

1 **Referee 1**

2 We thank the referee for the comments on our manuscript, which helped improving our study.  
3 We hope that our answers and the modifications are satisfactory.

4 **Page 2, Line 9: I don't understand "substitute observed by modeled fluxes". Substitute**  
5 **modeled fluxes for observed fluxes?**

6 We reformulated: "In addition we create synthetic observations using modeled fluxes instead of  
7 the observed ones, to explore the potential to infer prior uncertainties from model-model  
8 residuals".

9 **Page 2, Line 11: Was the random error added to observed or modeled "tower" fluxes? If**  
10 **this was added to the observed fluxes, why? There is already random error in the**  
11 **measurements.**

12 We clarified: " a random measurement noise was added to the modeled tower fluxes".

13 **Page 2, Lines 14-15: "This difference... " isn't clear. Do the large biases exist with respect**  
14 **to 5PM, or the other models? And how does a large bias cause a long temporal**  
15 **autocorrelation?**

16 We reformulated: "This difference is caused by a few sites with large biases between the data  
17 and the 5PM model."

18 Regarding the second question: we also computed the temporal autocorrelation time excluding  
19 those sites with a large model-data mismatch bias (see section 5.1, page 9408), and this was  
20 found to be less than half of the temporal autocorrelation time using all sites.

21 Note that unfortunately the abstract contained the wrong numbers which were inconsistent with  
22 those in table 2 and in the text. We corrected this, and page 9394 lines 13,14 now reads: 30 and  
23 70 days.

24 **Page 3, line 10. I don't understand the term "regularized." Can this be defined? It is used**  
25 **in more than one place, and I don't recognize what concept is being communicated.**

26 "Regularization" is a standard term in statistics and refers to a process of introducing additional  
27 constraints in order to solve an ill-posed problem. See for instance:

28 Hansen PC, Oleary DP, The use of the L-curve in the regularization of discrete ill-posed  
29 problems. SIAM journal of scientific computing, 14 (6) 1487-1503, 1993.

30 We added: "In this way, the solution of the otherwise ill-posed problem is regularized in the  
31 sense that the optimization problem becomes one with a unique solution."

1 **Page 4, lines 4-6. Coarser scale inversions may not explicitly utilize correlation lengths, but**  
2 **they are implicitly imposing a large correlation length (which may be entirely**  
3 **inappropriate). Please clarify.**

4 We clarified: "...(Houweling et al., 2004; Rödenbeck et al., 2003b). For the former case large  
5 correlation scales are implicitly assumed since fluxes within a grid-cell are fully correlated. For  
6 regional..."

7 **Page 4, lines 13-14. I don't understand this description. How can you derive a spatial**  
8 **correlation in the prior flux error from a coarse resolution inversion?**

9 We reformulated: "In some regional studies the same correlations are used as in large scale  
10 inversions..."

11 **Page 4, lines 22-23. This is not correct. The pattern of fluxes was not used to evaluate the**  
12 **spatial correlation length. Lauvaux et al, 2012 tested the spatial correlation length scale by**  
13 **cross-validation of the posterior CO2 mole fractions. CO2 observations were reserved from**  
14 **the inversion, and the correlation length that provided the best fit to the reserved CO2 data**  
15 **was identified as the best choice.**

16 We clarified: "...based on cross-validation of the simulated against observed CO2 mole  
17 fractions. The simulated mole fractions were derived using the influence functions and the  
18 inverted fluxes"

19 **Page 5, lines 1-4. This sentence is unintelligible.**

20 We clarified: "A recent study by Broquet et al. (2013) obtained good agreements between the  
21 statistical uncertainties as derived from the inversion system and the actual misfits calculated by  
22 comparing the posterior fluxes to local flux measurements at the European and 1-month scale"

23 **Page 4, lines 12-13. Make it clear that flux measurement sensitivities are areas, not lengths.**  
24 **You are describing dimensions of an area measurement. This is not clear as written. I also**  
25 **recommend that you note that the resolution of an inversion system is not necessarily the**  
26 **same as the true resolution of an inverse flux estimate.**

27 We clarified: "While typical inversion systems have a resolution ranging from tens of kilometers  
28 up to several degrees (hundreds of km), with the true resolution of the inverse flux estimates  
29 being even coarser, the spatial representativity of the flux observations typically covers an area  
30 with a radius of around a kilometer"

31

32 **Page 4, line 15. There are many, many studies of flux upscaling with towers and**  
33 **spatial databases. The paper would benefit from a somewhat expanded review of this**  
34 **literature.**

1 We added: "...one applied by Chevallier et al., (2012). Typical approaches to up-scale site level  
2 fluxes deploy for example model tree algorithms, a machine learning algorithm which is trained  
3 to predict carbon flux estimates based on meteorological data, vegetation properties and types  
4 (Jung et al., 2009, Xiao et al., 2008), or neural networks (Papale and Valentini 2003).  
5 Nevertheless eddy covariance measurements provide..."

6 References added:

7 Jung, M., Reichstein, M., and Bondeau, A.: Towards global empirical upscaling of FLUXNET  
8 eddy covariance observations: validation of a model tree ensemble approach using a biosphere  
9 model, *Biogeosciences*, 6, 2001-2013, doi:10.5194/bg-6-2001-2009, 2009.

10

11 Papale, D. and Valentini, R., A new assessment of European forests carbon exchanges by eddy  
12 fluxes and artificial neural network spatialization. *Global Change Biology*, 9: 525–535. doi:  
13 10.1046/j.1365-2486.2003.00609.x, 2003

14

15 Xiao, J. F., Zhuang, Q. L., Baldocchi, D. D., et al.: Estimation of net ecosystem carbon exchange  
16 for the conterminous United States by combining MODIS and AmeriFlux data, *Agr. Forest  
17 Meteorol.*, 148(11), 1827–1847, 2008.

18

19 **Page 7, line 25. What does “across entire daily course” mean?**

20 We mean flights were made at various times of the day to cover the daily course on average.

21 We clarified: "... flown 52 and 54 times respectively, covering the daily course. Exact routes..."

22 **Page 9, lines 13-14. Why is the year of EC data used to optimize respiration parameters  
23 singled out? What about the other parameters and years of EC data?**

24 For the heterotrophic respiration no calibration with many years was ever done here; so we  
25 followed an ad hoc procedure and chose the year 2007 to calibrate, which is the year for which  
26 the calculations presented in the paper have been done, to prevent any bias caused by a  
27 systematic error in the respiration. The calibration pertains by the way only to the “amplitude” of  
28 the respiration.

29

1 For other parameters, the calibrations come from Groenendijk et al., (2011) and are based on 4.5  
2 years per site on average.

3 We modified : “The “5 parameter model” (5PM) (Groenendijk et al., 2011), also used in  
4 atmospheric inversions (Tolk et al., 2011, Meesters et al., 2012), is a physiological model  
5 describing transpiration, photosynthesis, and respiration. It uses MODIS LAI (leaf area index) at  
6 10km resolution, meteorological data (temperature, moisture, and downward shortwave radiative  
7 flux, presently from ECMWF at 0.25 degrees resolution), and differentiates PFTs for different  
8 vegetation types and climate regions. 5PM fluxes are pertained to station locations at hourly  
9 temporal resolution. The optimization has been done with EC-data from Fluxnet as described  
10 (except for heterotrophic respiration) in Groenendijk et al., 2011. Regarding the heterotrophic  
11 respiration, an ad hoc optimization using Fluxnet EC-data from 2007 was performed since no  
12 previous optimization was available.”

13 **Page 10, line 3. English. I suggest, “it is the only model with spatial resolution (10 km)**  
14 **comparable to ...”**

15 We corrected: “VPRM was used since it is the only model with spatial resolution (10 km)  
16 comparable to ...”.

17 **Page 10, lines 6-8. Number of flights is repeated. Resolution of aircraft flight data**  
18 **contradicts earlier text that stated 2 km resolution.**

19 Repetition is deleted: “Each grid point was sampled 52 times in forests, and 54 in agricultural  
20 land”.

21 We clarified: “... short flight distances. Aircraft NEE data, natively at 2 km resolution along the  
22 track, have been aggregated into 10 km segments, to maximize the overlap with the VPRM  
23 grid...”.

24 **Page 10, lines 15-16. Do you neglect the impact of the observation errors? Or is it just that**  
25 **you cannot separate these errors from errors in the model? The observation errors are**  
26 **already part of the observations.**

27 As the observation errors are already part of the observations we do consider them in the model-  
28 observation analysis. The word “neglect” here only means that the impact on the correlation  
29 lengths is not studied here. The impact is studied later in the model-model analysis where we  
30 consider or neglect the observation error by adding or not a random measurement error to the  
31 fluxes used as reference.

32 We furthermore clarified: “therefore e-folding correlation length estimations do include the  
33 observation error term.”.

1 **Page 11, equations (2) and (3) and the “nugget” effect: Mathematically, when  $k=0$ ,**  
2 **equation (2) = 1. What is this nugget effect? How can it exist for  $k=0$ ? The numerator and**  
3 **denominator are identical when  $k=0$ .**

4 “Nugget” is a standard term in spatial statistics. The nugget effect can be attributed to  
5 measurement errors or spatial sources of variation at distances smaller than the sampling interval  
6 or both. See also Wackernagel, H., 2003, Multivariate Geostatistics, Springer.

7 Page 9403 line 15. We clarified: “... also known as the nugget effect. The nugget effect is driven  
8 by measurement errors and variations at distances (spatial or temporal) smaller than the sampling  
9 interval. For this...”.

10 **Page 11, line 19. What does “distributed along the entire daily course” mean? And how can**  
11 **aircraft data span 36 days? Are the times in between the aircraft flights neglected?**

12 We removed that unclear wording: “distributed along the entire daily course”.

13 36 days is the duration of the flight campaign. Within this period, flights were planned and made  
14 to equally sample each time of the day, along the campaign. Since flights are obviously  
15 discontinuous in time, they constitute a sub-sampling thus are insufficient to compute a daily  
16 average at each grid-cell. Our approach was therefore to keep all individual fluxes to constrain  
17 the temporal autocorrelation, assuming that fluxes are overall not biased toward specific times of  
18 the day, given the sub-sampling is equally distributed.

19 And yes, temporal gaps that are not sampled are not used in the analysis.

20 **Page 12, equation (4). Again, I don’t see how a nugget is relevant with the normalized**  
21 **correlation equation (2).**

22 We think that this is answered with the explanation of the nugget effect.

23 **Page 12-13. “... assuming that the involved prior errors for each model are identical in a**  
24 **sense that they share the same statistics and (are?) not correlated.” What does this mean?**

25 We reformulated: “, assuming that models have independent prior errors”.

26 **Page 13, line 8. Richardson et al., 2008, only deals with random sampling error. More**  
27 **recent papers add gap filling errors and friction velocity screening errors to the**  
28 **observational error assessment.**

29 We do not use gap filled data. The Richardson et al. (2008) reference was used to provide a  
30 somewhat realistic estimation of the observation error to be added in the model-model  
31 comparison. Including other sources of error such as resulting from gap filling and friction  
32 velocity screening would lead too far for this context.

1 **There are many years of flux tower data. Only 2007 observations are used. Why? Could**  
2 **the results be strengthened with additional years of observations?**

3 As we are interested in the daily time scale we used only 1 year data. Nevertheless by using  
4 additional years we would not expect to gain more information or significantly different  
5 correlation scales. For example Chevallier et al., 2012 used multiple years and we still have  
6 comparable results.

7 page 9397 line 18-19. We have clarified in this regard, also in response to the second reviewers  
8 comments by adding the following paragraph:

9 “Hilton et al., (2012) studied also the spatial model – data residual error structure using a  
10 geostatistical method. Hilton’s study is focused on the seasonal scale, i.e. investigated residual  
11 errors of seasonally aggregated fluxes. However, the state space (variables to be optimized  
12 considering also their temporal resolution) of current inversion systems is at high temporal  
13 resolution (daily or even three-hourly optimizations). Further, the statistical consistency between  
14 the error covariance and the state space is crucial. Thus the error structure at the daily time-scale  
15 is of interest here, and can be used in atmospheric inversions of the same temporal resolution.  
16 Similar to Hilton’s study we select an exponentially decaying model to fit the spatial residual  
17 autocorrelation.”

18 **Page 14, lines 3-6. Your equations state (but do not define explicitly in the equations) that**  
19 **you are going to evaluate model-data, and model-model differences. The results, however,**  
20 **begin by stating standard deviations in observed NEE. I cannot tell what you have**  
21 **computed. Please clarify the methods and the associated results so that these values can be**  
22 **interpreted by the reader. I think you are presenting time-constant, spatially varying**  
23 **standard deviations across sites, then summed over all times. But the paper should not**  
24 **make me guess what you are computing here, and the methods say nothing about this**  
25 **computation.**

26 This is just the standard deviation of the daily fluxes across all stations and all times to show how  
27 observed and modeled fluxes vary spatially and temporally. As this is a very general metric, we  
28 disagree that this should be stated explicitly in the methods section.

29 **Page 14, line 11. Why is this “in line” with the spatial standard deviation? Model-data**  
30 **differences at a given site are not necessarily related in any way to differences in observed**  
31 **fluxes across sites.**

32 We mean that the fact that the standard deviation of the residuals (modeled - observed) is only  
33 slightly smaller than the standard deviation of the fluxes themselves is in line with the small r-  
34 square values.

1 We have clarified the corresponding sentence: “Those values are only slightly smaller than the  
2 standard deviations of the observed or modeled fluxes themselves. This fact is in line with the  
3 generally low fraction of explained variance with r-square values ...”.

4 **Page 14, line 13-14. What are “site specific correlations” that are presented as a single  
5 value? If these are site specific correlations, what site is being presented? Again, the  
6 methods are not sufficiently clear.**

7 We clarified: “When using site-specific correlations (correlations computed for each site, then  
8 averaged over all sites), ...” .

9 **Page 14, lines 15-18. This sentence’s English needs work. Further, the statistical  
10 comparison (see above comment) isn’t clear. Finally, the model-data difference, if I  
11 understand it, also includes a temporal component since it is summed over time. As best I  
12 can tell, the same data go into both calculations, so I don’t understand how the  
13 authors can draw this fuzzy conclusion about ability to simulate temporal variability better  
14 than spatial variability. Clear, targeted work on this topic has been published elsewhere  
15 but has been neglected in the introduction to this paper. This is also not clearly a main  
16 focus of this paper. I would suggest that you either expand the paper to address this topic  
17 properly, or delete this discussion.**

18 We added a clarification of the computation in Page 9406, line 13, and we also reformulate  
19 regarding also comment from reviewer 2:

20 “When using site-specific correlations (correlations computed for each site, then averaged over  
21 all sites), the average fraction of explained variance increases to 0.38, 0.36, 0.35, and 0.42, for  
22 VPRM10, VPRM1, ORCHIDEE and 5PM, respectively. Note that for deseasonalized time-series  
23 (using a 2nd order harmonic, not shown) the same picture emerges with increased averaged site  
24 specific correlation compared to correlations using all sites. This indicates better performance for  
25 the models to simulate temporal changes (not only seasonal, but also synoptic) at the site level.  
26 Further, the differences between site-specific...”.

27 **Figure 2 is not clear. Is each point a different site? How is this figure related (or not) to the  
28 “site specific correlations” noted on lines 13-14 of Page 14?**

29 We clarified the “site-specific” term in the previous comment. This figure refers to the averaged,  
30 site-specific correlation, for sites that have the same vegetation type. The x-axis lists the different  
31 vegetation types (7 in total). Each site is characterized according to its representative vegetation.  
32 The total number of sites sharing a common vegetation type is shown under the vegetation type.  
33 Hence each of the bars are vegetation (according to x position in the plot) and model (according  
34 to color code) specific.

35 We also clarified the caption: “... Box and whisker plot for for site-specific correlation..”.

1 **Figure 3 has the same problem as Figure 2. Please specify what distribution (sites?) are**  
2 **being illustrated by the box and whisker plots. In addition, the sign of the bias is never**  
3 **defined. Finally, it would be useful to provide a conversion to gC m<sup>-2</sup> yr<sup>-1</sup>.**

4 We clarified as in figure 2.

5 Page 9406, line 24 With respect to the second comment we clarified: “Figure 3 shows the  
6 distribution of bias (defined as modeled – observed fluxes) for different...”.

7 Regarding the conversion, we added in the figure caption: “(for conversion to gC m<sup>-2</sup> yr<sup>-1</sup>  
8 reported values in y axis should be multiplied by 378,7694)”.

9 **Figure 4 has too many lines of similar color and tiny size to be read clearly.**

10 We changed the color for the sites which were excluded from the analysis to better contrast with  
11 the remaining sites. However we note that this figure is not meant to be analyzed for each site  
12 individually (thin red and blue lines) but we rather want to show the average characteristics of  
13 the sites. Then the reader should concentrate on the all-site data (thick black and grey lines) and  
14 on the exponential fit (thick green and dark green lines).

15 **Page 14, line 28. All site? Flux site? Sub-site? The terminology surrounding Figure 4 needs**  
16 **to be cleaned up. I cannot tell what is being plotted. The text appears to contradict the**  
17 **figure caption, and the terminology changes enough to be quite confusing.**

18 “All-site” temporal autocorrelation is explained in page 9403 line 6. The “site data” as given in  
19 page 9402 equation 1 refers to the autocorrelation for each individual site and this is exactly what  
20 is plotted in figure 9 with the thin red lines. These show the temporal autocorrelations for each  
21 site (53 in total). Sub-site is explained in the caption of figure 4. It is computed according to “all-  
22 site” but excluding those sites with large model data bias. To make this more clear, we have  
23 reformulated the figure caption:

24 “Temporal lagged autocorrelation from model-data daily averaged NEE residuals for all models.  
25 Thin red lines correspond to different sites, while the blue thin lines reveal the sites with a bias  
26 larger than +/-2.5  $\mu\text{mol m}^{-2} \text{s}^{-1}$ . The thick black line shows the all-site autocorrelation, and the  
27 thick grey line indicates the all-site autocorrelation but for a sub-set that excludes sites with large  
28 model-data bias (“sub-site”). The dark green line is the all-site exponential fit, and the light green  
29 line shows the all-site autocorrelation excluding the sites with large bias. The exponential fits use  
30 lag times up to 180 days.”

31 **Page 15, lines 2-10. Do not describe the figures in the text. The figures present these results.**  
32 **Discuss the significance of the figures in the text.**

33 We have shortened the text by deleting: “The all-sites correlation for the VPRM model at 10 km  
34 resolution remains positive for lags < 104 days and for lags > 253 days. Weak negative



1 correlations were found in between with minimum value -0.03. In contrast we found only  
2 positive correlation for VPRM at 1 km resolution for the whole year with a minimum value of  
3 0.002. Similarly, ORCHIDEE follows the same patterns with positive correlations for lags < 76  
4 days and for lags > 291. Minimum correlation was found to be -0.09. For 5PM model we also  
5 found only positive correlations. The minimum value was found to be 0.08.”

6 However, some detail is needed in order to compare with the corresponding results by Chevallier  
7 et al., (2012).

8 **The exponential fit appears to be a poor choice. Based on Figure 4, most of the site**  
9 **autocorrelations are below the fit lines at lag times of 30-70 days. Thus the model used to**  
10 **fit these curves is quite biased. It is also more biased for some model-data comparisons**  
11 **than for others. This makes the comparison of decay times misleading, to the point**  
12 **perhaps of being meaningless. I would not publish results based on such a biased**  
13 **approximation of the site-based results. I think this functional fit must be changed. At**  
14 **minimum, the quality of fit must be made very clear, and the logic for keeping this**  
15 **function, despite its relatively poor fit, articulated.**

16 We note that the fit is performed on the all-site and sub-site curves (black and grey lines). The  
17 referee should not compare the fit (light and dark green lines) with the thin red lines.

18 We find  $r^2$  values between the all-site autocorrelation and the exponential fit of 0.94, 0.94, 0.92,  
19 and 0.89 for VPRM10, VPRM1, ORCHIDEE and 5PM respectively. The standard deviations of  
20 the residuals (or RMSE for root mean square error) are 0.040, 0.036, 0.059, and 0.043 for the  
21 different models. Expressed as NRMSE (normalized RMSE, i.e. RMSE divided by the range of  
22 the autocorrelation), this results in values ranging from 0.061 to 0.092, which indicates relative  
23 errors in the fit of less than 10%. We therefore disagree with the reviewer that the fit is poor.  
24 Further we highlight also the importance of this fitting model, which is quite simple and is using  
25 only few parameters, which is a critical point for proper implementation into the inversion  
26 systems.

27 Following the suggestion of the reviewer for a better articulation of our arguments, we added the  
28 following indications:

29 Page 9407 line 14 We added: “However the correlogram exhibits a nugget effect (values ranging  
30 from 0.31 to 0.48 for the different models)”

31 Page 9407 line 17. We added: “The fit has a root mean square error ranging from 0.036 to 0.059  
32 for the different biosphere models. The normalized RMSE (i.e. RMSE divided by the range of  
33 the autocorrelation) results in values ranging from 0.061 to 0.092 indicating relative errors in the  
34 fit of less than 10%.”

35

1 **Page 15 lines 13-15. If the correlation at zero lag is not 1, either your calculations are**  
2 **flawed or equation (2) does not represent your methods.**

3 In spatial statistics nugget effect is a standard term which describes the sharp decrease of the  
4 correlation for infinitesimal temporal separation distances. The value which drops is  $1-a$  and this  
5 is the nugget as described in eq. 3.

6 **Page 15. A root mean square error of a functional fit to an autocorrelation curve is not very**  
7 **meaningful. Evaluating the quality of multiple potential fits to find the best fit to the data**  
8 **would be meaningful and improve the analysis.**

9 RMS error is a measure of agreement. Considering the relatively simple exponential model  
10 which contains only 2 parameters to be optimized, the agreement is satisfactory good.

11 We would like to refer the referee to Chevallier et al. (2012), where a polynomial model was  
12 fitted with much more complexity (5<sup>th</sup> order polynomial with 6 parameters), and the  
13 corresponding RMSE of 0.01 was not much smaller than the RMSE of 0.036 to 0.059 of our  
14 simpler (2 parameter) model.

15 **Page 16, line 5. “not applicable”**

16 We corrected: “...measurements are not applicable”.

17 **Page 16, lines 9-11. The sites were screened because the bias was greater than 2.5  $\mu\text{mol m}^{-2}$**   
18 **s<sup>-1</sup>. But now the text says that the bias for these sites was not greater than 2.5  $\mu\text{mol m}^{-2}$  s<sup>-1</sup>,**  
19 **just “larger than average.” I am confused. Please clarify.**

20 This threshold value is exceeded for all of those sites only for 5PM modeled fluxes. For the rest  
21 of the models the value was over the threshold only for some sites. Nevertheless the bias even  
22 though not exceeding always the threshold value yet, was larger than the averaged bias.

23 We clarified: “... threshold of  $2.5 \mu\text{mol m}^{-2} \text{ s}$  simultaneously for each individual model...”.

24 **There is no functional fit to the aircraft data (figure 5). Given the poor quality of the fits in**  
25 **Figure 4, and I am not convinced that there really is a difference between the two data sets.**  
26 **I would be much more convinced by a comparison of the mean or median values, binned by**  
27 **lag time.**

28 We replaced the wrong figure with the appropriate one which contains the model fit. The fit  
29 found to have 13 days e-folding length with values between 10 and 16 days within 95%  
30 confidence interval. Hence we disagree that this difference is not significant. We also disagree  
31 that the fit is biased.

32

1 **Page 16, line 29. What is the purpose of the root mean square error?**

2 The use of RMSE is very common as a general purpose error metric for numerical predictions.  
3 As RMSE has the useful property of being in the same units as the response variable, we can  
4 then evaluate how good the model performs.

5 We also added the normalized RMSE (divided by the autocorrelation range) estimations:

6 "... 5PM, respectively. The normalized RMSE is found to have values ranging from 0.05 to  
7 0.084 indicating relative errors of the fit less than 9%."

8 **Figure 6. The exponential fit is consistently below the median at distances of 200-400km. I**  
9 **would argue that your correlation computation shows consistently positive values out to**  
10 **approximately 200-400km, which is consistent with Hilton et al., (2013), who performed a**  
11 **similar calculation for North American flux tower sites and model-data differences using**  
12 **VPRM. Again, your exponential fit appears to be biased. I do not believe that quoting the**  
13 **results of a biased fit is sound.**

14 The reviewer is right for the case of VPRM model which might be an exception. A careful look  
15 to the other two models (ORCHIDEE and 5PM) shows that autocorrelation values are well  
16 centered around the exponential for distances longer than 200km. We note also in Hilton's paper  
17 that the 1<sup>st</sup> bin is at 500km with no information for the smaller scales. In our study we used bins  
18 of approximately 100km.

19 **Hilton et al., (2013), published a paper using very similar methods using North American**  
20 **flux towers, a much longer time series of data, and more evaluation of the robustness of**  
21 **the resulting length scales. This was published in Biogeosciences. The results contradict**  
22 **the results presented here in that Hilton et al (2013) found significantly larger length scales**  
23 **for their variogram fits. The similarity is so great that the Hilton et al (2013) paper really**  
24 **should be cited and evaluated with respect to these results.**

25 Hilton et al. (2013) calculated the length scales by considering seasonal mean residuals. In our  
26 study we used daily averaged residuals since this is the temporal scale used in the state space for  
27 regional inverse models. This largely coarser time resolution used in Hilton et al. (2013) is likely  
28 the driver of the differences on the spatial scales.

29 Regarding the robustness of the fit, Hilton et al., 2013 compared with the AIC criterion whether  
30 the exponential fit or the pure nugget (which means no spatial coherence) is better. We note that  
31 they did not fit different error models to evaluate which model was fitting better. We added also:

32 Page 9397 line 18-19 "Hilton et al., (2012) studied also the spatial model – data residual error  
33 structure using a geostatistical method. Hilton's study is focused on the seasonal scale, i.e.  
34 investigated residual errors of seasonally aggregated fluxes. However, the state space (variables  
35 to be optimized considering also their temporal resolution) of current inversion systems is at high

1 temporal resolution (daily or even three-hourly optimizations). Further, the statistical consistency  
2 between the error covariance and the state space is crucial. Thus the error structure at the daily  
3 time-scale is of interest here, and can be used in atmospheric inversions of the same temporal  
4 resolution. Similar to Hilton's study, we select an exponentially decaying model to fit the spatial  
5 residual autocorrelation."

6 Page 9413 line 17 additional paragraph: "Only weak spatial correlations for model-data residuals  
7 were found, comparable to those identified by Chevallier et al. (2012) limited to short lengths up  
8 to 40 km without any significant difference between the biospheric models (31 - 40 km). Hilton  
9 et al. (2012) estimated spatial correlation lengths of around 400km. However we note that  
10 significant differences exist between this study and Hilton et al. (2012) regarding the methods  
11 that were used and the landscape heterogeneity of the domain of interest. With respect to the first  
12 aspect the time resolution is much coarser (seasonal averaged flux residuals) compared to the  
13 daily averaged residuals used here. Furthermore spatial bins of 300 km were used for the  
14 autocorrelation analysis, which is far larger than the approximate bin width of 100 km that were  
15 used in our study. Regarding the second aspect North America has a more homogenous  
16 landscape compared to the European domain. The scales for each ecosystem type (e.g. forests,  
17 agricultural land etc.) are drastically larger than those in Europe as can be seen from MODIS  
18 retrievals (Friedl et al., 2002)."

19 **Figure 7. How is the confidence interval computed? This is not a simple case of computing**  
20 **the standard deviation of a Gaussian. Please explain the methods. I am still dubious of the**  
21 **value of the exponential fit, but in any case the methodology for the confidence interval**  
22 **estimate must be explained.**

23 Page 9404 line 19 We clarified by adding following paragraph: "Confidence intervals for the  
24 estimated model parameters were computed based on the profile likelihood (Venzon and  
25 Moolgavkar, 1987) as implemented within the "confint" function from MASS package inside the  
26 R statistical language."

27 **Figure 7 points out something that is lacking from the primary results reported in the**  
28 **abstract – uncertainty bounds. These results suggest that the uncertainties in the computed**  
29 **correlation lengths are very large. This should be reported in the abstract.**

30 Page 9394 line 16 We added : "... up to few tens of km but with uncertainties up to 100% of this  
31 estimation".

32 **The standard spatial statistical method for Figures 6 and 7 would be a variogram. Why**  
33 **have the authors chosen a different approach?**

34 The choice of the correlogram over the variogram was made since a) Chevallier et al. (2012) also  
35 used the correlogram for a similar analysis, and b) it was simpler to implement in the code.

1

2 **Page 18, lines 5-6. English. ‘it difficult to determine ... where the asymptote lies’ perhaps?**

3 Corrected, now reads : “making it difficult to identify where the asymptote lies.”.

4 **Figure 8 illustrates again how poorly the exponential model fits the data. And the**  
5 **exponential model is not shown on the figure, which is inconsistent with figures 6 and 7.**  
6 **D=35km with a 95% confidence interval of 26-46 km is clearly biased given that none of**  
7 **the aircraft data reaches 1/e of the zero correlation anywhere within that range. The**  
8 **exponential model is poor and should not be used, or only with serious caveats about the**  
9 **biased nature of the fit.**

10 We corrected: the figure now shows the model fit. However, we disagree in this point with the  
11 reviewer, the fit is not biased when looking at the residuals in Fig. 8.

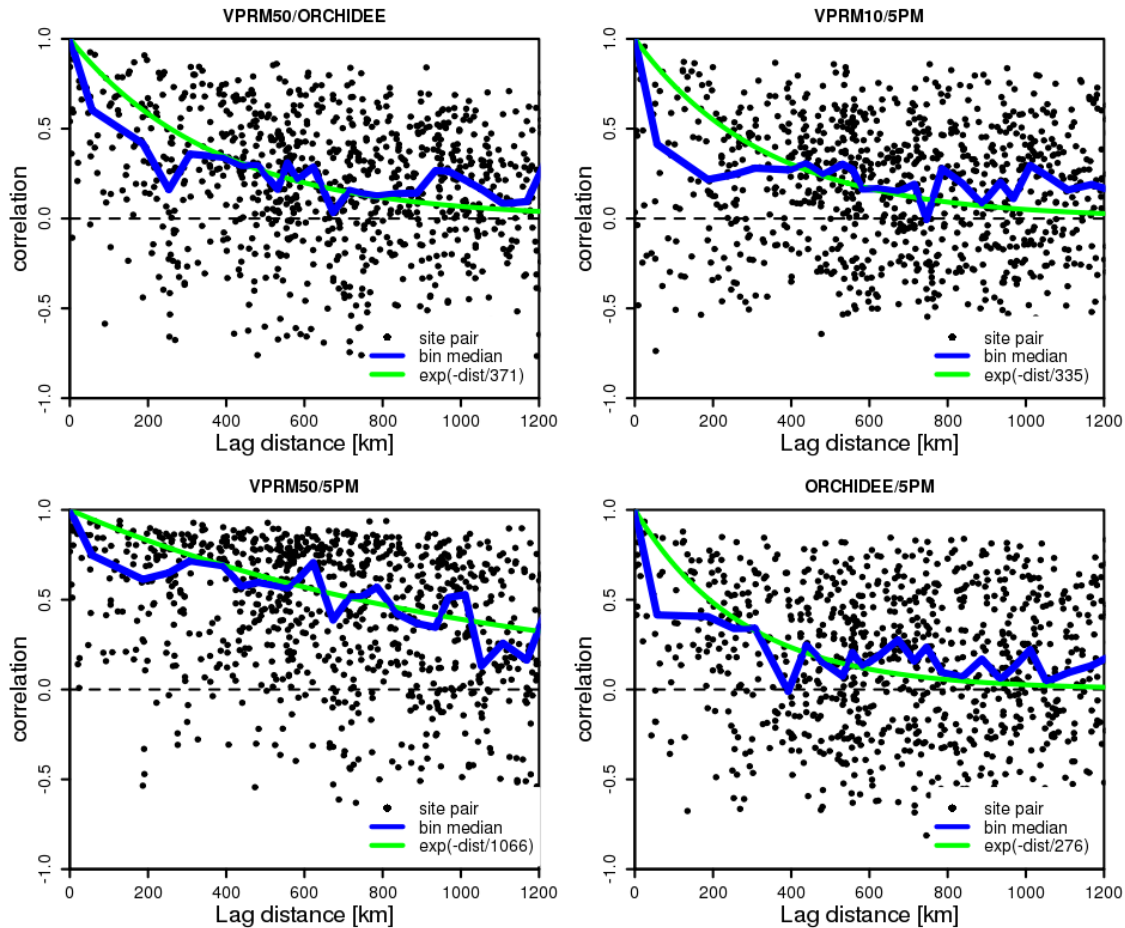
12 **Some of the colors in Figure 9 are nearly indistinguishable when used to plot very thin**  
13 **lines. Please either reduce the content of the figure or find a way to distinguish the**  
14 **different model-model pairs more clearly.**

15 We reduced the information and we added also comparisons between VPRM50/ORCHIDEE  
16 following also the comment from referee #2 regarding the incompatible model resolution.

17 **Figure 9: Why are the individual points not shown, as for the model-data comparison? I**  
18 **can understand reducing the information shown, but I am concerned about the quality of**  
19 **the exponential fit, and it is impossible to evaluate from this figure.**

20 In the figure below, which is similar to figure 7 we present the correlations of differences for all  
21 the different model-model combinations (model-data for figure 7) together with the respective  
22 exponentially decaying model fits. Figure 9 is not meant to present the goodness of fit neither to  
23 show all paired correlations.

24 Nevertheless we plotted the spatial autocorrelation for the paired models in order to evaluate the  
25 exponential fit (referred to as figure 11 in the attached files).



1

2 Above: Distance correlogram for the daily net ecosystem exchange (NEE) differences between  
 3 pairs of models using all sites. Black dots represent the different site pairs; the blue line  
 4 represents the median value of the points per 100-km bin, and the green line shows an  
 5 exponential fit. Results are shown for differences between VPRM at a resolution of 50 km vs.  
 6 ORCHIDEE (top left), between VPRM at a resolution of 1 km and 5PM (top right), between  
 7 VPRM at a resolution of 50 km and 5PM (bottom left), and between ORCHIDEE and 5PM  
 8 (bottom right).

9

10 **Figure 10. Again, the blue-green-purple lines are difficult to distinguish (all dark), and the**  
 11 **red-green lines will be indistinguishable for those who are red-green colorblind.**

12 We corrected the plot and further we reduced the amount of information. We also added the  
 13 VPRM50/ORCHIDEE comparisons.

14

1 **Discussion and conclusions. The results are compared only to Chevallier’s publications,**  
2 **and that comparison is limited to the temporal autocorrelation. There is insufficient effort**  
3 **to put these results into the context of prior work on this topic.**

4 Page 9413, lines 17 We added a new paragraph: “Only weak spatial correlations for model-data  
5 residuals were found, comparable to those identified by Chevallier et al. (2012) limited to short  
6 lengths up to 40 km without any significant difference between the biospheric models (31 - 40  
7 km). Hilton et al. (2012) estimated spatial correlation lengths of around 400km. However we  
8 note that significant differences exist between this study and Hilton et al. (2012) regarding the  
9 methods that were used and the landscape heterogeneity of the domain of interest. With respect  
10 to the first aspect the time resolution is much coarser (seasonal averaged flux residuals)  
11 compared to the daily averaged residuals used here. Furthermore spatial bins of 300 km were  
12 used for the autocorrelation analysis, which is far larger than the approximate bin width of 100  
13 km that were used in our study. Regarding the second aspect North America has a more  
14 homogenous landscape compared to the European domain. The scales for each ecosystem type  
15 (e.g. forests, agricultural land etc.) are drastically larger than those in Europe as can be seen from  
16 MODIS retrievals (Friedl et al., 2002).”

17 **Page 21, lines 26-28. The observational errors are in your calculations. It is not neglected.**  
18 **It cannot be isolated and removed, but it is not neglected.**

19 We deleted part of the sentence, which now reads: “Of note is that the eddy covariance  
20 observation error has no significant impact on the error structure, as the addition of an  
21 observation error to the analysis of model-model differences had only minor influence on the  
22 error structure.,.

23 **End of page 24 beginning of page 25. This is an interesting discussion. Again, the methods**  
24 **used for this interesting calculation are opaque. Please explain. Propagating these error**  
25 **estimates is not trivial. How was this done?**

26 The theoretical approach for this calculation is based on Rodger (2000) by introducing an  
27 aggregation operator. We explain with the following equation where “ $\times$ ” represents matrix  
28 multiplication notation:

29  $u \times Q_c \times u^T$ .  $Q_c$  is the full prior error covariance matrix with dimensions equal to the product  
30 between number of regions (grid-cells) and number of timesteps. With our setup this translates  
31 into  $184 \cdot 104 \cdot 8 \cdot 365 = 55877120^2$  elements (lon·lat·timesteps·days). The main diagonal contains  
32 the model-data difference scaled down to account for the difference in spatial resolution of the  
33 state space. The off-diagonal elements contain the spatial and temporal correlations.  $u$  is a scalar  
34 operator that aggregates the full covariance to the target quantity (i.e. domain-wide and full  
35 year).

1 Page 9416 line 24. We added: "... over longer time periods. To aggregate the uncertainty to large  
2 temporal and spatial scales, we used the following equation (after Rodgers, 2000):

$$3 \quad Ua = u \times Q_c \times u^T \quad (7)$$

4 Where "×" denotes matrix multiplication,  $Q_c$  is the prior error covariance matrix and  $u$  a scalar  
5 operator that aggregates the full covariance to the target quantity (e.g. domain-wide and full  
6 year)."

7 Reference added:

8 "Rodgers, C., D. Inverse methods for Atmosphere Sounding: Theory and Practice, World Sci.,  
9 River Edge, N. J., 2000" page 30 line 20

10

11 **Page 22 line 25-28 Same section: Comparing your aggregated error estimate to the**  
12 **range of existing continental-scale flux estimates (e.g. Peylin et al, 2013) would be more**  
13 **useful than the very limited analysis presented.**

14 We agree with the reviewer and we added: "This value is also 8 times smaller when comparing it  
15 to the variance of the signal between 11 global inversions reported in Peylin et al., (2013) which  
16 was found to be 0.45 GtC/y, proving that the aggregated uncertainties are unrealistically small."

17 **Same section: I agree that this analysis (pending evaluation of the unknown methods)**  
18 **would strongly suggest that the total continental-scale, annual flux errors are seriously**  
19 **underestimated, and I agree that this is an important issue to point out. This should be**  
20 **part of the abstract, as it leads to significant uncertainty regarding the validity of the**  
21 **correlation lengths. The current abstract suggests no such uncertainty regarding the**  
22 **conceptual model promoted in this paper.**

23 Page 9394 line 16. We added: "Propagating this error structure to annual continental – scale  
24 yields an uncertainty of 0.06 Gt C and strongly underestimates uncertainties typically used from  
25 atmospheric inversion systems, revealing the existence of another potential source of errors.  
26 Long spatial e-folding correlation lengths up to several...".

27 **Page 25, paragraph starting with "Exponentially decaying ...". This paragraph begins to**  
28 **give reasons for using an exponential model for the correlations. Some notes below about**  
29 **this discussion:**

30 **1) This discussion belongs earlier in the paper. It presents the logic for using this fit.**

31 We moved this paragraph and now is located in page 9412 early in the discussion.



1 **2) The reasons given are entirely reasons of simplicity and convenience, not accuracy of the**  
2 **fit. I would suggest that the best job of describing the correlations should be the primary**  
3 **goal of this paper. Considering how to simplify these correlation functions to make them**  
4 **convenient is another problem. I have already noted that I believe the exponential fit is so**  
5 **poor that it is significantly misleading. The analysis would be improved by evaluating**  
6 **different fits and finding what fits are best given the data.**

7 Regarding the simplicity and convenience we will refer the reviewer to Hilton et al., 2013. In  
8 that study two models were used. An exponentially decaying model and one that uses a pure  
9 nugget effect. The nugget only model is equivalent to an exponentially decaying model with the  
10 length scale of zero, which means no spatial correlation can be detected. So the assessment of  
11 whether the nugget only model or the exponentially decaying model in Hilton is appropriate,  
12 could simply be done, by assessing if the length scale in the exponentially decaying model is  
13 significantly different from zero.

14 Further this study tries to give insights in the error structure targeting to describe prior  
15 uncertainties with a relevant way that the atmospheric inverse modeling community may benefit  
16 from it. The exponential model is widely used by this community. Describing the correlations  
17 with a pure mathematical way which makes a convenient fit but is not being used from the  
18 inverse modeling community is also of less importance. However we do not believe that the fits  
19 are poor.

20 **3) Lines 20-25. Computational simplicity is not a good reason to use the wrong correlation**  
21 **length. This is a disturbing discussion.**

22 We disagree with the reviewer at this point that the correlation lengths are wrong, see also  
23 responses to previous comments. The exponential model might be simple but performs  
24 satisfactory good. Further despite the hesitation of the reviewer regarding the importance of  
25 computational efficiency, this is a major issue for regional scale inverse modeling.

26 **Equation (7) does not exist in the paper.**

27 Page 9417 line 23 Corrected, the sentence now reads “Using the same hyperbolic equation for  
28 the spatial correlation...”.

29 **Page 26, lines 3-4. English needs work. And I’m not sure what is meant by “for the short**  
30 **spatial scales.” And what studies have already used the correlation lengths derived from**  
31 **this study? I’m not sure that the ‘future work’ needs to be part of this paper**

32 With this sentence “...short spatial scales” we mean that numerous studies in atmospheric  
33 inverse modeling use significantly larger spatial correlation scales than those derived in this  
34 study. However most of the inversion systems already use temporal correlation scales of around  
35 1 month, which is in line with our findings.