

1           **Comment on “Comparison of low-frequency internal climate variability in**  
2                                   **CMIP5 models and observations”**

3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

by

Sergey Kravtsov<sup>1</sup>

University of Wisconsin-Milwaukee

28 June 2017

Submitted to the *Journal of Climate*

---

<sup>1</sup>*Corresponding author address:* Dept. of Mathematical Sciences, Atmospheric Science Group, University of Wisconsin-Milwaukee, P. O. Box 413, Milwaukee, WI 53201. *E-mail:* [kravtsov@uwm.edu](mailto:kravtsov@uwm.edu).

1 **Abstract**

2 In a recent article (*J. Climate*, doi:10.1175/JCLI-D-16-0712.1), Cheung et al. (2017)  
3 [hereafter, C2017] apply a semi-empirical methodology to isolate internal climate  
4 variability (ICV) in CMIP5 models and observations. The essence of their  
5 methodology is to subtract scaled CMIP5 multi-model ensemble-mean (MMEM)  
6 from individual model simulations and from the observed time series of several  
7 surface-temperature indices. C2017 detect large differences in both the magnitude and  
8 spatial patterns of the observed and simulated ICV, as well as large differences  
9 between the historical (simulated) ICV and pre-industrial (PI) control CMIP5  
10 simulations. Here it is shown that subtraction of the scaled MMEM from CMIP5  
11 historical simulations produces a poor estimate of the modeled ICV due to the  
12 difference between the scaled MMEM and a given model's true forced signal  
13 masquerading as ICV. The resulting phase and amplitude errors of the ICV so  
14 estimated are large, which compromises most of the C2017 conclusions. By contrast,  
15 an alternative methodology based on forced signals computed from individual model  
16 ensembles produces a much more accurate estimate of the ICV in CMIP5 models,  
17 whose magnitude is consistent with the PI control simulations and is much smaller  
18 than any of the semi-empirical estimates of the observed ICV on decadal and longer  
19 time scales.

20

## 1        **1. Introduction**

2            Accurate identification of internal (unforced) climate variability (hereafter,  
3        ICV) is essential for understanding and predicting long-term climate change apparent,  
4        for example, in a variety of the Northern Hemisphere surface-temperature indices.  
5        One of the ways to isolate ICV in climate-model historical simulations is to utilize  
6        ensembles of simulations available through the Coupled Model Intercomparison  
7        Project phase 5 (CMIP5: Taylor et al. 2012). Here, averaging historical simulations  
8        over multiple realizations of a climate-index time series of a given model provides an  
9        estimate of the forced signal, while differencing each simulation with the forced  
10       signal so computed provides an estimate of the simulated ICV. Steinman et al.  
11       (2015a,b) correctly noted that forced signals inferred via such single-model ensemble  
12       means (SMEM) are contaminated by insufficient averaging of ICV — due to a small  
13       number of historical realizations typically available in individual model ensembles, —  
14       and argued that additional averaging across the entire multi-model ensemble (multi-  
15       model ensemble mean: MMEM) is required for proper estimation of the true forced  
16       signal. They claimed that the MMEM rescaled, via linear regression, — to match  
17       either the observed or simulated climate change — provides unbiased estimates of the  
18       forced signal in observations and CMIP5 historical simulations, and can thus be used  
19       to obtain unbiased estimates of (both observed and simulated) residual ICV. Along  
20       the same lines, Frankcombe et al. (2015) introduced two-factor and three-factor  
21       scaling methods for identification of the forced signal and ICV; these methods  
22       utilized, in addition to the MMEM based on historical simulations incorporating all  
23       forcing factors, the MMEMs from the historical simulations forced by the  
24       greenhouse-gas or natural forcing only. In this comment, we will refer to the single-

1 factor, two-factor and three-factor estimates of the forced signal and the associated  
2 estimates of ICV as the MMEM1, MMEM2 and MMEM3 methods, respectively.

3       Recently, Kravtsov et al. (2015) and Kravtsov and Callicutt (2017) showed  
4 that Steinman et al. (2015a,b) arguments are flawed, and it is impossible to  
5 characterize a rich spectrum of forced signals realized in the individual CMIP5  
6 models using the single rescaled MMEM signal. In particular, these authors  
7 demonstrated that subtraction of the rescaled MMEM from the ensembles of  
8 individual-model simulations does not produce independent (uncorrelated)  
9 realizations of this model's ICV, which are instead contaminated (and often  
10 dominated) by the common difference between this model's true forced signal and the  
11 rescaled MMEM. They also showed that smoothed (5-yr low-pass filtered) SMEM of  
12 an individual model provides a much more accurate estimate of this model's true  
13 forced signal. Kravtsov and Callicutt (2017) further used the appropriately rescaled  
14 SMEMs of individual models in conjunction with the observed time series of a given  
15 climate index to estimate forced signal (and ICV) in observations. Meanwhile,  
16 Kravtsov (2017) developed a spectral variance-inflation method to correct for the  
17 errors associated with insufficient SMEM averaging and to produce unbiased  
18 estimates of ICV in CMIP5 historical simulations. We will hereafter jointly refer to  
19 the above methods for estimating the observed and simulated ICV over the historical  
20 period as KC2017.

21       The above criticisms of the MMEM1 method also fully apply to a recent work  
22 by Cheung et al. (2017) [hereafter C2017], who used it as the main tool to compare  
23 the observed and CMIP5 simulated ICV of surface temperature over the course of the  
24 twentieth century. These authors found considerable mismatches in the magnitude of  
25 the observed and modeled ICV on decadal and longer time scales, as well as

1 mismatches in the ICV between CMIP5 historical runs and pre-industrial (PI) control  
2 runs, with the simulated ICV over the historical period being generally much lower  
3 than the observed ICV, but higher than the variability in the PI control runs. The main  
4 point of this comment is to show that these results are quantitatively inaccurate due to  
5 biases of the MMEM1 method, whereas KC2017 method provides means to a much  
6 more accurate estimation of both the amplitude and the phase of the true ICV in  
7 historical simulations.

8 Further presentation is organized as follows. The data and analysis  
9 methodology are described in section 2. Section 3 analyzes the performance of the  
10 various MMEM methods, as well as KC2017 method in isolating ICV in the  
11 simulations from the Community Earth System Model (CESM) Large Ensemble  
12 Project (LENS: Kay et al. 2015): [http://www.cesm.ucar.edu/projects/](http://www.cesm.ucar.edu/projects/community-projects/LENS/)  
13 [community-](http://www.cesm.ucar.edu/projects/community-projects/LENS/)  
14 [projects/LENS/](http://www.cesm.ucar.edu/projects/community-projects/LENS/), in which a large ensemble size permits a fairly accurate  
15 determination of the true forced (and internal) variability. The application of these  
16 methods to CMIP5 multi-model ensemble is discussed in section 4, and conclusions  
17 are presented in section 5. Supplemental Materials contain further information in  
18 support of this paper’s main results, as well as the link to the data and code for all of  
19 the analyses performed here.

## 20 **2. Data and methods**

21 We utilized the same historical (see **Table 1**), historicalNat (natural forcing  
22 only), historicalGHG (greenhouse-gas forcing only) and PI control CMIP5  
23 simulations (1880–2005) as in Kravtsov (2017). The sub-ensemble of CMIP5  
24 historical simulations (17 models, 111 simulations) was smaller than the ensemble

1 used in C2017 since we only considered the models with four or more historical  
2 simulations available to be able to apply the KC2017 method; note, however, that the  
3 MMEM based on this sub-ensemble is virtually indistinguishable from the MMEM  
4 based on the entire CMIP5 ensemble (not shown). We also utilized the 40-member  
5 historical ensemble (1920–2005) and long PI control simulations from the LENS  
6 project (Kay et al. 2015). The observed surface-temperature data we used was also  
7 identical to the data used in Kravtsov (2017) [and C2017]; one notable difference with  
8 C2017 here is our use of the twentieth-century reanalysis (20CR; Compo et al. 2011)  
9 to compute the Northern Hemisphere mean temperature (instead of GISTEMP data in  
10 C2017). Our results do not depend on the choice of surface-temperature data set.

11       Following Steinman et al. (2015a) and C2017, we considered Atlantic  
12 Multidecadal Oscillation (AMO), Pacific Multidecadal Oscillation (PMO) and  
13 Northern Hemisphere Multidecadal Oscillation (NMO) indices; the AMO and PMO  
14 indices were computed as sea-surface temperature averaged over the North Atlantic  
15 and North Pacific ( $0^{\circ}$ – $60^{\circ}$ N), respectively, and the NMO was computed as a weighted  
16 sum of the  $0^{\circ}$ – $60^{\circ}$ N surface air temperature averaged over ocean and land, with the  
17 weights set to 0.61 and 0.39, respectively.

18       We computed forced signals and ICV estimates (in both observations and  
19 model simulations) using MMEM1, MMEM2, MMEM3 and KC2017 methods. In  
20 applying KC2017 method to LENS simulations, we considered random sub-  
21 ensembles of five simulations (which is a typical size of an individual model  
22 ensemble in CMIP5 models) containing a given simulation to compute 5-yr low-pass  
23 filtered ensemble mean as a first-guess estimate of the forced signal and the residual  
24 ICV in that simulation, then applied this procedure to all other simulations. The  
25 resulting 40 estimates of ICV were further variance-bias corrected using the

1 procedure developed in Kravtsov (2017). For the LENS ensemble, we also computed  
2 the “true” forced signal as the grand (40-member) ensemble mean, and the associated  
3 “true” residual ICV for each of the 40 historical realizations. As in C2017, we further  
4 computed low-pass filtered versions of the AMO, PMO and NMO ICV using the  
5 data-adaptive filter of Mann (2008) with cutoff frequencies of 0.1, 0.05 and 0.025  
6 cycles per year (or corresponding smoothing time scales of 10, 20 and 40 years).

7

### 8 **3. Quantifying performance of different methods for isolating internal** 9 **variability in historical simulations using CESM LENS ensemble**

10 **Figure 1** shows the magnitude of internal variability in LENS simulations  
11 estimated by different methods (top row), as well as the root-mean-square (rms) error  
12 of each method with respect to the “true” internal variability based on the subtraction  
13 of the grand ensemble mean from individual simulations (bottom row), as a function  
14 of smoothing time scale. It turns out that the rescaled CMIP5 MMEM is fairly close to  
15 the true forced signal in the CESM model, and the MMEM-based methods generally  
16 provide decent estimates of ICV in CESM historical runs (see **Fig. S1** of  
17 Supplemental Materials). Still, the KC2017 method is more accurate than MMEM1  
18 method in identifying low-frequency ICV, with the rms error of the 40-yr low-pass  
19 filtered ICV estimated using KC2017 method being 75% of the MMEM1 error for  
20 PMO (Fig. 1b), and about 40% of the MMEM1 error for AMO and NMO (Figs. 1d, f).  
21 The MMEM2 and MMEM3 methods perform worse than both MMEM1 and KC2017  
22 methods across the entire range of time scales and for all indices. The MMEM1 and  
23 MMEM2 methods tend to overestimate the magnitude of the true ICV at low  
24 frequencies, whereas the MMEM3 method slightly underestimates this magnitude

1 (Figs. 1a,c,e). Note that the amplitudes of ICV estimated by all methods are overall  
2 consistent; hence, large rms errors of the MMEM2 and MMEM3 methods (Figs.  
3 1b,d,f) are due to a combination of both amplitude and, most importantly, phase errors  
4 with respect to the true ICV (in other words, the magnitude of the estimated ICV may  
5 match the magnitude of the true ICV, but it may still be poorly correlated with the  
6 true ICV; see examples in Fig. S1 of Supplemental Materials). In summary, KC2017  
7 method provides the best estimate of the ICV in CESM LENS simulations, whereas  
8 the performance of MMEM-based methods strongly depends on how close a given  
9 type of MMEM is to the true forced signal of this model.

10

11 **4. Internal variability in observations and CMIP5 simulations: Comparison**  
12 **of different estimation methods**

13 The main problem with using MMEM-based methods for estimation of ICV in  
14 CMIP5 models, as in C2017, is illustrated in **Fig. 2**, which shows estimates of ICV in  
15 the five available historical simulations of GFDL CM3 model. Based on the results of  
16 section 3, we assume that the KC2017 method provides the best estimate of the ICV;  
17 Kravtsov et al. (2015) showed that this method indeed leads to the independent  
18 (uncorrelated) ICV estimates in different realizations of a given model. By contrast,  
19 all of the GFDL CM3's realizations of ICV estimated by using MMEM1 and  
20 MMEM2 methods clearly have the same shape (and are thus strongly correlated) in  
21 the second half of the twentieth century, where the true forced signal of the GFDL  
22 CM3 model is very different from its MMEM-based estimate (not shown). The  
23 internal variability estimated using these methods is thus contaminated (and in this  
24 case dominated) by the mismatch between the true forced signal of GFDL CM3  
25 model and the MMEM. The MMEM3 method does much better, and its estimated



1 ICV is closer to the KC2017 estimates (because the MMEM3 based forced signal is  
2 closer to that estimated using KC2017 method), but there are still considerable  
3 differences in the phase and amplitude of ICV estimated by the MMEM3 and  
4 KC2017 methods.

5 **Figure 3** summarizes the performance of the four methods discussed in this  
6 paper across the entirety of the CMIP5 sub-ensemble considered. We concentrate here  
7 on the results for the 40-yr low-pass filtered ICV, as the results utilizing 10-yr and 20-  
8 yr smoothing time scale are qualitatively similar. The results for the PMO, AMO and  
9 NMO indices are also consistent. First off, irrespective of the method used to infer  
10 either observed or simulated ICV, the magnitude of multidecadal ICV in CMIP5  
11 simulations is much lower than the observed magnitude, consistent with Frankcombe  
12 et al. (2015), Kravtsov and Callicutt (2017), Kravtsov (2017) and C2017. Among the  
13 observed estimates, KC2017 method leads to the largest estimated magnitude of the  
14 ICV, and MMEM3 method — to the smallest magnitude, while the magnitude  
15 estimated by the MMEM1 and MMEM2 methods is in between. For the simulated  
16 variability, the magnitude of the ICV inferred via MMEM1 and MMEM2 methods is  
17 the largest, and exceeds the magnitude of the variability in PI control runs by a factor  
18 of  $\sim 1.7$ – $1.8$  for PMO and NMO indices, and by a factor of 1.25 for AMO index,  
19 consistent with C2017. Both KC2017 and MMEM3 methods give magnitudes of  
20 simulated ICV consistent with the PI control runs, with KC2017-based magnitude  
21 being slightly larger, and MMEM3-based magnitude — slightly smaller than the  
22 magnitude in the PI control simulations.

23 The amplitude bias in the MMEM1 and MMEM2 based estimates of ICV in  
24 historical CMIP5 simulations stems from misidentification of the ICV illustrated in

1 Fig. 2, as the individual models' forced signals inferred using these methods are  
2 systematically different from the true forced signals in these models, which are best  
3 approximated by the KC2017 method. Of the three MMEM-based methods  
4 considered, the MMEM3 method provides the estimates of the forced signal and ICV  
5 that best match the KC2017 estimates (Figs. 3b,d,f), but even in this case, the  
6 mismatch between these two methods on multidecadal time scales (for the 40-yr low-  
7 pass filtered data) is characterized by a large rms error on the order of the standard  
8 deviation of the simulated ICV (compare brown bars in Figs. 3b,d,f with the green or  
9 brown bars in Figs. 3a,c,e, respectively).

10 In general, the MMEM-based methods are expected to provide faithful  
11 estimates of the true ICV in a given model if this model's true forced signal — well  
12 approximated by the KC2017 method — is close to the corresponding scaled  
13 MMEMs, and poor estimates otherwise. We thus developed a similarity index to  
14 characterize the individual models in terms of how well their MMEM-based estimates  
15 of the forced signal and ICV match the corresponding KC2017 estimates. To do so,  
16 we utilized 40-yr low-pass filtered MMEM1, MMEM2, MMEM3 and KC2017  
17 estimates of ICV and computed the average rms error of the three MMEM-based  
18 methods with respect to the KC2017 method for individual model sub-ensembles. For  
19 a given climate index (AMO, PMO or NMO) and a given MMEM method (MMEM1,  
20 MMEM2 or MMEM3), this procedure resulted in 17 values of the rms error (one  
21 value per model). We then assigned the consistency index value of +1 to the models  
22 characterized by the lowest rms (below the 33<sup>rd</sup> percentile of the 17 rms values), the  
23 consistency index value of -1 to the models characterized by the highest rms (above  
24 the 66<sup>th</sup> percentile of the 17 rms values) and the index value of 0 to all other models.  
25 Finally, we averaged the resulting consistency index values across the estimates

1 corresponding to the MMEM1, MMEM2 and MMEM3 methods, and then across the  
2 consistency index estimates for AMO, PMO and NMO time series. The resulting  
3 consistency index has a range between  $-1$  and  $+1$ . The value of  $-1$  corresponds to the  
4 situation when a given model is characterized by a high rms error with respect to the  
5 KC2017 method for all MMEM-based methods and for all three climate indices  
6 (hence, for such a model, the MMEM-based estimates of the forced signal and  
7 internal variability are poor). Conversely, the consistency index value of  $+1$  would  
8 indicate that the corresponding model's true forced signal is close to its MMEM-  
9 based estimate. We further clustered the models in three groups by selecting the  
10 models with consistency indices below  $-0.34$ , above  $+0.34$ , and in between, and  
11 assigning to these models the values of  $-1$  (true forced signal inconsistent with  
12 MMEM),  $+1$  (true forced signal consistent with MMEM), and  $0$  (intermediate group  
13 in terms of consistency between MMEM-based and KC2017-based forced signals).

14 The consistency index results (Table 1) indicate that MMEM-based forced  
15 signal estimation methods produce good estimates of the ICV in GISS models (except  
16 for Rp3), as well as in MRI-CGCM3, poor estimates of ICV in CanESM2,  
17 CSIROmk360, GFDL-CM3, GISSERp3, HadGEM2-ES and IPSL-CM5A-LR  
18 models, and intermediate-quality estimates in CCSM4, CNRMCM5, GFDL-CM2p1,  
19 HadCM3 and MIROC5 models. Figure 2 gives an example of the ICV in a model  
20 from the 'poor' group; see **Figs. S2 and S3** of Supplemental Materials for the  
21 analogous examples from 'good' and 'intermediate' groups. Thus, CMIP5 models are  
22 characterized by a wide variety of forced signals, which MMEM-based estimation  
23 methods fail to capture (compare with Kravtsov et al. 2015; Kravtsov and Callicutt  
24 2017; Kravtsov 2017). Therefore, application of these methods (and especially  
25 MMEM1 method) to characterize ICV in historical CMIP5 simulations — as was

1 done in C2017 — is not appropriate.

2

### 3 **5. Conclusions**

4 We compared the performance of four methods for estimating internal climate  
5 variability (ICV) in historical simulations of CMIP5 climate models. The methods  
6 considered included a suite of methods based on subtraction of scaled multi-model  
7 ensemble means (MMEM) from individual model simulations or observations  
8 (MMEM1, MMEM2, MMEM3 methods; Steinman et al. 2015a,b, Frankcombe et al.  
9 2017; Cheung et al. 2017 [C2017]), as well as the KC2017 method (Kravtsov et al.  
10 2015; Kravtsov and Callicutt 2017; Kravtsov 2017), which works with smoothed  
11 ensemble means of individual models. We found that:

- 12 1. KC2017 method provides the best estimates of the low-frequency (10-yr+)  
13 ICV in historical simulations of individual models, and results in independent  
14 (uncorrelated) time series of ICV in individual simulations of a given model.  
15 The magnitude of ICV so computed is consistent with that of the variability in  
16 the PI control runs.
- 17 2. The performance of MMEM-based methods depends on whether the scaled  
18 MMEM signal is close to this model's true forced signal or not. In the latter  
19 case, the estimated ICV is contaminated, and often dominated by the  
20 difference between the true forced signal and its MMEM-based estimate,  
21 resulting in large phase and amplitude errors of the ICV so inferred. The  
22 resulting estimates of the ICV in individual simulations of a given model are  
23 not independent (correlated). For MMEM1 method, the magnitude of the  
24 estimated ICV is larger than that in the PI control runs.

- 1        3. Out of the 17 CMIP5 models considered here, only 6 models are characterized  
2        by the forced signals well represented by the MMEM methods. For the rest of  
3        the models, the differences between the scaled MMEMs and their true forced  
4        signal estimates obtained by the KC2017 method are large, and the resulting  
5        MMEM-based estimates of the ICV are strongly biased as per item 2 above.  
6        Hence, the ICV in the ensemble of historical CMIP5 simulations cannot be  
7        reliably determined using MMEM-based methods.
- 8        4. Out of the three MMEM-based methods applied to the CMIP5 ensemble, the  
9        three-factor scaling method (MMEM3) performs, on average, better than one-  
10       factor and two-factor methods (MMEM1 and MMEM2), but is still  
11       characterized by large errors relative to the KC2017 method.
- 12       5. All methods indicate that the estimated observed ICV at decadal and longer  
13       time scales has a much larger magnitude than the CMIP5 simulated ICV. The  
14       disparity between the observed and simulated ICV is most pronounced in the  
15       KC2017 estimates of this variability.

16       Our results demonstrate that application of the MMEM1 method to infer ICV in  
17       historical CMIP5 simulations in C2017 work is inadequate, and leads to large phase  
18       and amplitude errors of ICV so estimated, which invalidates most of these authors'  
19       conclusions.

20       **Acknowledgements.** We acknowledge the World Climate Research Programme's  
21       Working Group on Coupled Modeling, which is responsible for CMIP, and we thank  
22       the climate modeling groups for producing and making available their model output.  
23       This research was supported by the NSF grant AGS-1408897. All data and MATLAB  
24       © scripts for this paper are available for downloading using the link listed in the  
25       Supplemental Materials.

1 **References**

- 2 Cheung, A. H., M. E. Mann, B. A. Steinman, L. M. Frankcombe and M. H. England,  
3 S. K. Miller, 2017: Comparison of low-frequency internal climate variability  
4 in CMIP5 models and observations. *J. Climate*, **30**, 4763–4776,  
5 doi:10.1175/JCLI-D-16-0712.1.
- 6 Compo, G. P., et al., 2011: The twentieth century reanalysis project. *Quart. J. Royal*  
7 *Meteor. Soc.*, **654**, 1–28, doi: 10.1002/qj.776.
- 8 Frankcombe, L. M., M. H. England, M. E. Mann, and B. A. Steinman, 2015:  
9 Separating internal variability from the externally forced climate response. *J.*  
10 *Climate*, **28**, 8184–8202, doi: <http://dx.doi.org/10.1175/JCLI-D-15-0069.1>.
- 11 Kay, J. E., et al., 2015: The Community Earth System Model (CESM) Large  
12 Ensemble Project: A community resource for studying climate change in the  
13 presence of internal climate variability. *Bull. Amer. Meteor. Soc.*, **96**, 1333–  
14 1349.
- 15 Kravtsov, S., 2017: Pronounced differences between observed and CMIP5-simulated  
16 multidecadal climate variability in the twentieth century. *Geophys. Res. Lett.*,  
17 published online 15 June 2017, doi: 10.1002/2017GL074016,  
18 <http://onlinelibrary.wiley.com/doi/10.1002/2017GL074016/full>.
- 19 Kravtsov, S., and D. Callicutt, 2017: On semi-empirical decomposition of  
20 multidecadal climate variability into forced and internally generated  
21 components. *Internat. J. Climatol.*, published online April 2017, doi:  
22 10.1002/joc.5096, <http://onlinelibrary.wiley.com/doi/10.1002/joc.5096/full>.
- 23 Kravtsov, S., M. G. Wyatt, J. A. Curry, and A. A. Tsonis, 2015: Comment on  
24 “Atlantic and Pacific multidecadal oscillations and Northern Hemisphere  
25 temperatures.” *Science*, **350**, 1326, doi: 10.1126/science.aab3570.

- 1 Mann, M. E., 2008: Smoothing of climate time series revisited. *Geophys. Res. Lett.*,  
2 **35**, L16708, doi:10.1029/2008GL034716.
- 3 Steinman, B. A., M. E. Mann and S. K. Miller, 2015a: Atlantic and Pacific  
4 multidecadal oscillations and Northern Hemisphere temperatures. *Science*,  
5 **347**, 988, doi: 10.1126/science.1257856.
- 6 Steinman, B. A., L. M. Frankcombe, M. E., Mann, S. K. Miller, and M. H. England,  
7 2015b: Response to comment on “Atlantic and Pacific multidecadal  
8 oscillations and Northern Hemisphere temperatures.” *Science*, **350**, 1326, doi:  
9 10.1126/science.aac5208.
- 10 Taylor, K. E., R. J. Stouffer and G. A. Meehl, 2012: An Overview of CMIP5 and the  
11 experiment design. *Bull. Am. Meteorol. Soc.*, **93**, 485–498.
- 12

1 **Table captions**

2 **Table 1.** Classification of the CMIP5 models considered in terms of consistency  
3 between their internal variability estimates computed using MMEM-based  
4 methods and KC2017 method. The first two columns give model acronyms  
5 and the number of historical simulations for a given model. The consistency  
6 indices themselves — defined to range from  $-1$  (low consistency between the  
7 two methods) to  $+1$  (high consistency between the two methods) — are listed  
8 in the last two columns, with the first of these columns showing the raw index,  
9 and the second — its ternary version. See text for details.



1 **Figure captions**

2 **Figure 1:** Performance of different methods for isolating internal variability in 40  
3 historical simulations (1920–2005) of the Community Earth System Model  
4 (CESM) Large Ensemble Project (LENS). The methods are one-factor  
5 (MMEM1), two-factor (MMEM2) and three-factor (MMEM3) methods based  
6 on subtraction of scaled CMIP5 multi-model ensemble-mean (MMEM)  
7 forced-signal estimates from individual LENS simulations, as well as KC2017  
8 method, applied here to random 5-member sub-ensembles of LENS  
9 simulations. Top row: standard deviation of the raw and low-pass-filtered  
10 estimates of internal variability, as a function of smoothing time scale. In  
11 addition to estimates for all methods above, also included here are the  
12 estimates for the “true” internal variability based on subtracting LENS grand  
13 ensemble mean from individual simulations, as well as estimates based on  
14 random 86-yr-long snippets from CESM pre-industrial (PI) control run (see  
15 panel legends). Bottom row: root-mean-square (rms) errors of internal  
16 variability estimated using the four methods with respect to the “true” internal  
17 variability (see panel legends), again as a function of smoothing time scale.  
18 Bars show the mean value and error bars — the standard deviation of each  
19 quantity across the 40-member LENS ensemble. (a, b) Results for PMO index;  
20 (c, d) results for AMO index; (e, f) results for NMO index.

21 **Figure 2:** Estimates of internal variability in historical simulations of the GFDL CM3  
22 model using (from top to bottom) MMEM1, MMEM2, MMEM3, KC2017  
23 methods; also shown are five random samples of internal variability from the  
24 PI control run (bottom row). All time series were 40-yr low-pass filtered using

1 data-adaptive filter of Mann (2008) and windowed using tapers from Kravtsov  
2 and Callicutt (2017) to reduce spurious end effects. Left column: results for  
3 AMO; middle column: results for PMO; right column: results for NMO.

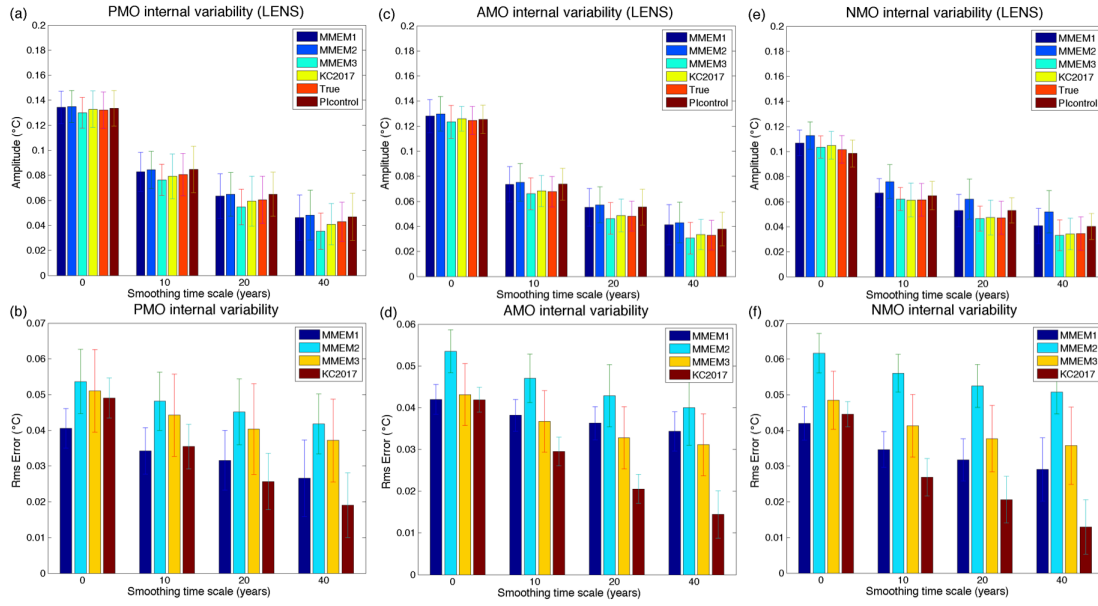
4 **Figure 3:** Internal variability in observations and CMIP5 simulations. Top row:  
5 magnitude (standard deviation) of the observed and simulated internal  
6 variability estimated using four different methods, as well as estimates based  
7 on random 126-yr-long snippets from PI control runs (see panel legends), as a  
8 function of smoothing time scale. Bottom row: rms errors of internal  
9 variability estimated using various MEM methods with respect to the  
10 KC2017 estimate of the internal variability (see panel legends), again as a  
11 function of smoothing time scale. Bars show the mean value and error bars —  
12 the standard deviation of each quantity across the entire ensemble of  
13 simulations. (a, b) Results for PMO index; (c, d) results for AMO index; (e, f)  
14 results for NMO index.

15

1 **Table 1.** Classification of the CMIP5 models considered in terms of consistency  
2 between their internal variability estimates computed using MMEM-based  
3 methods and KC2017 method. The first two columns give model acronyms  
4 and the number of historical simulations for a given model. The consistency  
5 indices themselves — defined to range from  $-1$  (low consistency between the  
6 two methods) to  $+1$  (high consistency between the two methods) — are listed  
7 in the last two columns, with the first of these columns showing the raw index,  
8 and the second — its ternary version. See text for details.  
9

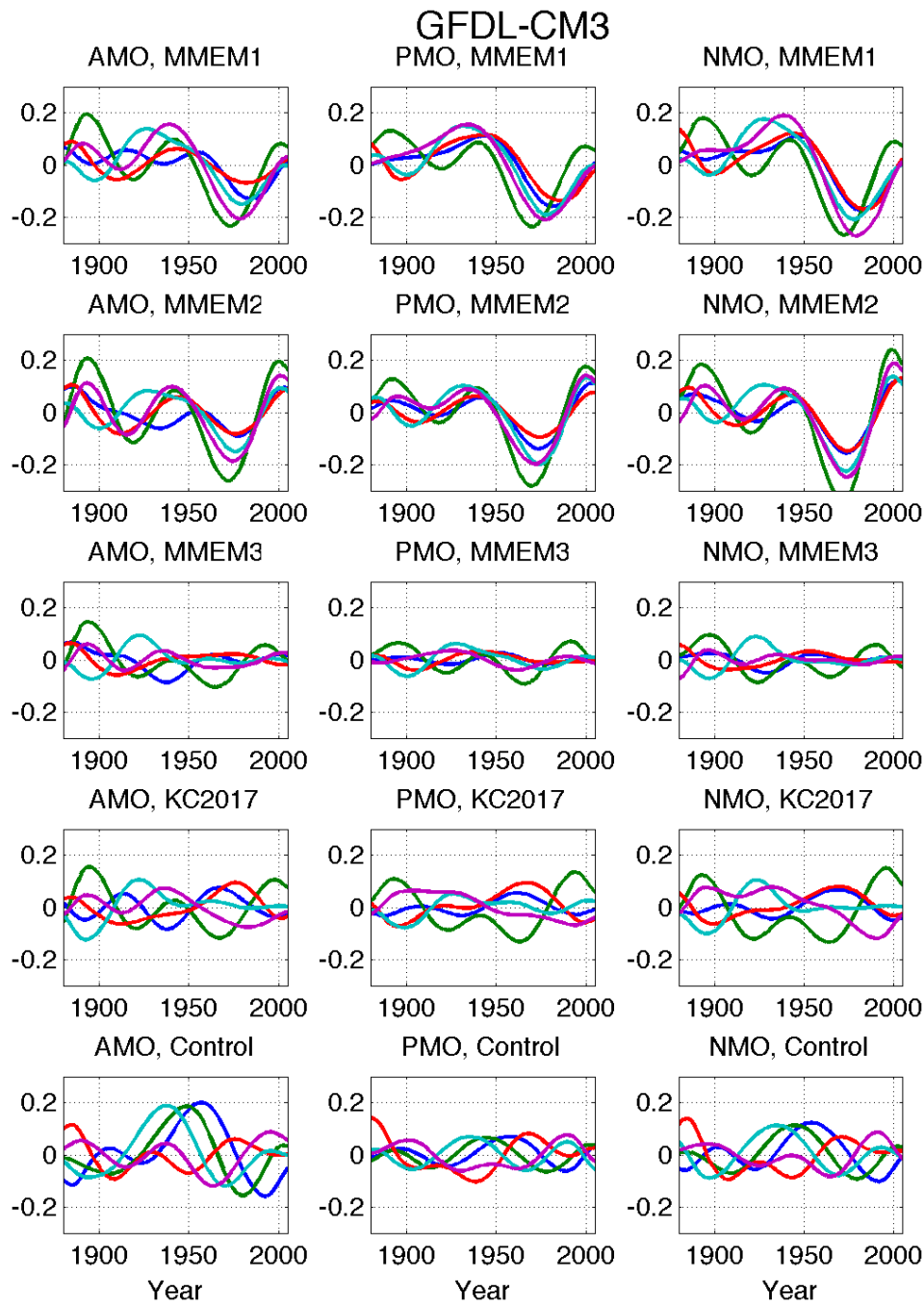
Model acronym	Number of realizations	Consistency between MMEM and KC17 methods	
		Raw index	Ternary version
CanESM2	5	-0.44	-1
CCSM4	6	0	0
CNRM-CM5	10	0.22	0
CSIRO-Mk3-6-0	10	-1	-1
GFDL-CM2p1	10	-0.33	0
GFDL-CM3	5	-0.78	-1
GISS-E2-Hp1	6	0.67	1
GISS-E2-Hp2	6	0.89	1
GISS-E2-Hp3	6	0.89	1
GISS-E2-Rp1	6	1	1
GISS-E2-Rp2	6	0.67	1
GISS-E2-Rp3	6	-0.44	-1
HadCM3	10	-0.33	0
HadGEM2-ES	5	-0.55	-1
IPSL-CM5A-LR	6	-0.78	-1
MIROC5	5	-0.11	0
MRI-CGCM3	3	0.44	1

10



1 **Figure 1:** Performance of different methods for isolating internal variability in 40  
 2 historical simulations (1920–2005) of the Community Earth System Model  
 3 (CESM) Large Ensemble Project (LENS). The methods are one-factor  
 4 (MMEM1), two-factor (MMEM2) and three-factor (MMEM3) methods based  
 5 on subtraction of scaled CMIP5 multi-model ensemble-mean forced-signal  
 6 estimates from individual LENS simulations, as well as KC2017 method,  
 7 applied here to random 5-member sub-ensembles of LENS simulations. Top  
 8 row: standard deviation of the raw and low-pass-filtered estimates of internal  
 9 variability, as a function of smoothing time scale. In addition to estimates for  
 10 all methods above, also included here are the estimates for the “true” internal  
 11 variability based on subtracting LENS grand ensemble mean from individual  
 12 simulations, as well as estimates based on random 86-yr-long snippets from  
 13 CESM pre-industrial (PI) control run (see panel legends). Bottom row: root-  
 14 mean-square (rms) errors of internal variability estimated using the four  
 15 methods with respect to the “true” internal variability (see panel legends),  
 16 again as a function of smoothing time scale. Bars show the mean value and  
 17 error bars — the standard deviation of each quantity across the 40-member  
 18 LENS ensemble. (a, b) Results for PMO index; (c, d) results for AMO index;  
 19 (e, f) results for NMO index.

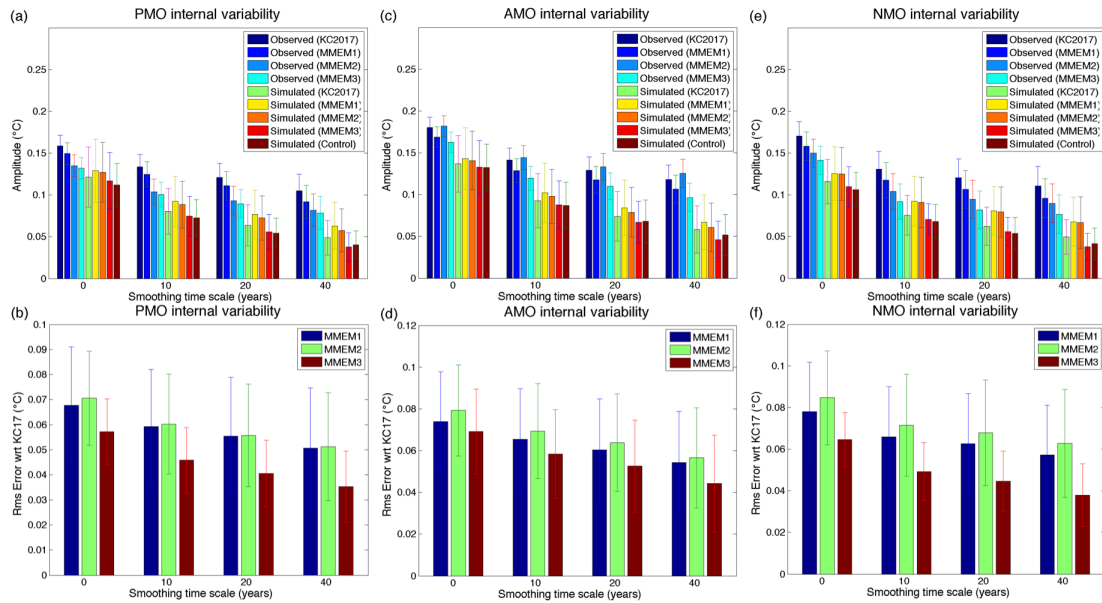
20  
21



1

2 **Figure 2:** Estimates of internal variability in historical simulations of the GFDL CM3  
 3 model using (from top to bottom) MMEM1, MMEM2, MMEM3, KC2017  
 4 methods; also shown are five random samples of internal variability from the  
 5 PI control run (bottom row). All time series were 40-yr low-pass filtered using  
 6 data-adaptive filter of Mann (2008) and windowed using tapers from Kravtsov  
 7 and Callicutt (2017) to reduce spurious end effects. Left column: results for  
 8 AMO; middle column: results for PMO; right column: results for NMO.

9



1  
 2  
 3 **Figure 3:** Internal variability in observations and CMIP5 simulations. Top row:  
 4 magnitude (standard deviation) of the observed and simulated internal  
 5 variability estimated using four different methods, as well as estimates based  
 6 on random 126-yr-long snippets from PI control runs (see panel legends), as a  
 7 function of smoothing time scale. Bottom row: rms errors of internal  
 8 variability estimated using various MMEM methods with respect to the  
 9 KC2017 estimate of the internal variability (see panel legends), again as a  
 10 function of smoothing time scale. Bars show the mean value and error bars —  
 11 the standard deviation of each quantity across the entire ensemble of  
 12 simulations. (a, b) Results for PMO index; (c, d) results for AMO index; (e, f)  
 13 results for NMO index.  
 14