

Response to Rypdal and Rypdal (*italics*):

This referee report will comment on the revised manuscript and the authors' reply to the first referee report. Martin Rypdal has joined in to speed up and secure the quality of the review process.

We have been asked to help the editor make a final decision about publication of this paper in ESD. Our general position is that the issue is important, and although we find many weaknesses in this paper, we have no desire to prevent it from being published, provided our comments are published as part of the public discussion. We find similar weaknesses in other of Shaun Lovejoy's publications, and it may be better that these issues are publicly discussed.

Thus, our recommendation is that the editor reads our comment and use it as a background for the decision.

The subadditivity in the ZC model

In the original manuscript the analysis leading to the subadditivity in the ZC-model was based on an incorrect neglect of a cross-term leading to Eq. (5). The authors were asked to use the full expression and correct the result in Fig. 3 according to this. In their reply the authors acknowledge the error, state that they have corrected it in Fig. 3, but claim that the correction does not alter the conclusion. However, the actual revision they have made is rather disturbing. In the text the incorrect Eq. (5) and the incorrect argument for it is maintained, and so is the conclusion that the "theoretical additive fluctuation level is reduced by a factor ≈ 2.5 ." Then they add the following paragraph: "It should be noted that the latter holds assuming independence (pink curve in Fig. 3c) of the solar and volcanic forcing. For comparison, the purple curve in Fig. 3c illustrates the results obtained when analyzing the series constructed by directly summing the two response series (instead of assuming statistical independence). It is clearly seen that the basic result still holds but it is a little less strong (a factor of ≈ 2). The reason for the difference is that the cancellation of the cross terms assumed by statistical independence is only approximately valid on simple realizations, especially at the lower frequencies where the statistics are worse."

***Authors:** It was not "incorrect" to neglect the cross terms! The cross terms are exactly zero under the stated hypothesis that the system is linear and the forcings are statistically independent! There is no "error" here! However, since there is a single realization, there will be fluctuations away from zero (note that for every lag Δt , there are generally still many fluctuations, and for the cross terms, these tend to cancel). We happily added the lower curve in the figure to allow the readers to see the size of this statistical fluctuation. That is why it is quite useful to show both curves and to compare them. Neither is "wrong" they are just valid under the slightly different hypotheses that were clearly stated. Perhaps the referee was misled by a typo, the last sentence should have read: "on single realizations", not "on simple realizations".*

In Fig. 3b they keep the incorrect curves for additive response and the ratio, and do not plot the correct curve, and in Fig. 3c they also keep the incorrect curve for the ratio, but

here they also add the correct curve.

Authors: *see the above, there are no errors.*

For unknown reasons they plot this ratio as a logarithm, so the reader will have to use a calculator to find out what the correct ratio is. The mean \log_{10} ratio in the narrow scale range 250-1000 yr the ratio is ≈ 1.4 . In the revised text in Section 10 3.4 (line 339) the authors claim that this corrected ratio is ≈ 2 , and in the caption of Fig. 3 that it is ≈ 1.6 . In the concluding section (line 522) it has again risen to “a factor ≈ 2 -2.5.”, This way of “correcting” a flawed analysis is unacceptable.

Authors: *Again, there was no flaw. However, we accept that the magnitude of the effect was somewhat qualitatively judged. In the new version of figure 3c, reference lines with actual ratios (not logs) is given. Ratios of 2 and 1.5 for the two different curves are quite accurate as any reader can now judge by visual inspection.*

As a minimum the authors should take the following actions: 1. Eq. (5) and the arguments leading to it should be taken out of the manuscript. The text gives the reader the false impression that linearity implies additivity of variances.

Authors: *It was quite clearly stated that the additivity of the variances requires an assumption of statistical independence of the forcings. This has been reinforced by some new wording.*

The truth is that linearity implies additivity of the responses, and this is what should be tested, and nothing else. Moreover, making this approximation does not simplify anything, so there is no scientific reason to make it.

Authors: *From the data we have, the best way to test the linearity of the response is to consider the statistics, and – unless the referee believes that solar and volcanic forcings are causally linked - assuming statistical independence is not unreasonable.*

In Fig. 3b the incorrect curve for the additive response and the ratio should be replaced by the correct one, and in Fig. 3c the incorrect curve should be removed. In all panels in Fig. 3 it would be easier to grasp the content if the vertical axes are linear, not logarithmic. The authors must refrain from cheating with numbers. The true result of the correct analysis is that the rms-ratio is ≈ 1.4 ; not 2.5, not 2, and not 1.6. The paper should present only the correct number 1.4.

Authors: *See above.*

A linear framework cannot ignore internal noise from stochastic forcing. But the problems do not end here, because the output data from the ZC model contains a noise contribution from internal variability in addition to the model response to forcing. Ideally, if the ensemble size is infinite, the ensemble mean could be interpreted as a “deterministic” response, but as is clearly shown in Fig. 1 of Mann et al. (2005) the

magnitude of the high-frequency noise depends on the number of realisations averaged over (they show results for 5, 20, and 100 realisations), and even after averaging over 100 realisations the RMS of the noise is as high as the deterministic response. For instance, the 40 yr moving average in Mann et al., Fig. 1c (the maroon curve) can be interpreted as the deterministic response to solar forcing (the blue curve), and the difference between the red and the maroon curve is the noise that remains after averaging the internal variability over 100 realisations.

This noise turns out to contribute to the subadditivity and completely explains the factor 1.4, even if the deterministic response is strictly linear!

Proof

(...not reproduced due to problems with the math symbols....)

Authors: *We thank the referee for pointing out the possibility that the internal variability might still be important even at centennial scales and even after averaging over 100 realizations. We do not contest the mathematics, but we do contest his assertion that the internal variability is responsible for the multicentennial response of the models under solar, volcanic or combined solar and volcanic forcings. While it is true that a definitive answer to this requires running the model in “control mode” so as to capture only the internal variability (as was done in for the GISS model, see fig. 4), there are nevertheless several reasons why the internal variability is almost certainly smaller than the response due to the forcings:*

- i) *We can get a typical order of magnitude of the internal variability from the GISS model, fig. 4; we see that for a single realization - without averaging over 100 realizations as in fig. 3a – that the typical centennial variability is $\approx \pm 0.05K$ and decreasing with a power law with exponent $\approx \xi(2)/2 \approx -0.2$. After averaging for 100 realizations, we expect this to decrease by $(100)^{0.5} = 10$ i.e. to $\pm 0.005K$. This is much smaller than the centennial scale variability of the ZC responses in fig. 3a (from the graph, these are about $\approx \pm (10^{-1.2})/2 \approx \pm 0.03K$.*
- ii) *We can use the fact that a) the observed responses are upper bounds on the internal variability and b) that the internal variability must decrease with scale (otherwise the model’s climate diverges rather than converges for long times). Exponents near the GISS value $\xi(2)/2 \approx -0.2$ are common, see e.g. [Lovejoy et al., 2013]. From fig. 2, we see that the ZC solar response at ≈ 20 years is $\pm 0.03K$, so this is an upper bound for the internal variability at all scales longer than ≈ 20 years. However, over the range ≈ 50 -500 years (relevant for the subadditivity conclusion), the solar response variability is considerably larger than this noise value: from the graph, $\approx \pm (10^{0.8})/2 \approx \pm 0.08K$.*

We conclude that it is unlikely that the internal variability is strong enough to account for the results. We have nevertheless added some comments to this effect in the paper.

Multifractality and linearity in responses

The equations that constitute the climate models studied in this paper are nonlinear.

The noise (internal variability) in the models (as appearing in control runs) is a result of nonlinear dynamical processes. However, the issue studied in this paper is whether the responses to external forcing add up. It is not easy for us to see the relevance of all the multifractal formalism presented in Section 4, but the essence seems to be a mathematical corollary claiming that linearity in the response implies that the intermittency function $K(q)$ is the same for forcing and response (lines 400-410). However, as was pointed out in the first referee report, this result depends on a particular form of the linear response function, namely that it is a power law, which is necessarily an idealization, since it leads to infinite responses on long time scales. In their reply the authors claim that there is no cut-off of the power-law tail of the response function. This is in direct contradiction to their claim that GCMs do not reproduce low-frequency (multicentennial) variability. If the linear response function has a power-law tail up to multicentennial scales, then a white-noise stochastic forcing would produce an internal variability with a power-law spectrum up to these scales, and this would be reflected in the spectrum of control runs. This is what the authors claim does not happen in GCMs on lines 82-83. In the referee report the authors were asked to test their results by using also other response functions, but they have refused to do that.

Authors:

In the text we have used the appropriate conditionals, the statements are correct. For the analysis, all we claim is that if the response is linear that it cannot change the nonlinear part of the moment scaling exponent $\xi(q)$, that it can only make a linear change in $\xi(q)$, not a nonlinear one.

The lack of accounting for internal variability also arises in the multifractal analysis. If the data output from the climate model is one realisation (or the average over limited number of realisations), and this output is modeled as a linear response to a forcing, then it is necessary to assume the existence of a stochastic forcing. Otherwise, the linear model cannot account for the internal noise in control runs of GCMs and the noisy output in the ZC model with smooth solar forcing observed in Fig. 1c of Mann et al. (2005). Even more striking is Fig. 1a of Mann et al. (2005) for volcanic forcing, where the 100-realisation mean is completely dominated by this noise for all times except 1-2 years after major eruptions. But this stochastic forcing is not included in the computation of structure functions of the forcing in Fig. 6. Thus, the linear model put to test is a model that cannot produce internal variability, but it is tested against data that is known to exhibit such variability.

Authors:

Again, the referee would only be correct if the noise of the average of 100 realizations was comparable to the model response, yet at short times (several years), there is no argument that the model response – even without averaging over 100 realizations - is strongly displays responses from volcanism.

The multifractal analysis of the climate model output is mostly an analysis of the internal noise, and hence there is no reason to find a linear connection between forcing and output. The weaknesses discussed above are all symptoms of a missing hypothesis-testing

strategy. The linearity hypothesis is not clearly formulated, which creates confusion with respect to the rôle of internal variability. There are of course always some deviation from linearity, but is no attempt in this paper to test if the observed deviations are significant. This is particularly problematic for the multifractal analysis. Here results are notoriously unconvincing unless they are followed up by Monte Carlo simulations of linear response models with a range of response functions and stochastic forcing. Such simulations will provide information about biases and statistical errors, and could be used to reject a linear response hypothesis if this hypothesis is false.

Authors: *We agree that systematic model studies are needed including in particular control run outputs.*

Other comments

The first referee-report recommended a drastic shortening of this review-style paper, and therefore refrained from commenting on the extensive discussions of issues not directly related to the new results presented. The authors have not followed this recommendation (which was also made by referee #3), so we feel that we in this final referee report should also make some comments on these parts of the manuscript.

Authors: *The referee was of the opinion that there was too much review and explanatory material. We believe that the community is not fully aware of many aspects of nonlinear geophysics approaches to the climate and to climate models, so that this material is useful. It was apparently also the opinion of at least one of the other referees.*

However, to demonstrate good faith, we did remove several pages of text to follow the referee's opinion.

Section 1.1

Line 36-44: Here the authors give the impression that assuming a linear response is synonymous with neglecting feedbacks. This is of course nonsense. Feedbacks, like the ice-albedo feedback, can in many cases, and on different time scales, be accounted for by constant feedback factors, i.e., in a linear framework. But there are of course also situations where this is impossible, in particular at times when the climate system undergoes transitions. The validity of a linear approximation is not primarily a question of time-scale, but of the structural stability of the system. For instance, if the climate system is close to a tipping point, the nonlinearities that are responsible for the bifurcation will be important and detectable in the response to external forcing as well as to "internal" stochastic forcing. The papers by other authors referred to do not assume linearity as a general feature of the climate system, but as reasonable approximation for global responses to realistic forcing in the present stable Holocene climate state.

Authors: *The referee is being overly sensitive. In section 1.1 we specifically discuss "strongly nonlinear interactions/feedbacks" giving the celebrated example of Daisyworld. These cannot be understood with constant feedback factors.*

Section 1.2

This subsection is unnecessary as a background for assessing the results presented in the paper. Our understanding of the different spectral regimes is very different from what is presented here. For instance, in a recent paper just accepted for ESD (where Shaun Lovejoy was a referee) we demonstrate that the spectra of global temperature can be described as a relatively non-intermittent background noise with spectrum of the form f^{-1} , superposed on a succession of rapid (nonlinear) transitions. The authors end the subsection with stating that opinions diverge on the value of the global transition scale. This is not a correct description of the disagreement, since we don't accept the very notion of a transition scale in the Holocene climate.

Authors: *We disagree but we do cite disagreement, in particular, the principle point of disagreement, the value of the transition scale in the Holocene and its possible regional variations.*

Section 1.3

Line 82-83: "From the point of view of the GCMs, the low-frequency (multicentennial) variability arises exclusively as a response to external forcings, although potentially - with the addition of (known or currently unknown) slow processes such as land-ice or biogeochemical processes - new internal sources of low-frequency variability could be included." This is a gross overstatement. In a recent paper (Fredriksen and Rypdal, 2015) we analyse an ECHO-G control run and a large number of long control runs in the CMIP5 ensemble. ECHO-G shows near-perfect long-memory power-law scaling up to a millennium, while most CMIP5 models exhibit a flattening tendency of the spectrum on scales larger than a century. Østvand et al. (2014) analysed millennium-long control runs and runs with full forcing of the ECHO-G and COSMOS models. Also analysed were forced runs of HadCM3, and LOVECLIM models (no controls were available), and here the residuals after subtracting a forced, deterministic response were also subject to scaling analysis. The results showed no enhanced variability on the multicentennial scales by inclusion of forcing, and the internal variability showed no deviation from power-law scaling with $0.75 < \alpha < 1$ on the scales available in these experiments. This demonstrates that GCMs differ among each other when it comes to variability on multicentennial scales. Various techniques are used to deal with phenomenon of model drift, which will influence variability on large time scales. But these issues with the GCMs cannot be summarised as the authors do on lines 82-83.

Authors: *I think we basically agree here – see our publication [Lovejoy et al., 2013]. The point is that we attribute the relatively weak low frequency response of the models to a problem with the models whereas the referees are of the opinion that weak (power law or nearly power law multicentennial, multimillennial) model response is realistic.*

Section 3.1 lines 187-192

Here the authors suggest to use as a test of model performance to check the equality $\sigma_{Tsim} = \sigma_{Tobs}$. We don't disagree with that, but what the authors do not specify is that in order to compare statistics of models with the statistics of observations, σ_{Tsim} should be computed from single realizations of model output, and not from ensemble averages

over many model runs as they do in their analysis of the ZC model.

Authors: *Yes, of course. The issue of the internal variability – the part affected by the averaging – is now discussed in the text.*

References

Fredriksen, H.- B., and Rypdal, K.: Scaling of Atmosphere and Ocean Temperature Correlations in Observations and Climate Models, *J. Climate*, e-view, doi: <http://dx.doi.org/10.1175/JCLI-D-15-0457.1>, 2015.

Mann, M. E, Cane, M. A., Zebiak, S. E., and Clement, A.: Volcanic and solar forcing of the tropical pacific over the past 1000 years, *J. Climate*, 18, 447-456, 2005.

DOI: 10.1175/JCLI-D-15-0183.1.

Østvand, L., Nilsen, T., Rypdal, K., Divine, D., and Rypdal, M.: Long-range memory in internal and forced dynamics of millennium-long climate model simulations, *Earth. Syst. Dynam.*, 5, 295-308, doi:10.5194/esd-5-295-2014, 2014.

Lovejoy, S., D. Schertzer, and D. Varon (2013), Do GCM's predict the climate.... or macroweather?, *Earth Syst. Dynam.* , 4, 1–16 doi: 10.5194/esd-4-1-2013.