affects the output. Therefore knowledge about the process  
490 governed by that parameter is key to understanding variability in the output. The opposite  
491 may not be true: a strong but non-linear effect may yield a low PCC. Also, the importance of  
492 a parameter is not an intrinsic property: it depends on the distribution of that parameter.  
493 Whenever Bayesian calibration reduces the variance of a parameter, the contribution of that  
494 parameter to output variability is expected to decrease. However, PCC-analysis is not a  
495 variance-decomposition method (Saltelli et al., 2000), so it may not be able to show that  
496 effect, and indeed, for two out of the four models posterior values of PCC were generally  
497 larger than the prior values.

498 With these caveats, the results of the PCC-analysis did reveal commonalities between  
499 the models. Across the models, N-emissions were mainly correlated with N<sub>2</sub>O-NO  
500 partitioning, soil respiration mainly with decomposition but also tree productivity, and soil

501 water content with leaf area dynamics. These agreements between the models suggest that  
502 some of the differences between complex process-based models may not overly affect their  
503 behaviour. The models compared here all represent, albeit in very different ways, the linkages  
504 between the ecosystem C-, N- and H<sub>2</sub>O-cycles. Therefore there are inevitable similarities in  
505 the overall feedback structure of the models, imposed by constraints of stoichiometry and  
506 mass-balance of the three biogeochemical cycles. These similarities may outweigh details of  
507 process representation and parameterisation (Van Oijen et al., 2004; Van Oijen et al., 2010b).

508         The prior PCC-values were, as expected, indicative of which parameters were most  
509 informed by the data in the subsequent calibration. For all four models, the relative decrease  
510 of marginal variance tended to be greatest for those parameters that had a strong *a priori*  
511 correlation with N<sub>2</sub>O-emission, the output variable for which the largest number of data  
512 points were available. The parameter variance reduction accounted for by the prior PCC-  
513 values ranged from 25% (COUP, MoBiLE) to 29% (BASFOR) and even 35% (DAYCENT).  
514 PCC-analysis might thus play a useful role in parameter screening before calibration.

515

#### 516 4.3 *Bayesian model comparison*

517 Comprehensive model comparison requires taking into account parameter uncertainty. A  
518 complex model might, in principle, be able to predict complex biogeochemical time series  
519 more closely than a simple model. But if it is unclear what the parameterisation of the  
520 complex model for good prediction should be, then its predictive capacity is reduced. A  
521 simple model whose parameters are well-known might then perform better. Regarding model  
522 complexity, there is a trade-off between the need to represent the intricacies of the real world  
523 and the need to minimise parameter uncertainty. We thus need a method for comparing the  
524 behaviour of the models not just at the modes of the parameter distributions or at the  
525 maximum likelihood estimates, but across their whole parameter distributions. Bayesian

526 model comparison is such a method. This method has been used before in environmental  
527 modelling, in the elegant study by Tuomi et al. (2008) who compared different functions  
528 describing the impact of temperature on soil respiration, but to our knowledge this is the first  
529 application to complex process-based ecosystem models. Formally, the relative magnitudes  
530 of the integrated likelihoods equate to relative probabilities for the individual models of being  
531 correct, conditional on a correct model being present in the comparison. However, in  
532 environmental modelling all models are incorrect in some way, so we prefer to use the  
533 integrated likelihoods as a guide towards plausible model structures rather than as  
534 probabilities of correctness (Gelman and Shalizi, 2010; Kass and Raftery, 1995).

535         Bayesian model comparison requires that a common likelihood function is used with  
536 all models – because the criterion for comparison is the integrated likelihood of the models,  
537 where the integration is over the prior parameter distribution. Therefore only those data can  
538 appear in the likelihood function, which are part of the output set of each model. For  
539 example, in the current study, we could not include vertical profiles of soil temperature in the  
540 likelihood function because the simplest model BASFOR has only one soil compartment.  
541 Three of the models had remarkably similar values of the integrated likelihood, the exception  
542 being the low value for DAYCENT, which was thus identified as the least plausible model.  
543 There was only a slight preference for MoBiLE-DNDC. We analysed these prior integrated  
544 likelihood values further by considering the underlying distributions of likelihoods associated  
545 with the four different categories of output variable ( $N_2O$ , NO, respiration, water) (Fig. 4).  
546 The lower value of the integrated likelihood for DAYCENT can be seen to be mainly due to  
547 poorer performance for  $N_2O$ . The similarity between the integrated likelihoods for the other  
548 models was seen to extend to the underlying distributions of category-wise likelihoods (Fig.  
549 4).

550

551 4.4 *Analysis of model-data mismatch*

552 The NRMSE and  $r$  statistics provided useful additional information beyond the formal  
553 Bayesian model comparison. When going from prior to posterior,  $r$  did not improve (apart  
554 from Q5) in any of the models, but NRMSE improved throughout (Table 4). The calibration  
555 thus was more successful in reducing the average magnitude of the differences between  
556 simulations and measurements, than in aligning the distribution of the outputs over time.  
557 There clearly remain difficulties for all models in simulating the large interannual variation in  
558 nitrogenous emission characteristics. The models were also similar in that the posterior  
559 NRMSE was lowest for soil water content and highest for the monthly values of N<sub>2</sub>O-  
560 emission (Table 4). It was easier to simulate water than nitrogen dynamics.

561 The decomposition of the MSE for the modes of the prior and posterior parameter  
562 distributions gave further information. The MSE-decomposition is only possible for long time  
563 series, i.e. the monthly data (Table 1). There were not enough annual quantiles available to  
564 allow the reliable estimation of the variance and phase shift terms. The calibration reduced  
565 the MSE for the parameter modes of all four models, and all four categories of output  
566 variables (Fig. 5), confirming the effectiveness of the calibration. However, the analysis  
567 revealed large differences between the models. The simplest model, BASFOR, had the  
568 highest variance error for N<sub>2</sub>O and NO. This suggests that a simple model may not be able to  
569 respond quickly enough to changes in the environment that affect nitrogenous emissions. The  
570 low integrated likelihood of DAYCENT has already been attributed to poor simulation of  
571 N<sub>2</sub>O-emission (Fig. 4), and the MSE-decomposition showed that this was mainly due to a  
572 large phase shift for N<sub>2</sub>O emission (Fig. 5). In fact, DAYCENT had very low bias and  
573 variance error for N<sub>2</sub>O, so it was able to capture general characteristics (mean, ‘peakiness’) of  
574 the time series of emission very well, but not the timing of emission events.

575 Most data were in the central Q5-Q95 ranges of the prior output distributions of the

576 models (Table 3). However, perhaps surprisingly, most data were to be found in the upper  
577 tails ( $>Q95$ ) of the posterior output distributions (Table 3). This is because of a trade-off  
578 between the different variables in the calibration. Such trade-off is inevitable with models  
579 that are imperfect and cannot capture the whole range of behaviour of all variables  
580 simultaneously. Moreover, there is also the distinct possibility of systematic error in the  
581 measurement of the different categories of output variables. Where the models were unable to  
582 reconcile all the data, the calibration tended to modify the settings of the scale parameters  
583 (Fig. 2). It is important to establish whether the likely underlying cause of the revision of the  
584 scale factors was indeed systematic measurement error or model structural error. One way of  
585 attempting this is to consider the differences between the models in their posterior estimates  
586 for the scaling factors. *A posteriori*, all models suggested that the measurements for CO<sub>2</sub>-  
587 emissions were unrealistically high (Fig. 2:  $\gamma_{CO_2}$  mostly  $<1$ ), but only the BASFOR model  
588 suggested the same for the N-emissions ( $\gamma_{N_2O}$  and  $\gamma_{NO} \ll 1$ ). There thus is some doubt about  
589 the CO<sub>2</sub>-measurements, and likewise about the capacity of BASFOR to simulate N-dynamics.

590 Note that, in our approach, the data-scaling factors are intended to represent the  
591 idiosyncrasies of specific datasets, so we cannot expect the calibrated scaling values to apply  
592 elsewhere. However, if calibration of a model using data from multiple sites were shown to  
593 consistently lead to scaling factors different from unity, we would expect the error to be  
594 predominantly the model's rather than that of the data. Here we only had data from one site  
595 available, so no strong conclusions can be drawn.

596

#### 597 4.5 *The Bayesian framework*

598 The three-operation Bayesian framework proposed here - calibration, comparison, analysis of  
599 model-data mismatch – was shown to work well in this study. Most of the techniques  
600 employed in each of these operations are novel in their application to complex dynamic

601 models of forest biogeochemistry, in particular in combination with the use of data from  
602 processes that vary strongly over time (Luo et al., 2009). The application of the framework  
603 showed that model parameter uncertainty could be reduced in all models, irrespective of their  
604 level of complexity. However, it also showed that the models still suffer from structural  
605 deficiencies, even if we allow for the possibility of errors in the data, and that the deficiencies  
606 are stronger for the nitrogen cycle than for the carbon and water cycles.

607         There are limitations associated with this study and a caveat has to be made regarding  
608 the use of the results. We only used data from one site, Höglwald, and the general  
609 applicability of the models to European forests needs to be tested with data from other sites.  
610 When data from these new sites become available, the posterior distributions found here will  
611 become prior distributions in further calibration. Another limitation of the study was the use  
612 of parameter screening, which was considered necessary for the three models with the highest  
613 computational demand, but there may be environmental conditions under which simulated  
614 forest biogeochemistry is sensitive to the excluded parameters. Also, uncertainty concerning  
615 environmental drivers such as weather conditions was ignored as it was expected to be small  
616 compared to the structural and parameter uncertainties. This assumption needs to be verified.

617         Wang et al. (2009) called for the development of an integrated Bayesian framework  
618 that can account for the different sources and types of error arising in environmental  
619 modelling. The Bayesian framework proposed here is an integrated one in the sense that its  
620 three operations were linked methodologically and in that its three operations provide  
621 complementary information. Methodologically, the Bayesian calibration made use of the  
622 same likelihood function as the Bayesian model comparison. Calibration was also linked to  
623 comparison in the use of MCMC to explore parameter space, with the thus generated sample  
624 being used in the model comparison for estimating the integrated likelihood. Because this  
625 estimate, based on the harmonic mean of the sampled likelihoods, can be unstable (Chib and

626 Jeliaskov, 2001), we also used the method of directly sampling from the prior.  
627 Methodological links further existed between the calibration and the analysis of model-data  
628 mismatch, in that the sample generated by the calibration was used to calculate the NRMSE  
629 across the posterior parameter distribution, rather than just for one parameter vector.

630 More importantly than these methodological links, the three operations in the  
631 framework also complemented each other in how they help improve the modelling.  
632 Calibration reduced parameter uncertainty, model comparison reduced uncertainty about the  
633 relative plausibility of the different models, and the analysis of model-data mismatch showed  
634 which parts of the models needed most improvement, which in this case was the nitrogen  
635 dynamics.

636

637

## 638 **5. Conclusions**

- 639 • Bayesian calibration can be used to reduce parametric uncertainty of complex  
640 dynamic models for forest biogeochemistry.
- 641 • Bayesian calibration allows for the use of datasets that contain long time series of gas  
642 emissions with high intra- and interannual variability, and with both random and  
643 systematic error.
- 644 • Data need to be compared with models at the appropriate temporal scale. This may  
645 involve, as shown here, monthly averaging and the calculation of annual frequency  
646 distributions. These transformations, and the use of heavy-tailed likelihood functions  
647 that account for uncertainty about random and systematic measurement errors, can  
648 help prevent collapse of the parameter distributions in the calibration.
- 649 • Bayesian model comparison can be used to calculate the relative conditional  
650 probabilities of models being correct, irrespective of the type and complexity of the

651 considered models.

- 652 • Bayesian model comparison treats models as black boxes, so it can only identify  
653 which models are implausible, but it cannot identify any specific model deficiencies.
- 654 • Analysis of model-data mismatch can help identify model weaknesses by  
655 decomposition of the MSE and by showing how the NRMSE and the correlation  
656 coefficient  $r$  vary for the different processes simulated by the models.
- 657 • Together, the three operations of Bayesian calibration, Bayesian model comparison,  
658 and analysis of model-data mismatch, constitute a promising framework for  
659 uncertainty reduction and improvement of complex dynamic models in forest  
660 biogeochemistry.
- 661 • This was confirmed by the case-study analysed here, in which four different  
662 parameter-rich process-based models of forest biogeochemistry were confronted with  
663 long time-series of biogeochemical data. Parameter uncertainties were reduced in all  
664 models and the relative model plausibilities were quantified, with MoBiLE-DNDC  
665 having a slight preference over the other models. The simplest model, BASFOR, was  
666 shown to underestimate variance of nitrogenous emissions even after calibration. The  
667 model of intermediate complexity, DAYCENT, simulated the time series well but  
668 with large phase shift. COUP and MoBiLE-DNDC were able to remove most bias  
669 through calibration.
- 670 • The calibration not only reduced parameter and model uncertainty, but also identified  
671 possible systematic error in the measurement of soil respiration, to which all models  
672 assigned data-scaling factors less than unity with high posterior probability.
- 673 • There remain patterns in the data - in particular infrequent events of very high  
674 nitrogenous emission rate - that are unexplained by any of the models, even after  
675 calibration. Given the intensive exploration of parameter space in the calibration, this

676 is unlikely to be due to incorrect model parameterisation.

- 677 • The analysis showed that the models still suffer from structural deficiencies, even if  
678 we allow for the possibility of errors in the data. The deficiencies are stronger for the  
679 nitrogen cycle than for the carbon and water cycles.

680

681

## 682 **Acknowledgements**

683 We thank the European Union for financial support to carry out this work in projects  
684 NitroEurope and Carbo-Extreme, and we are grateful to our colleagues in these projects for  
685 discussion. J.B.Y. wishes to acknowledge the help of M. Richards in programming the  
686 routines for the Bayesian analysis of DAYCENT. We express our thanks to two anonymous  
687 reviewers for their constructive comments on the original manuscript.

688

689

## 690 **References**

691

692 Butterbach-Bahl, K., Rothe, A. and Papen, H., 2002. Effect of tree distance on N<sub>2</sub>O and CH<sub>4</sub>-fluxes from soils in  
693 temperate forest ecosystems. *Plant and Soil*, 240(1): 91-103.

694 Chib, S. and Jeliazkov, I., 2001. Marginal likelihood from the Metropolis-Hastings output. *Journal of the*  
695 *American Statistical Association*, 96(453): 270-281.

696 Clark, J.S., 2005. Why environmental scientists are becoming Bayesians. *Ecology Letters*, 8(1): 2-14.

697 de Bruijn, A.M.G., Butterbach-Bahl, K., Blagodatsky, S. and Grote, R., 2009. Model evaluation of different  
698 mechanisms driving freeze-thaw N<sub>2</sub>O emissions. *Agriculture, Ecosystems & Environment*, 133(3-4):  
699 196-207.

700 de Vries, W. et al., 2009. The impact of nitrogen deposition on carbon sequestration by European forests and  
701 heathlands. *Forest Ecology and Management*, 258(8): 1814-1823.

702 Del Grosso, S.J. et al., 2001. Simulated interaction of carbon dynamics and nitrogen trace gas fluxes using the

703 DAYCENT model. In: M. Schaffer, et al. (Editor), Modeling Carbon and Nitrogen Dynamics for Soil  
704 Management. CRC Press, Boca Raton, pp. 303-332.

705 Fox, A. et al., 2009. The REFLEX project: Comparing different algorithms and implementations for the  
706 inversion of a terrestrial ecosystem model against eddy covariance data. Agricultural and Forest  
707 Meteorology, 149(10): 1597-1615.

708 Gelman, A. and Shalizi, C.R., 2010. Philosophy and the practice of Bayesian statistics. Working paper,  
709 <http://arxiv.org/abs/1006.3868>.

710 Hyvönen, R. et al., 2007. The likely impact of elevated [CO<sub>2</sub>], nitrogen deposition, increased temperature, and  
711 management on carbon sequestration in temperate and boreal forest ecosystems. A literature review.  
712 New Phytologist, 173: 463-480.

713 Jaynes, E.T., 2003. Probability Theory: The Logic of Science. Cambridge University Press, Cambridge, 758 pp.

714 Kahle, H.P. et al. (Editors), 2008. Causes and Consequences of Forest Growth Trends in Europe. Brill, Leiden,  
715 xiv + 261 pp.

716 Kass, R.E. and Raftery, A.E., 1995. Bayes Factors. Journal of the American Statistical Association, 90(430):  
717 773-795.

718 Keats, A., Cheng, M.T., Yee, E. and Lien, F.S., 2009. Bayesian treatment of a chemical mass balance receptor  
719 model with multiplicative error structure. Atmospheric Environment, 43(3): 510-519.

720 Kesik, M. et al., 2005. Inventories of N<sub>2</sub>O and NO emissions from European forest soils. Biogeosciences, 2(4):  
721 353-375.

722 Klemetsson, L. et al., 2008. Bayesian calibration method used to elucidate carbon turnover in forest on drained  
723 organic soil. Biogeochemistry, 89(1): 61-79.

724 Kobayashi, K. and Salam, M.U., 2000. Comparing simulated and measured values using mean squared  
725 deviation and its components. Agronomy Journal, 92(2): 345-352.

726 Kreutzer, K., Butterbach-Bahl, K., Rennenberg, H. and Papen, H., 2009. The complete nitrogen cycle of an N-  
727 saturated spruce forest ecosystem. Plant Biology, 11(5): 643-649.

728 Kroon, P.S. et al., 2010. Uncertainties in eddy covariance flux measurements assessed from CH<sub>4</sub> and N<sub>2</sub>O  
729 observations. Agricultural and Forest Meteorology, 150(6): 806-816.

730 Li, C.S., Aber, J., Stange, F., Butterbach-Bahl, K. and Papen, H., 2000. A process-oriented model of N<sub>2</sub>O and  
731 NO emissions from forest soils: 1. Model development. Journal of Geophysical Research-  
732 Atmospheres, 105(D4): 4369-4384.

733 Li, C.S., Frohling, S. and Frohling, T.A., 1992. A model of nitrous-oxide evolution from soil driven by rainfall  
734 events.1. Model structure and sensitivity. *Journal of Geophysical Research-Atmospheres*, 97(D9):  
735 9759-9776.

736 Luo, Y. et al., 2009. Parameter identifiability, constraint, and equifinality in data assimilation with ecosystem  
737 models. *Ecological Applications*, 19(3): 571-574.

738 Luysaert, S. et al., 2010. The European carbon balance. Part 3: forests. *Global Change Biology*, 16(5): 1429-  
739 1450.

740 McCulloch, R.E. and Rossi, P.E., 1992. Bayes factors for nonlinear hypotheses and likelihood distributions.  
741 *Biometrika*, 79(4): 663-676.

742 Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H. and Teller, E., 1953. Equation of state  
743 calculations by fast computing machines. *Journal of Chemical Physics*, 21: 1087-1092.

744 Michalak, A.M. et al., 2005. Maximum likelihood estimation of covariance parameters for Bayesian  
745 atmospheric trace gas surface flux inversions. *Journal of Geophysical Research-Atmospheres*,  
746 110(D24): 16.

747 Moala, F.A. and O'Hagan, A., 2010. Elicitation of multivariate prior distributions: A nonparametric Bayesian  
748 approach. *Journal of Statistical Planning and Inference*, 140(7): 1635-1655.

749 Morris, M.D., 1991. Factorial sampling plans for preliminary computational experiments. *Technometrics*, 33(2):  
750 161-174.

751 Norman, J. et al., 2008. Simulation of NO and N<sub>2</sub>O emissions from a spruce forest during a freeze/thaw event  
752 using an N-flux submodel from the PnET-N-DNDC model integrated to CoupModel. *Ecological*  
753 *Modelling*, 216(1): 18-30.

754 Ogle, K., 2009. Hierarchical Bayesian statistics: merging experimental and modeling approaches in ecology.  
755 *Ecological Applications*, 19(3): 577-581.

756 Ogle, K. and Barber, J.J., 2008. Bayesian data-model integration in plant physiological and ecosystem ecology.  
757 *Progress in Botany*, 69: 281-311.

758 Papen, H. and Butterbach-Bahl, K., 1999. A 3-year continuous record of nitrogen trace gas fluxes from  
759 untreated and limed soil of a N-saturated spruce and beech forest ecosystem in Germany. 1. N<sub>2</sub>O  
760 emissions. *Journal of Geophysical Research-Atmospheres*, 104(D15): 18487-18503.

761 Parton, W.J. et al., 1993. Observations and modeling of biomass and soil organic-matter dynamics for the  
762 grassland biome worldwide. *Global Biogeochemical Cycles*, 7(4): 785-809.

- 763 Prinn, R.G., 2000. Measurement equation for trace chemicals in fluids and solution of its inverse. In: P.  
764 Kasibhatla et al. (Editors), *Inverse Methods in Global Biogeochemical Cycles*. Geophysical Mongraph.  
765 American Geophysical Union, Washington DC, pp. 3-18.
- 766 Raupach, M.R. et al., 2005. Model-data synthesis in terrestrial carbon observation: methods, data requirements  
767 and data uncertainty specifications. *Global Change Biology*, 11(3): 378-397.
- 768 Rothe, A., 1997. Einfluss des Baumartenanteils auf Durchwurzelung, Wasserhaushalt, Stoffhaushalt und  
769 Zuwachsleistung eines Fichten-Buchen-Mischbestandes am Standort Höglwald. *Forstliche*  
770 *Forschungsberichte*: 163 pp.
- 771 Rougier, J.C., (in press). Formal Bayes methods for model calibration with uncertainty. In: K. Beven and J. Hall  
772 (Editors), *Applied Uncertainty Analysis for Flood Risk Management*. Imperial College Press / World  
773 Scientific, London. Draft version available at <http://www.maths.bris.ac.uk/~mazjcr/FRMbox2-4.pdf>.
- 774 Saltelli, A., Chan, K. and Scott (Eds.), E.M. (Editors), 2000. *Sensitivity Analysis*. Wiley, Chichester, xv + 475.  
775 pp.
- 776 Sivia, D.S., 2006. *Data Analysis: A Bayesian Tutorial*. Second edition. Oxford University Press, Oxford, 260  
777 pp.
- 778 Sutton, M.A. et al., 2007. Challenges in quantifying biosphere-atmosphere exchange of nitrogen species.  
779 *Environmental Pollution*, 150: 125-139.
- 780 Sutton, M.A. et al., 2008. Uncertainties in the relationship between atmospheric nitrogen deposition and forest  
781 carbon sequestration. *Global Change Biology*, 14(9): 2057-2063.
- 782 Svensson, M. et al., 2008. Bayesian calibration of a model describing carbon, water and heat fluxes for a  
783 Swedish boreal forest stand. *Ecological Modelling*, 213(3-4): 331-344.
- 784 Tuomi, M., Vanhala, P., Karhu, K., Fritze, H. and Liski, J., 2008. Heterotrophic soil respiration--Comparison of  
785 different models describing its temperature dependence. *Ecological Modelling*, 211(1-2): 182-190.
- 786 Van Oijen, M. et al., 2008a. Evaluation of past and future changes in European forest growth by means of four  
787 process-based models. In: H.P. Kahle et al. (Editors), *Causes and Consequences of Forest Growth*  
788 *Trends in Europe*. Brill, Leiden, pp. 183-199.
- 789 Van Oijen, M. et al., 2008b. Methodology for the application of process-based models to analyse changes in  
790 European forest growth. In: H.P. Kahle et al. (Editors), *Causes and Consequences of Forest Growth*  
791 *Trends in Europe*. Brill, Leiden, pp. 67-80.
- 792 Van Oijen, M., Cannell, M.G.R. and Levy, P.E., 2004. Modelling biogeochemical cycles in forests: state of the

793 art and perspectives. In: F. Andersson, Y. Birot and R. Päivinen (Editors), Towards the Sustainable Use  
794 of Europe's Forests - Forest Ecosystem and Landscape Research: Scientific Challenges and  
795 Opportunities. EFI Proceedings. European Forest Institute, pp. 157-169.

796 Van Oijen, M., Dauzat, J., Harmand, J.-M., Lawson, G. and Vaast, P., 2010a. Coffee agroforestry systems in  
797 Central America: II. Development of a simple process-based model and preliminary results.  
798 *Agroforestry Systems*, 80(3): 361-378.

799 Van Oijen, M., Rougier, J. and Smith, R., 2005. Bayesian calibration of process-based forest models: bridging  
800 the gap between models and data. *Tree Physiology*, 25(7): 915-927.

801 Van Oijen, M., Schapendonk, A. and Höglind, M., 2010b. On the relative magnitudes of photosynthesis,  
802 respiration, growth and carbon storage in vegetation. *Annals of Botany*, 105(5): 793-797.

803 Van Oijen, M. and Thomson, A., 2010. Toward Bayesian uncertainty quantification for forestry models used in  
804 the United Kingdom Greenhouse Gas Inventory for land use, land use change, and forestry. *Climatic  
805 Change*, 103(1): 55-67.

806 Wang, Y.-P., Trudinger, C.M. and Enting, I.G., 2009. A review of applications of model-data fusion to studies  
807 of terrestrial carbon fluxes at different scales. *Agricultural and Forest Meteorology*, 149(11): 1829-  
808 1842.

809 Werner, C., Butterbach-Bahl, K., Haas, E., Hickler, T. and Kiese, R., 2007. A global inventory of N<sub>2</sub>O  
810 emissions from tropical rainforest soils using a detailed biogeochemical model. *Global Biogeochem.  
811 Cycles*, 21(3): GB3010.

812 Wu, X. et al., 2010. Environmental controls over soil-atmosphere exchange of N<sub>2</sub>O, NO, and CO<sub>2</sub> in a temperate  
813 Norway spruce forest. *Global Biogeochem. Cycles*, 24(2): GB2012.

814 Yeluripati, J.B. et al., 2009. Bayesian calibration as a tool for initialising the carbon pools of dynamic soil  
815 models. *Soil Biology & Biochemistry*, 41(12): 2579-2583.

816

817

818

819

Table 1. Overview of data used for calibration. Variables are annual quantiles of daily emission rates where indicated, and monthly averages otherwise. The period of measurement indicates the years of first and last measurement, and  $n$  is the total number of data points over that period. The columns marked *min*, *mean* and *max* show the extreme values and the mean of the  $n$  values.

Variable	Unit	Period of measurement	$n$	<i>min</i>	<i>mean</i>	<i>max</i>
N <sub>2</sub> O (Q10)	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2003	6	0.08	0.21	0.31
N <sub>2</sub> O (Q50)	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2003	6	0.27	0.52	0.85
N <sub>2</sub> O (Q90)	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2003	6	0.55	2.32	9.14
NO (Q10)	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2002	5	1.65	2.33	2.79
NO (Q50)	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2002	5	4.81	6.64	8.30
NO (Q90)	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2002	5	10.65	13.65	18.52
N <sub>2</sub> O	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2003	70	0.03	0.99	16.55
NO	kg N ha <sup>-1</sup> y <sup>-1</sup>	1994-2003	61	1.74	7.72	24.04
Rsoil	mg C m <sup>-2</sup> h <sup>-1</sup>	1995-1997	36	23.9	113.3	255.2
Water	% (vol)	1994-1996	25	27.1	33.6	37.0

820

821

822

823

824

825

826

827

828

829

Table 2. Overview of the models.				
Property	BASFOR	COUP	DAYCENT	MoBiLE-DNDC
# State variables (trees, soil)	14 (6,8)	56 (8,48)	20 (10,10)	38 (10,28)
# Parameters (in calibration)	48 (48)	>300 (23)	>300 (17)	67* (26)
Time step	Daily	Hourly	Daily	Daily (but less than 1 minute for diffusion processes)
Inputs: Environmental time series	Radiation, temperature, precipitation, humidity, wind speed, N-deposition, [CO <sub>2</sub> ]	Radiation, temperature, precipitation, humidity, wind speed, N-deposition	Radiation, temperature, precipitation, humidity, wind speed, N-deposition, [CO <sub>2</sub> ]	Radiation, temperature, precipitation, N- deposition, [CO <sub>2</sub> ]
Inputs: environmental constants	Soil water retention curve, rooting depth	Soil water retention curve, rooting depth	Soil water retention curve, rooting depth	Layer-specific texture, bulk density, field capacity, pH

831 \* For the soil chemistry module only.

832

Table 3. Summary of model output, with prior and posterior uncertainties. All values shown refer to time series averages, to be compared with the last-but-one column of Table 1. For each model, the table entries are three quantiles (5, 50 and 95%) of the output distributions generated by the prior and posterior parameter distributions. In **bold**: posterior distributions for which the posterior width (Q5-Q95) is at least an order of magnitude less than for the prior.

Var.	Unit	Dist.	BASFOR			COUP			DAYCENT			MoBiLE-DNDC		
			Q5	Q50	Q95	Q5	Q50	Q95	Q5	Q50	Q95	Q5	Q50	Q95
N <sub>2</sub> O (Q10)	kg N	Prior	0.01	0.05	1.78	0.01	0.03	0.13	0.01	0.02	0.02	0.09	0.68	2.82
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	<b>0.02</b>	<b>0.03</b>	<b>0.04</b>	0.03	0.05	0.07	0.01	0.01	0.02	<b>0.13</b>	<b>0.16</b>	<b>0.21</b>
N <sub>2</sub> O (Q50)	kg N	Prior	0.03	0.18	3.83	0.02	0.10	0.49	0.34	0.41	0.50	0.19	1.64	6.82
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	<b>0.07</b>	<b>0.12</b>	<b>0.17</b>	0.06	0.10	0.16	0.29	0.33	0.39	<b>0.25</b>	<b>0.30</b>	<b>0.34</b>
N <sub>2</sub> O (Q90)	kg N	Prior	0.12	1.13	8.08	0.16	0.69	4.49	2.00	2.13	2.28	0.40	3.21	12.85
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	<b>0.29</b>	<b>0.46</b>	<b>0.65</b>	<b>0.11</b>	<b>0.18</b>	<b>0.33</b>	1.90	1.99	2.08	<b>0.41</b>	<b>0.51</b>	<b>0.62</b>
NO (Q10)	kg N	Prior	0.03	0.10	8.99	3.42	7.24	14.26	0.72	0.75	0.78	0.27	2.88	12.22
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	<b>0.05</b>	<b>0.11</b>	<b>0.18</b>	1.60	2.10	2.77	0.71	0.72	0.75	<b>1.84</b>	<b>2.28</b>	<b>2.72</b>
NO (Q50)	kg N	Prior	0.09	0.45	12.21	5.95	13.02	28.90	1.16	1.26	1.37	0.92	5.83	17.91
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	<b>0.25</b>	<b>0.79</b>	<b>1.43</b>	<b>2.41</b>	<b>3.20</b>	<b>4.64</b>	1.08	1.15	1.23	<b>3.52</b>	<b>4.19</b>	<b>4.86</b>
NO (Q90)	kg N	Prior	0.31	2.38	19.99	10.56	22.88	48.30	0.51	0.57	0.63	1.89	9.39	23.70
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	0.99	2.69	4.57	4.41	6.22	8.70	0.48	0.51	0.55	5.61	6.65	8.02
N <sub>2</sub> O	kg N	Prior	0.05	0.39	4.42	0.07	0.27	1.44	1.55	1.82	2.11	0.24	1.81	7.42
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	<b>0.11</b>	<b>0.18</b>	<b>0.24</b>	<b>0.08</b>	<b>0.12</b>	<b>0.20</b>	1.40	1.55	1.74	<b>0.27</b>	<b>0.32</b>	<b>0.37</b>
NO	kg N	Prior	0.15	0.92	13.25	6.05	14.78	32.82	4.04	4.97	6.00	1.21	6.00	17.58
	ha <sup>-1</sup> y <sup>-1</sup>	Post.	0.43	1.16	1.99	2.96	4.11	5.68	3.54	4.05	4.74	<b>3.68</b>	<b>4.32</b>	<b>5.01</b>
CO <sub>2</sub>	mg C	Prior	44.2	69.4	106.8	76.5	97.0	127.0	25.3	48.2	73.0	46.4	66.3	96.5
	m <sup>2</sup> h <sup>-1</sup>	Post.	56.2	76.5	99.5	67.6	84.4	101.2	26.9	49.6	76.2	49.1	55.4	64.2
Water	%	Prior	28.3	31.6	33.7	33.2	34.5	36.1	32.3	32.3	32.3	34.5	34.5	34.5
	(vol)	Post.	28.7	31.0	33.0	33.3	34.6	35.9	32.3	32.3	32.3	34.5	34.5	34.5

Table 4. Comparison of data with model outputs: correlation coefficient ( $r$ ) and normalised root mean square error (NRMSE). The table shows quantiles (Q5, Q50, Q95) of the distributions of  $r$  and NRMSE induced by prior and posterior parameter distributions. In **bold**: posterior values that are improvements over the prior ( $r$  increased, NRMSE reduced).

Var.	Dist.	Statistic	BASFOR			COUP			DAYCENT			MoBiLE-DNDC		
			Q5	Q50	Q95	Q5	Q50	Q95	Q5	Q50	Q95	Q5	Q50	Q95
N <sub>2</sub> O (Q10)	Prior	$r$	-0.37	0.42	0.72	-0.83	-0.09	0.47	-0.36	-0.27	-0.19	-0.27	0.54	0.76
		NRMSE	0.44	1.07	11.95	0.41	0.93	1.91	0.52	1.01	1.95	0.42	2.47	12.95
	Post.	$r$	-0.75	-0.19	0.49	-0.60	-0.22	0.20	-0.39	-0.29	-0.21	<b>0.21</b>	<b>0.70</b>	<b>0.87</b>
		NRMSE	<b>0.41</b>	<b>0.47</b>	<b>0.53</b>	<b>0.31</b>	<b>0.38</b>	<b>0.48</b>	0.54	<b>0.77</b>	<b>1.05</b>	<b>0.20</b>	<b>0.39</b>	<b>1.32</b>
N <sub>2</sub> O (Q50)	Prior	$r$	-0.46	0.65	0.88	-0.62	0.28	0.77	0.45	0.60	0.70	-0.67	-0.14	0.33
		NRMSE	0.40	1.02	7.64	0.35	0.87	1.86	0.23	0.40	1.25	0.48	2.46	13.32
	Post.	$r$	-0.53	-0.04	0.78	-0.73	-0.37	0.38	0.45	<b>0.62</b>	<b>0.71</b>	<b>-0.06</b>	<b>0.15</b>	<b>0.50</b>
		NRMSE	<b>0.31</b>	<b>0.38</b>	<b>0.46</b>	<b>0.32</b>	<b>0.40</b>	<b>0.48</b>	<b>0.19</b>	<b>0.28</b>	<b>0.48</b>	<b>0.22</b>	<b>0.47</b>	<b>1.51</b>
N <sub>2</sub> O (Q90)	Prior	$r$	-0.51	0.32	0.80	-0.57	0.22	0.96	0.60	0.60	0.61	-0.58	-0.41	-0.08
		NRMSE	0.83	1.64	3.13	0.64	1.50	3.01	0.72	1.06	2.34	0.94	2.04	5.43
	Post.	$r$	<b>-0.35</b>	<b>0.55</b>	<b>0.86</b>	-0.68	-0.33	0.70	0.60	0.60	0.61	<b>-0.38</b>	<b>-0.25</b>	-0.13
		NRMSE	<b>0.67</b>	<b>0.75</b>	<b>0.86</b>	0.77	<b>0.82</b>	<b>0.91</b>	<b>0.71</b>	<b>0.82</b>	<b>1.13</b>	<b>0.78</b>	<b>1.34</b>	<b>3.05</b>
NO (Q10)	Prior	$r$	0.20	0.52	0.78	-0.11	0.20	0.49	-0.44	-0.44	-0.43	-0.73	0.37	0.55
		NRMSE	0.53	1.13	4.54	0.72	2.26	5.37	0.37	0.75	1.62	0.30	0.97	4.37
	Post.	$r$	<b>0.39</b>	<b>0.60</b>	<b>0.84</b>	<b>0.05</b>	<b>0.41</b>	<b>0.67</b>	-0.44	-0.44	-0.44	<b>-0.62</b>	-0.01	<b>0.66</b>
		NRMSE	<b>0.45</b>	<b>0.48</b>	<b>0.54</b>	<b>0.18</b>	<b>0.35</b>	<b>0.55</b>	<b>0.35</b>	<b>0.40</b>	<b>0.51</b>	<b>0.18</b>	<b>0.37</b>	<b>1.00</b>
NO (Q50)	Prior	$r$	-0.60	-0.32	0.37	-0.90	-0.68	-0.36	0.14	0.18	0.22	-0.92	-0.21	0.39
		NRMSE	0.50	1.04	1.96	0.53	1.26	3.58	0.38	0.82	1.71	0.24	0.67	1.84

	Post.	$r$	<b>-0.50</b>	-0.38	-0.25	<b>-0.78</b>	<b>-0.51</b>	<b>-0.20</b>	0.13	0.16	0.20	<b>-0.80</b>	-0.36	0.03
		NRMSE	<b>0.38</b>	<b>0.44</b>	<b>0.51</b>	<b>0.19</b>	<b>0.27</b>	<b>0.40</b>	<b>0.36</b>	<b>0.42</b>	<b>0.56</b>	<b>0.16</b>	<b>0.39</b>	<b>1.31</b>
NO (Q90)	Prior	$r$	-0.59	-0.34	0.11	-0.93	-0.70	-0.36	-0.59	-0.53	-0.43	-0.89	-0.30	0.27
		NRMSE	0.41	0.90	1.78	0.42	1.00	2.79	0.51	0.96	1.86	0.20	0.56	1.44
	Post.	$r$	<b>-0.51</b>	-0.35	-0.09	<b>-0.90</b>	<b>-0.53</b>	<b>-0.09</b>	-0.61	-0.58	-0.53	<b>-0.86</b>	-0.52	0.01
		NRMSE	<b>0.28</b>	<b>0.37</b>	<b>0.48</b>	<b>0.14</b>	<b>0.25</b>	<b>0.42</b>	<b>0.48</b>	<b>0.54</b>	<b>0.69</b>	<b>0.15</b>	<b>0.49</b>	1.44
N <sub>2</sub> O	Prior	$r$	-0.15	0.07	0.41	-0.14	0.25	0.66	-0.08	-0.08	-0.07	-0.24	-0.19	-0.06
		NRMSE	1.34	2.64	5.00	1.16	2.32	4.46	3.28	4.10	5.29	1.56	3.40	8.58
	Post.	$r$	<b>-0.13</b>	-0.04	0.21	-0.26	-0.05	0.36	-0.08	<b>-0.07</b>	-0.07	<b>-0.19</b>	<b>-0.12</b>	<b>-0.02</b>
		NRMSE	<b>1.14</b>	<b>1.20</b>	<b>1.35</b>	<b>1.14</b>	<b>1.20</b>	<b>1.33</b>	<b>2.79</b>	<b>3.18</b>	<b>3.74</b>	<b>1.19</b>	<b>2.00</b>	<b>4.47</b>
NO	Prior	$r$	-0.59	-0.32	0.26	0.42	0.55	0.65	0.69	0.69	0.70	0.09	0.52	0.67
		NRMSE	0.62	1.24	2.21	0.50	1.23	3.74	0.39	0.58	1.48	0.37	0.86	1.83
	Post.	$r$	<b>-0.43</b>	<b>-0.20</b>	0.05	0.42	0.50	0.58	0.69	0.69	0.70	<b>0.52</b>	<b>0.67</b>	<b>0.74</b>
		NRMSE	<b>0.50</b>	<b>0.54</b>	<b>0.60</b>	<b>0.28</b>	<b>0.35</b>	<b>0.51</b>	<b>0.31</b>	<b>0.36</b>	<b>0.42</b>	<b>0.26</b>	<b>0.61</b>	<b>1.66</b>
CO <sub>2</sub>	Prior	$r$	0.76	0.80	0.82	0.85	0.87	0.89	0.87	0.89	0.90	0.85	0.89	0.91
		NRMSE	0.27	0.64	1.49	0.25	0.44	1.25	0.24	0.63	1.67	0.17	0.55	1.53
	Post.	$r$	0.76	0.79	0.81	<b>0.86</b>	<b>0.88</b>	0.89	0.87	0.89	0.90	<b>0.86</b>	0.88	0.89
		NRMSE	0.28	<b>0.39</b>	<b>0.64</b>	<b>0.20</b>	<b>0.28</b>	<b>0.53</b>	<b>0.22</b>	<b>0.34</b>	<b>0.76</b>	0.17	<b>0.42</b>	<b>1.36</b>
Water	Prior	$r$	-0.16	0.03	0.55	0.36	0.38	0.39	0.53	0.53	0.53	0.55	0.56	0.56
		NRMSE	0.04	0.10	0.26	0.08	0.13	0.26	0.11	0.15	0.28	0.02	0.10	0.25
	Post.	$r$	-0.17	-0.03	0.27	0.36	0.38	0.39	0.53	0.53	0.53	0.55	0.55	0.56
		NRMSE	0.04	<b>0.07</b>	<b>0.13</b>	0.08	<b>0.10</b>	<b>0.17</b>	0.11	<b>0.14</b>	<b>0.21</b>	0.02	<b>0.09</b>	<b>0.24</b>

838

839

840

841

842           **FIGURES**

843

844 Fig. 1. Dots: Monthly averages of measured emissions of N<sub>2</sub>O, NO, CO<sub>2</sub> and of soil water  
845 content. The lines represent output of model BASFOR. Dashed red line: output for the  
846 mode of the prior parameter distribution. Thick black line: output for the mode of the  
847 posterior. Thin black lines: 5% and 95% quantiles of the posterior output distribution.

848

849 Fig. 2. Posterior marginal distributions for the four data-scaling factors that represent  
850 systematic multiplicative error in the data according to the different models.

851

852 Fig. 3. Differences between measurements and simulations (positive values indicating  
853 underestimates by the models) for time series with monthly averages. Top 4 panels:  
854 simulations using the mode of the prior parameter distribution. Bottom 4 panels:  
855 simulations with the posterior mode.

856

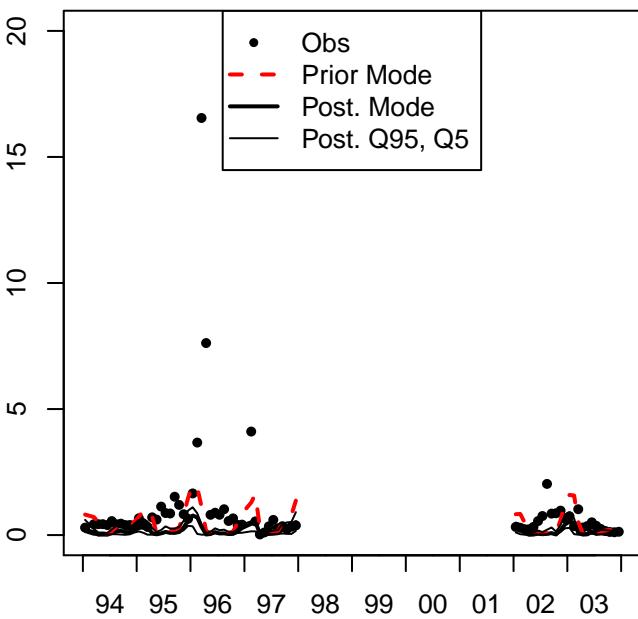
857 Fig. 4. Distributions of log-likelihoods for each of the four models, for the four categories of  
858 output variables. Grey: prior, black: posterior.

859

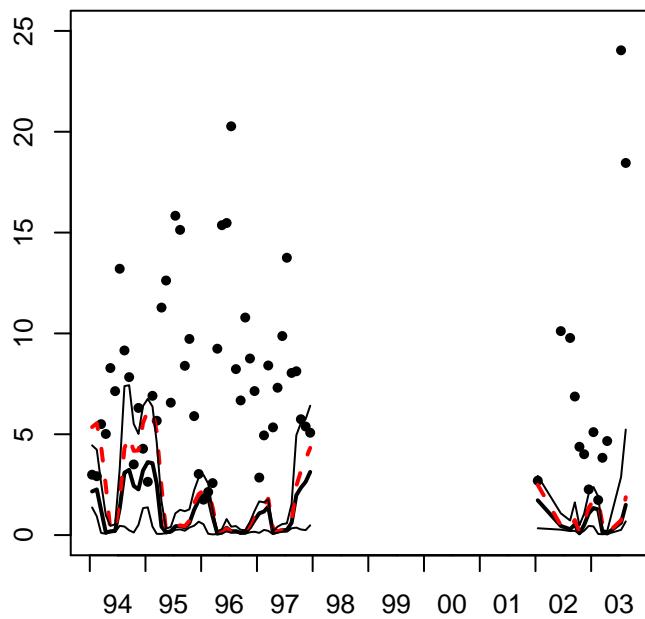
860 Fig. 5. Decomposition of the Mean Squared Error (MSE) associated with the modes of the  
861 prior and posterior parameter distributions, for the time series with monthly data. The  
862 MSE-values for the N<sub>2</sub>O and NO are in the same units (kg N ha<sup>-1</sup> y<sup>-1</sup> squared), the

863 MSE-values for soil respiration and soil water content are in squared  $\text{mg C m}^{-2} \text{ h}^{-1}$  and  
864 squared %, respectively.

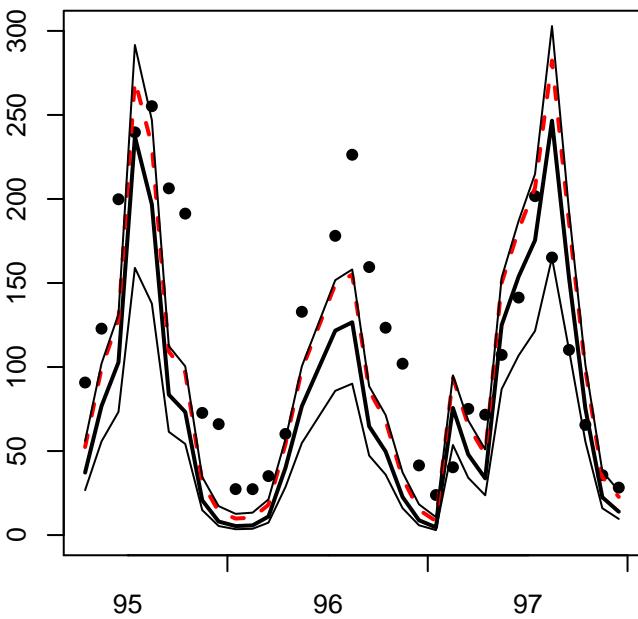
**Figure 1**  $\text{N}_2\text{O}$  emission ( $\text{kg N ha}^{-1}\text{yr}^{-1}$ )



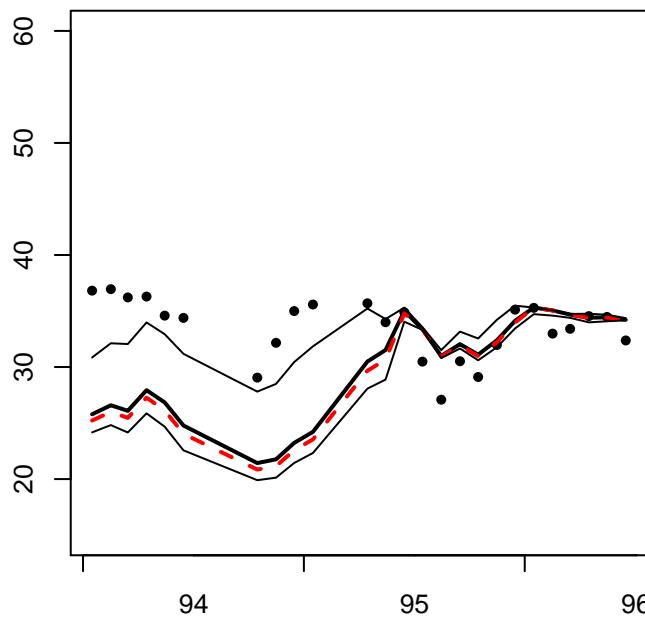
$\text{NO}$  emission ( $\text{kg N ha}^{-1}\text{yr}^{-1}$ )

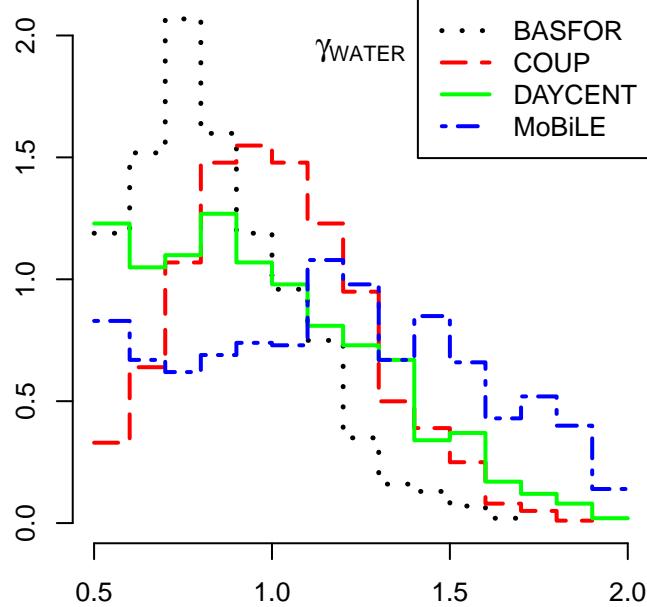
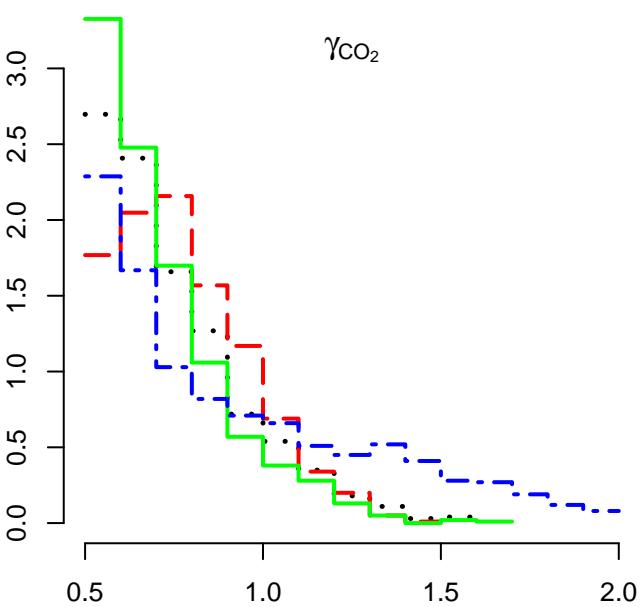
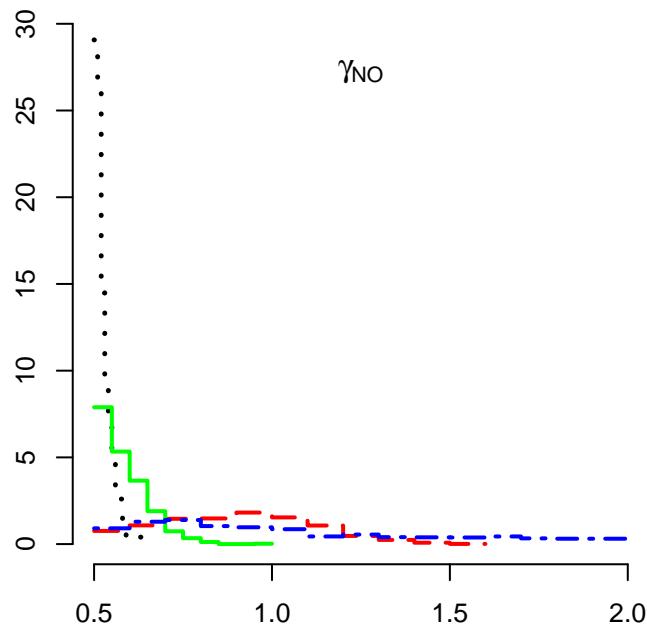
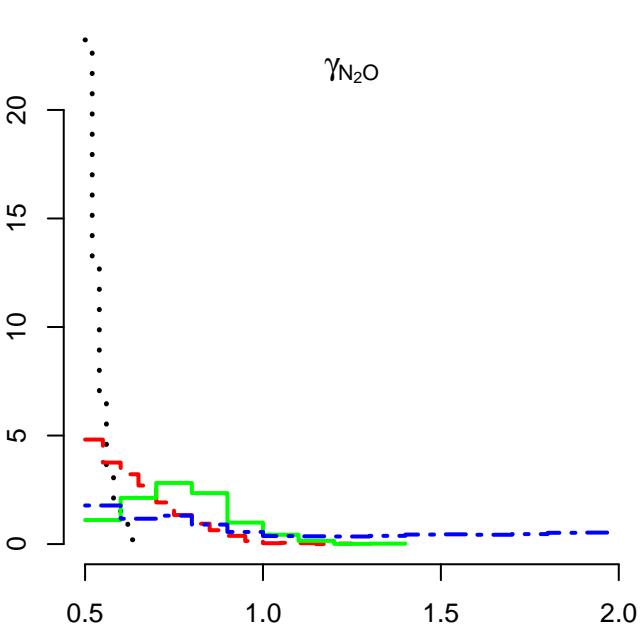


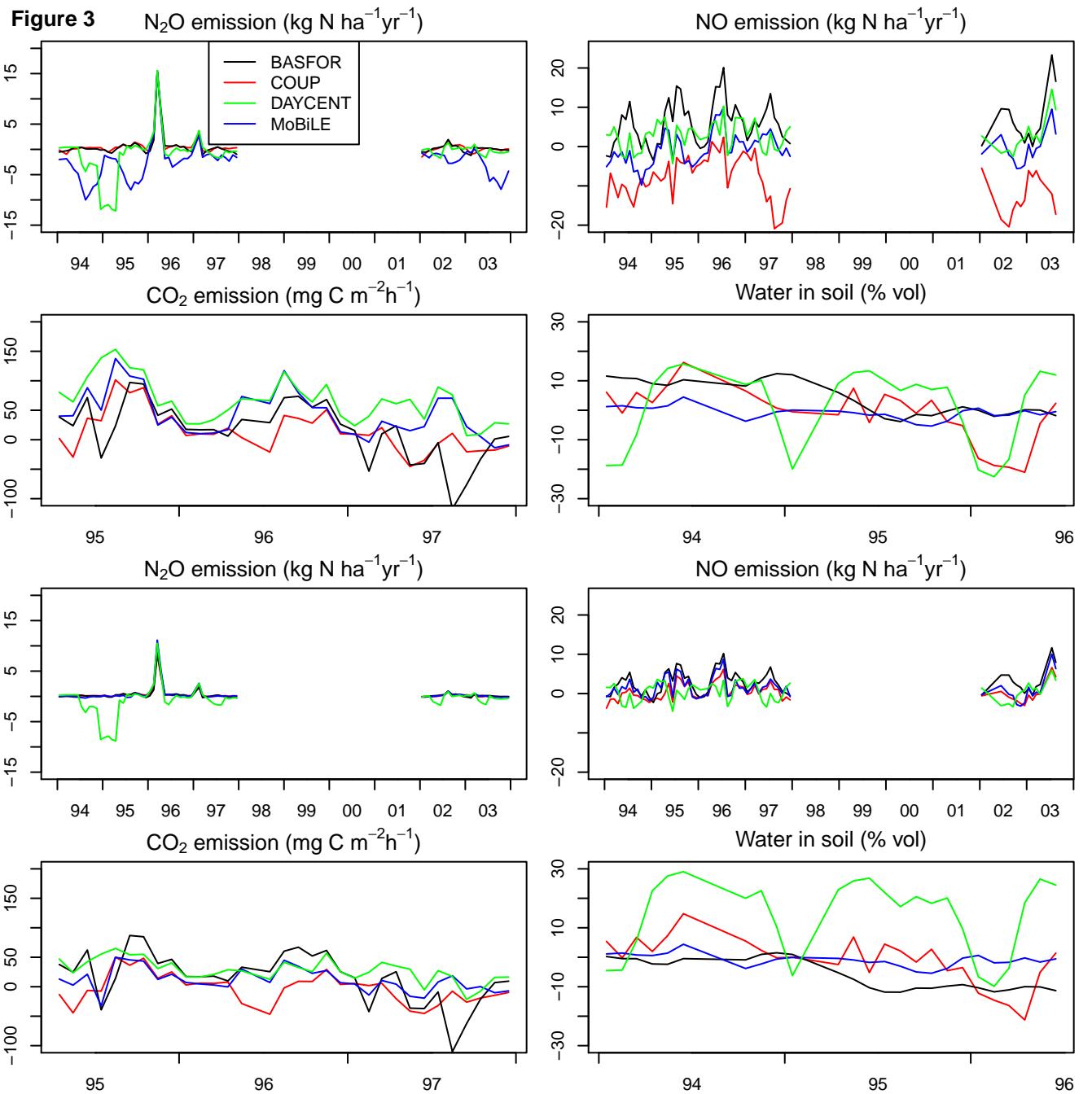
$\text{CO}_2$  emission ( $\text{mg C m}^{-2}\text{h}^{-1}$ )



Water in soil (% vol)



**Figure 2**



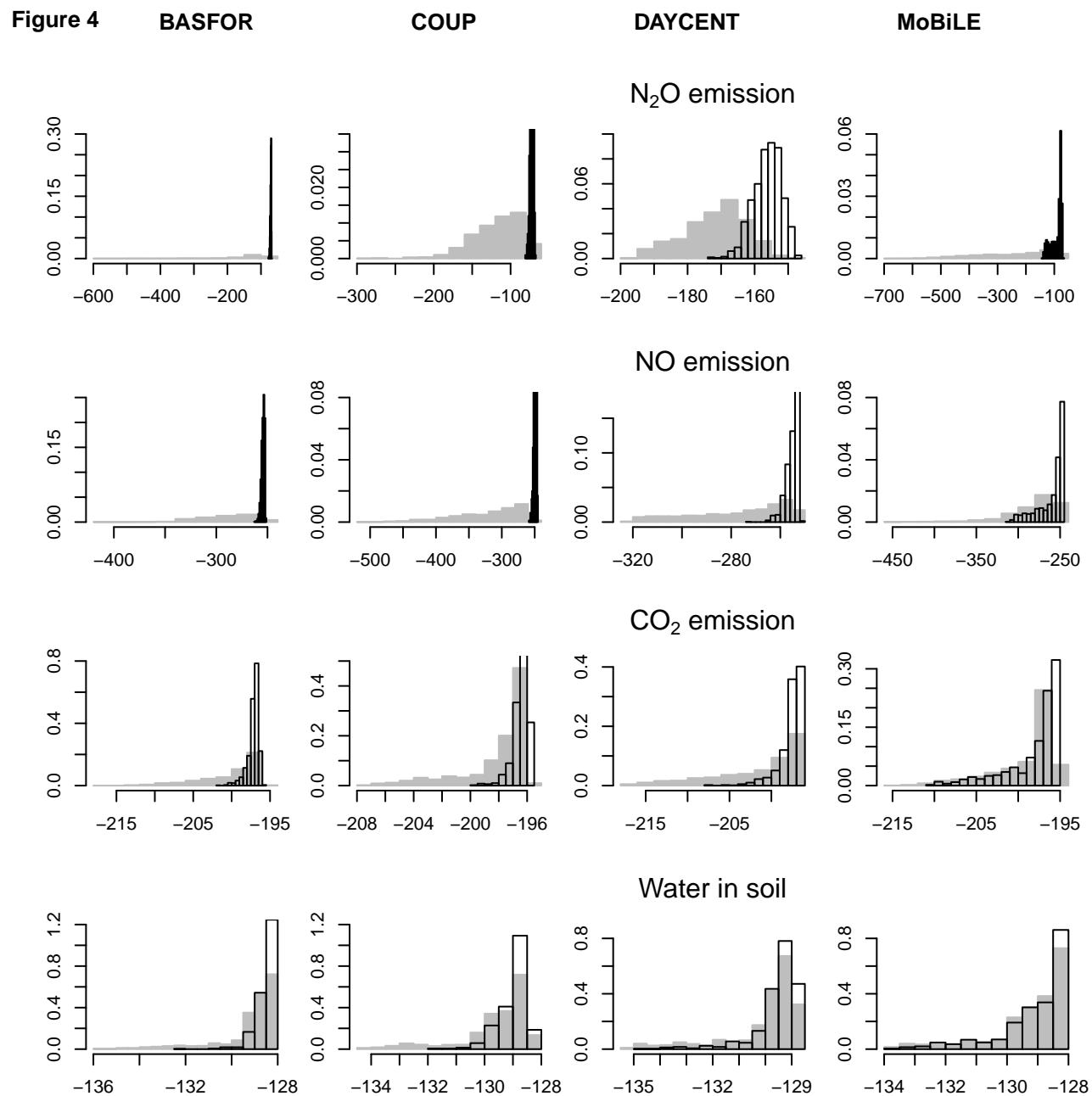


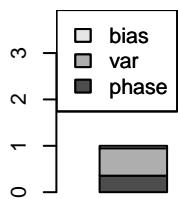
Figure 5

BASFOR

COUP

DAYCENT

MoBiLE

N<sub>2</sub>O emission

NO emission

CO<sub>2</sub> emission

Water in soil

