

A Representative Democracy to Reduce Interdependency in a Multimodel Ensemble

BENJAMIN M. SANDERSON

National Center for Atmospheric Research, Boulder, Colorado*

RETO KNUTTI

Institute for Atmospheric and Climate Science, ETH, Zurich, Switzerland

PETER CALDWELL

Lawrence Livermore National Laboratory, Livermore, California

(Manuscript received 22 May 2014, in final form 6 March 2015)

ABSTRACT

The collection of Earth system models available in the archive of phase 5 of CMIP (CMIP5) represents, at least to some degree, a sample of uncertainty of future climate evolution. The presence of duplicated code as well as shared forcing and validation data in the multiple models in the archive raises at least three potential problems: biases in the mean and variance, the overestimation of sample size, and the potential for spurious correlations to emerge in the archive because of model replication. Analytical evidence is presented to demonstrate that the distribution of models in the CMIP5 archive is not consistent with a random sample, and a weighting scheme is proposed to reduce some aspects of model codependency in the ensemble. A method is proposed for selecting diverse and skillful subsets of models in the archive, which could be used for impact studies in cases where physically consistent joint projections of multiple variables (and their temporal and spatial characteristics) are required.

1. Introduction

Today's Earth system models (ESMs) are great testament to collaborative scientific thinking. Millions of lines of computer code represent the pinnacle of understanding of the intricate coupled interactions of Earth's land, ocean, cryosphere, and atmosphere systems. Unlike the more simple atmospheric models of the past, few people (if any) now understand the models in their entirety and so the models themselves have become vehicles of a scientific consensus that we use to

project future climates and cannot be directly validated for decades to come. For some parts, such as the representation of the equations of fluid flow, understanding is mature and thus (relatively) uncontentious. However, other components, such as the effect of a changing climate on ecosystem dynamics, are sufficiently complex that any computational code must inevitably make significant approximations in order to even represent the bulk behavior of the system in any tractable fashion.

A given model is thus more than a computer program; it is a collection of axioms and beliefs about which processes might be important for evaluating how our environment might change and how those processes should be represented, and as such each model is a self-consistent entity. The challenge arises, however, when one wishes to combine the results of many models to attain some more comprehensive understanding of the uncertainties present in their individual implementation. Given a set of models of the climate system, assessing the value of adding another model clearly requires a consideration of whether the model is fit for purpose (e.g., the validity of its axioms, forcing data, and

 Denotes Open Access content.

*The National Center for Atmospheric Research is sponsored by the National Science Foundation.

Corresponding author address: Benjamin Sanderson, National Center for Atmospheric Research, 1850 Table Mesa Dr., Boulder, CO 80305.
E-mail: bsander@ucar.edu

tuning protocols). We would argue also that it is important to assess if the model provides new information: to measure how independent the new model is from those in the original set. In an extreme case, adding an exact duplicate of a model already in the set would not add value; rather, it would bias any combination of model results toward the results of the duplicated model (Caldwell et al. 2014).

Phase 5 of the Coupled Model Intercomparison Project (CMIP5; Taylor et al. 2012) is the largest archive of climate data the world has seen to date. Such multi-model ensembles (MMEs) have often been referred to as “ensembles of opportunity” (Tebaldi and Knutti 2007) because the range of models represent some sample of the systematic choices that developers face in the course of representing the climate system in the form of computer code. However, as has been noted before (Knutti 2010), this sample is far from perfect.

First, the models available may vary in their ability to resolve certain processes that might be observed in the Earth system. For any given process, a researcher may find relevant observations to rank models for their purposes, but the output of the ESMs is sufficiently high dimensional that any ranking is unlikely to be universal (Santer et al. 2009). In contrast to weather forecast models, ESMs can also rarely be validated out of sample and so there remains a risk that empirical components of ESMs can be calibrated using the only available observations; although this might be a pragmatic approach, it leaves little opportunity for assessing and contrasting model performance (Sanderson and Knutti 2012).

A second problem lies in the lack of independence of models, where independence is not meant in a statistical sense but in a more loose sense of models sharing ideas for parameterizations and simplifications or sharing actual computer code and therefore being biased in similar ways relative to reality. At the time of writing, 61 models are listed in the Earth System Grid database. This does not necessarily mean that each of these models provides an independent estimate of future climate change. Indeed, some of these codedependencies are trivial and can be accounted for by considering models submitted with different resolutions (e.g., MPI-ESM-MR and MPI-ESM-LR; see Knutti et al. 2013). Most institutions also produce model variants with a range of different configurations, with options for interactive atmospheric chemistry or carbon cycle (e.g., CMCC-CESM and CMCC-CM). Finally, different institutions can share model components: for example, the FIO-ESM model shares its atmosphere, ocean, sea ice, and land surface code with CCSM4 but adds a surface ocean wave parameterization. Submodel replication is common throughout the ensemble: for example, in the models considered for

this study over 25% use some variant of the Community Atmosphere Model (CAM3, CAM3.5, CAM4, or CAM5) to represent atmospheric processes. The GFDL Modular Ocean Model is similarly popular (MOM2.2, MOM4.0, and MOM4.1). Table 1 shows a broad illustration of shared model components in the CMIP5 models considered for this study.

This extensive model replication in the CMIP5 and its predecessors is not a problem per se; in fact, it seems natural to copy successful parts and build on the work of others, and it requires enormous effort to develop entirely new model components. Hence, each institution understandably focuses on certain aspects but copies other components. However, model replication presents a number of issues for model ensemble analysis. The first is simply a matter of representation: the assessment reports of the Intergovernmental Panel on Climate Change (IPCC) have often used the multimodel mean of the CMIP ensembles to represent a consensus view of model projections of future climate, but clearly this mean will be biased toward models that are highly replicated within the ensemble. Similarly, model agreement on the sign or magnitude of a change in future climate is often taken to imply confidence in a result (Tebaldi et al. 2011; Knutti and Sedáček 2013) but, if models are highly replicated within the ensemble, such agreement becomes less significant.

Another issue lies in the possible effect of replicated models in studies that attempt to constrain aspects of future climate change. If a researcher discovers a correlation between an observable quantity and some unknown climate parameter in a multimodel ensemble (e.g., Fasullo and Trenberth 2012; Qu and Hall 2013), the statistical significance of that correlation would be inflated if some points are repeated. This argument is developed in Caldwell et al. (2014), who show that, although a data-mining approach will yield more strong correlations between climate sensitivity and potentially observable fields than one would expect to see by chance in CMIP5, this may be attributable in part to model codedependencies.

This is the second in a series of papers examining interdependency in the CMIP ensembles. In Sanderson et al. (2015), we developed a distance metric that enabled both models and observations to be represented as points in a multidimensional space. We then showed that model properties could be interpolated within this space, allowing a resampling of model properties in a manner that was less sensitive to model replication and could take into account a measure of performance in reproducing observations. However, the approach of Sanderson et al. (2015) is also unable to provide full spatial and temporal variations in quantities. For example,

a farmer may not want an estimate of the change in average rainfall but a set of representative summers with full spatial and temporal information and the corresponding temperature, sunshine, and wind data. For such cases, it may be better to use the raw or bias-corrected model output directly, but that requires selecting a set of models to use.

It has been proposed before that subsets of larger ensembles may produce more statistically robust results; [Evans et al. \(2013\)](#) investigated this concept using subsets of a multiphysics ensemble of weather forecasting models. Perhaps the simplest approach to achieve this might be to take a single model from each institution, but there are numerous issues with this. First, although there are often similarities between models published by single institutions, such a crude approach would eliminate cases where significantly different models were produced by the same group. There are several examples of the latter case: the GISS-E2 models, for example, are published with two structurally different oceans. Furthermore, several groups [National Science Foundation (NSF)–DOE–NCAR (CESM), GFDL, and Met Office (UKMO), among others] publish both a “bleeding edge” model and a legacy model to the archive, where there might be significant structural changes between the releases. Finally, an institution-based pruning approach would not help identify models from different institutions that share a large fraction of their code.

It could be argued that one could account for many of these problems through careful consideration of model lineages, by documenting the basic parameterizations shared by different models or by assessing the fraction of common code between different models. This, however, would be a considerable undertaking, and the results would require a comprehensive understanding of each model’s code. First, although some models document and publish their code base in full before submitting simulations to the CMIP archive, this practice is far from universal. A model could in theory be defined by summarizing the parameterizations, their values, and other structural assumptions that have been employed in that model, but assessing the relative importance of each of those parameterizations in terms of model climatology or response to external forcing would require good prior intuition of the relationships between the parameterizations and the process to be studied, which might be possible in some but not necessarily all cases. Such an approach would clearly be worthwhile, and could greatly aid in the interpretation of differences in climate change projections, but it would be a monumental undertaking.

An alternative approach is to utilize output from the models themselves to establish codependencies. This approach has been demonstrated with some promise by

[Masson and Knutti \(2011, 2013\)](#), who used intermodel distances derived from spatial patterns of climatological temperature and precipitation to establish a hierarchical clustering of models that resembles a tree showing structural relationships one might expect from considering model lineages. As noted in [Masson and Knutti \(2011\)](#) and [Sanderson et al. \(2015\)](#), the distribution of intermodel distances shows recognizable structure, with models from the same institution and models with common heritage generally exhibiting similar patterns of mean state bias. However, the aforementioned studies did not establish any quantitative assessment of intermodel distance, which we attempt to address here.

To this end, we formalize an approach to use model similarity information to select models based on their skill and independence. This does not eliminate model interdependency but allows us to select a subset of models where the most glaring examples of model replication are no longer present. In [section 2a](#), we establish a method for identifying near neighbors in a climate model ensemble. In [section 2d](#), we use model similarity information to produce a weighting scheme that accounts for both model skill and model interdependence. [Section 2e](#) shows how this framework can be used to select a subset of models from an archive of climate models. Finally, [section 3b](#) demonstrates this method using the CMIP5 multimodel archive.

2. Method

a. Processing model output

In this study, as in our accompanying paper ([Sanderson et al. 2015](#)), we produce a matrix of intermodel distances in an EOF space derived from 30-yr mean climatological output from each model’s historical simulation conducted for CMIP5. The details of the construction of the distance matrix are identical to that of [Sanderson et al. \(2015\)](#). We use the historical and RCP8.5 experiments and the CMIP5 ensemble-member simulations in each case. In the special case of CCSM4, we also consider the sensitivity of the technique to internal variability by repeating the analysis with all available simulations in the CMIP5 archive (e.g., r1i1p1, r1i2p1, r1i2p2, r2i1p1, r3i1p1, r4i1p1, r5i1p1, and r6i1p1 for the historical runs and r1i1p1, r2i1p1, r3i1p1, r4i1p1, r5i1p1, and r6i1p1 for the RCP8.5 simulations).

The input data for this study are both processed and used to conduct an EOF analysis in a similar fashion to [Sanderson et al. \(2015\)](#). Minor differences in the intermodel distances occur because the former study considers both CMIP3 and CMIP5 models, which slightly changes the exact form of the EOFs. For each model, a

TABLE 1. Submodel components for the 38 CMIP5 models considered in this study. Expansions of common model, model component, and dataset acronyms are available at <http://www.ametsoc.org/Pubs/AcronymList>.

Model	Atmosphere	Land	Ocean	Ice	Source
NorESM1-ME	CAM4	CLM4	Miami Isopycnic Coordinate Ocean Model (MICOM)–Hamburg Model of the Ocean Carbon Cycle (HAMOCC)	CICE	https://verc.enes.org/ISENES2/models/earthsystem-models/ncc/noresm
NorESM1-M	CAM4	CLM4	MICOM–HAMOCC	CICE	https://verc.enes.org/ISENES2/models/earthsystem-models/ncc/noresm
MRI-CGCM3	MRI-AGCM3	HAL	MRI-COM3	—	http://www.mri-jma.go.jp/Publish/Technical/DATA/VOL_64/index_en.html
MPI-ESM-MR	ECHAM6	JSBACH	MPI-OM	—	http://www.mpimet.mpg.de/en/science/models/mpie-sm.html
MPI-ESM-LR	ECHAM6	JSBACH	MPI-OM	—	https://verc.enes.org/models/earthsystem-models/mpie-m/mpie-esm
MIROC5	FRCGC-AGCM	Minimal Advanced Treatments of Surface Interaction and Runoff (MATSIRO)	COCO	Bitz–Lipscomb	Watanabe et al. (2010)
MIROC-ESM-CHEM	FRCGC-AGCM		COCO	Bitz–Lipscomb	http://www.wcrp-climate.org/wgem/WGCM15/presentations/21Oct/KIMOTO_Japan.pdf
MIROC-ESM	FRCGC-AGCM	MATSIRO	COCO	Bitz–Lipscomb	http://www.wcrp-climate.org/wgem/WGCM15/presentations/21Oct/KIMOTO_Japan.pdf
IPSL-CM5B-LR	LMDZ (CM4)	ORCHIDE	NEMO-OPA	NEMO-LIM	http://icmc.ipsl.fr/index.php/icmc-models/icmc-ipsl-cm5
IPSL-CM5A-MR	LMDZ	ORCHIDE	NEMO-OPA	NEMO-LIM	http://icmc.ipsl.fr/index.php/icmc-models/icmc-ipsl-cm5
IPSL-CM5A-LR	LMDZ	ORCHIDE	NEMO-OPA	NEMO-LIM	http://icmc.ipsl.fr/index.php/icmc-models/icmc-ipsl-cm5
INM-CM4.0	INM-CM	INM-CM	INM-CM	INM-CM	Volodin et al. (2010)
FGOALS-g2	Gridpoint Atmospheric Model of IAP–LSAG, version 2.0 (GAMIL 2.0)	CLM3	LICOM2	CICE4 (LASG)	Li et al. (2013)
HadGEM2-ES	HadGAM2 (N96L38)	Top-down Representation of Interactive Foliage and Flora Including Dynamics (TRIFFID)	HadGOM2	—	http://cms.ncas.ac.uk/wiki/UM/Configurations/HadGEM2

TABLE 1. (Continued)

Model	Atmosphere	Land	Ocean	Ice	Source
HadGEM2-CC	HadGAM2(N96L60)	TRIFFID	HadGOM2	—	http://cms.ncas.ac.uk/wiki/UM/Configurations/HadGEM2
HadGEM2-AO	HadGAM2 (N96L38)	MOSES2	HadGOM2	—	http://cms.ncas.ac.uk/wiki/UM/Configurations/HadGEM2
GISS-E2-R	GISS	GISS	Russell ocean model	Russell ocean model	http://data.giss.nasa.gov/modelE/ar5/
GISS-E2-H	GISS	GISS	HYCOM	HYCOM	http://data.giss.nasa.gov/modelE/ar5/
GFDL-ESM2M	GFDL AM2.1	LM3	MOM4.1	Sea Ice Simulator (SIS)	http://cms.ncas.ac.uk/wiki/UM/Configurations/HadGEM2
GFDL-ESM2G	GFDL AM2.1	LM3	GOLD	SIS	http://www.gfdl.noaa.gov/earth-system-model
GFDL CM3	GFDL AM3	LM3	MOM4.1	SIS	http://www.gfdl.noaa.gov/earth-system-model
FIO-ESM	CAM3.5	CLM3	POP2	CICE4	http://www.wcrp-climate.org/wgcm/WGCM15/presentations/21Oct/WANG_WGCM.pdf
CanESM2	AGCM4	Canadian Land Surface Scheme (CLASS)	NCAR	—	Yang and Saenko (2012)
CSIRO Mk3.6.0	Gordon	CSIRO Atmosphere-Biosphere Land Exchange (CABLE)	MOM2.2	SIS	Jeffrey et al. (2013)
CNRM-CM5	ARPEGE-Climate	ISBA	NEMO-OPA	Global Experimental Leads and Sea Ice for Atmosphere and Ocean (GELATO)	http://www.cnrm-game.fr/spip.php?article126&lang=en
CMCC-CMS	ECHAM5	Surface Interactive Land Vegetation (SILVA)	OPA8.2	LIM	http://www.wcrp-climate.org/wgcm/WGCM16/Bellucci_CMCC.pdf
CMCC-CM	ECHAM5	SILVA	OPA8.2	LIM	http://www.cmcc.it/models/cmcc-cm
CMCC-CESM	ECHAM5	SILVA	OPA8.2	LIM	http://www.cmcc.it/models/cmcc-cm
CESM1(CAM5)	CAM5	CLM4	POP2	CICE4	https://www2.cesm.ucar.edu/models
CESM1(BGC)	CAM4	CLM4	POP2	CICE4	https://www2.cesm.ucar.edu/models
CCSM4	CAM4	CLM4	POP2	CICE4	https://www2.cesm.ucar.edu/models
BNU-ESM	CAM3.5	CLM (BNU)	MOM4.1	CICE4.1	http://www.wcrp-climate.org/wgcm/WGCM15/presentations/21Oct/WANG_WGCM.pdf
BCC_CSM1.1(m)	BCC_AGCM2.1	CLM3	MOM4	SIS	Wu et al. (2014)
BCC_CSM1.1	BCC_AGCM2.1	CLM3	MOM4	GFDL SIS	Wu et al. (2014)
ACCESS1.3	UKMO GA1.0	CABLE version 1.8	MOM4.1	CICE4.1	https://wiki.csiro.au/display/ACCESS/Home
ACCESS1.0	HadGEM2 r1.1	MOSES	MOM4.1	CICE4.1	http://www.cawcr.gov.au/publications/technicalreports/CTR_059.pdf

TABLE 2. Observational datasets used as observations throughout the study.

Field	Source	Reference	Years	Global normalization
TAS	HadCRUT3	Brohan et al. (2006)	1970–2000	2.09 K
PR	GPCP	Adler et al. (2003)	1979–2001	30.1 W m ⁻²
RSUT	CERES EBAF	NASA (2011)	2000–05	25.8 W m ⁻²
RLUT	CERES EBAF	NASA (2011)	2000–05	3.32 mm day ⁻¹
<i>T</i>	Atmospheric Infrared Sounder (AIRS)*	Aumann et al. (2003)	2002–10	0.28 K
RH	AIRS*	Aumann et al. (2003)	2002–10	12.12%

* The data used in this effort were acquired as part of the activities of NASA's Science Mission Directorate and are archived and distributed by the Goddard Earth Sciences (GES) Data and Information Services Center (DISC).

number of monthly, gridded diagnostic variables are considered to represent the climatology of the model. For each available model in the CMIP3 and CMIP5 ensembles, monthly climatologies are obtained from a single historical simulation by averaging monthly mean fields for the time period 1970–2000. Data are obtained for five two-dimensional fields [surface air temperature (TAS), total precipitation (PR), outgoing top-of-atmosphere shortwave radiative flux (RSUT), outgoing longwave top-of-atmosphere flux (RLUT), and sea level pressure (PSL)] and two three-dimensional fields [atmospheric temperature (*T*) and relative humidity (RH)]. Three-dimensional fields are zonally averaged. Corresponding observational monthly mean climatologies are obtained by averaging available years for each field type, as shown in Table 2.

Data from each model and dataset are regridded onto a 2.5° by 3.75° latitude–longitude grid, and zonal vertical fields are regridded onto a 2.5° latitude grid at 17 pressure levels. For each variable, values are area weighted. Vertically resolved fields are also weighted by the pressure difference between the top and bottom of the corresponding level. To usefully concatenate the multivariate field for EOF analysis, the variables must be normalized for each to represent a similar amount of variance in the multimodel ensemble. We normalize each observable field using values obtained from the observations. For two-dimensional fields, we calculate the intermonthly variance of tropical grid cells and take the average over the tropics to obtain a single normalization factor for each variable. For three-dimensional fields, we take the intermonthly variance of zonally averaged fields in the tropics between 700 and 400 hPa and then average the variances over the spatial domain to obtain the normalization factor. Normalization factors are calculated from the observations only, and the fields from each model are divided by the same factor (shown in Table 2). Each field is then reformulated into a single vector. If any elements of the vector in any single model or in the observations are missing, those particular elements are removed from all models. Each field vector is

then normalized by the number of remaining elements, and the 2D and 3D fields are concatenated into a single vector length n (where $n = 358\,248$ when all fields are utilized). Each of the m vectors is combined to form a matrix \mathbf{X}^{20c} (size m by n , where m is 36, comprising 36 CMIP5 model vectors). The ensemble mean value is calculated by averaging the m rows of the matrix, and this is subtracted from each row to yield the anomaly matrix $\Delta\mathbf{X}^{20c}$, such that

$$\Delta\mathbf{X}^{20c} = \mathbf{X}^{20c} - \overline{\mathbf{X}^{20c}}. \quad (1)$$

The analysis is also repeated with a number of different subsets of the entire set of variables. In these cases, the matrix $\Delta\mathbf{X}^{20c}$ is formed using only that subset, and the analysis continues in the same fashion.

The process is repeated to produce a similar matrix to represent the climate change between the historical simulation (1970–2000) and the RCP8.5 simulation (2070–2100). In this second analysis, the anomaly between the two 30-yr periods is taken to form the matrix $\Delta\mathbf{X}^{21c}$. The future analysis is also repeated with a number of different subsets of the entire set of variables. In these cases, the matrix $\Delta\mathbf{X}^{21c}$ is formed using only that subset, and the analysis continues in the same fashion.

b. Principal component analysis

We conduct a principal component analysis on the resulting matrix formed by combining the climatology vectors from each participating model, such that the EOF loadings define a t -dimensional space (where t is the truncation length of the principal component analysis) in which intermodel and observation–model Euclidean distances may be defined. The use of the EOF prefilter combines fields that are trivially correlated (i.e., adjacent grid cells) into a single mode. The results of the analysis do change in a subtle fashion with truncation length, and we discuss this sensitivity further in the subsection in section 3c, but for the initial analysis we use a truncation length of $t = 9$. This truncation length effectively provides enough degrees of freedom to

represent some subtle differences between related models in the resulting distance metric but not so many as to introduce excessive random noise into the calculation.

The PCA on any $\Delta\mathbf{X}$ can be performed by singular value decomposition and truncated to t modes, such that

$$\Delta\mathbf{X}^{20c} = \mathbf{U}^{20c} \boldsymbol{\lambda}^{20c} (\mathbf{V}^{20c})^T \quad (2)$$

for the present-day case (20c) and

$$\Delta\mathbf{X}^{21c} = \mathbf{U}^{21c} \boldsymbol{\lambda}^{21c} (\mathbf{V}^{21c})^T \quad (3)$$

for the future case (21c). The \mathbf{U}^{20c} and \mathbf{U}^{21c} (sized m by t) are matrices of model loadings, \mathbf{V}^{20c} and \mathbf{V}^{21c} (sized n by t) are spatial patterns of ensemble variability, and $\boldsymbol{\lambda}^{20c}$ and $\boldsymbol{\lambda}^{21c}$ (sized t by t) are diagonal matrices representing the variances associated with each mode.

The intermodel distances can then be measured in a Euclidean sense in the loadings matrices \mathbf{U}^{20c} and \mathbf{U}^{21c} , such that the distances between two models i and j can be expressed as

$$\delta_{ij}^{20c} = \left\{ \sum_{l=1}^t [\mathbf{U}^{20c}(i, l) - \mathbf{U}^{20c}(j, l)]^2 \right\}^{1/2} \quad (4)$$

for the present day and

$$\delta_{ij}^{21c} = \left\{ \sum_{l=1}^t [\mathbf{U}^{21c}(i, l) - \mathbf{U}^{21c}(j, l)]^2 \right\}^{1/2} \quad (5)$$

for the future. Model–observation distances $\delta_{i(\text{obs})}^{20c}$ that can obviously only be calculated for the present-day case are created using a climatological vector from an observational dataset $\mathbf{X}_{(\text{obs})}$ prepared in the same fashion as \mathbf{X}^{20c} ,

$$\Delta\mathbf{X}_{(\text{obs})n}^{20c} = \mathbf{X}_{(\text{obs})} - \overline{\mathbf{X}^{20c}}, \quad (6)$$

where $\overline{\mathbf{X}^{20c}}$ is the multimodel mean of \mathbf{X}^{20c} , with length n . This observational anomaly vector can be projected onto \mathbf{V}^{20c} to form an observational loading vector $\mathbf{U}_{(\text{obs})}$ (length t). The distance between each model and the observations can be then calculated in a similar fashion,

$$\delta_{i(\text{obs})}^{20c} = \left\{ \sum_{l=1}^t [\mathbf{U}^{20c}(i, l) - \mathbf{U}_{(\text{obs})}^{20c}(l)]^2 \right\}^{1/2}. \quad (7)$$

Finally, we calculate the variability expected in an initial condition (ic) ensemble by taking $n_{\text{ic}} = 8$ (historical) or $n_{\text{ic}} = 6$ (future) member CCSM4 ensemble for both the historical simulation and RCP8.5. In each case,

the data are processed in the same fashion as for the multimodel case to create an n_{ic} by n matrix, $\mathbf{X}_{\text{ic}}^{20c}$ and $\mathbf{X}_{\text{ic}}^{21c}$. We then take anomalies from the CMIP5 ensemble mean,

$$\Delta\mathbf{X}_{\text{ic}}^{20c} = \mathbf{X}_{\text{ic}}^{20c} - \overline{\mathbf{X}^{20c}} \quad \text{and} \quad (8)$$

$$\Delta\mathbf{X}_{\text{ic}}^{21c} = \mathbf{X}_{\text{ic}}^{21c} - \overline{\mathbf{X}^{21c}}. \quad (9)$$

These can also be projected onto \mathbf{V}^{20c} and \mathbf{V}^{20c} to form loading vectors $\mathbf{U}_{(\text{ic})}^{20c}$ and $\mathbf{U}_{(\text{ic})}^{21c}$ (size n_{ic} by t). The distance between initial condition ensemble members can be then calculated as before for the multimodel case.

c. Forming random ensembles

To compare intermodel distances in the CMIP5 archive with distances expected by chance, we create a set of 10^5 matrices of random data with the same dimensions as \mathbf{U}^{20c} and \mathbf{U}^{21c} (where m is 36). Each random distribution represents interpoint distances for all possible pairwise combinations m points (703 distances in this case). Our results are not sensitive to further increasing the number of random cases.

Each row of one of these random matrices is populated with draws from a Gaussian PDF with variance equal to that from the rows of \mathbf{U}^{20c} and \mathbf{U}^{21c} (all of the rows have equal variance in each case). As a result, data in these random matrices are independent in directions corresponding to both the EOF number and the model number. We desire matrices with an independent model dimension in order to test the likelihood that CMIP5 output was drawn from a set of independent models. Having independence in the field direction is appropriate because the columns of \mathbf{U}^{20c} and \mathbf{U}^{21c} are independent by construction.

Our assumption that the t -dimensional normal distribution is representative of an independent ensemble of climate projections is subject to some caveats; we are making the effective assumption that a normal distribution of models in the space defined by \mathbf{U}^{20c} or \mathbf{U}^{21c} is plausible and that there are no parts of that space that might represent an unphysical climate state. There are some justifications for this assumption; the random distributions are compared with the loading matrices \mathbf{U}^{20c} and \mathbf{U}^{21c} , which are themselves orthogonal basis sets defined by multimodel variability. As such, we are making the assumption that, if there are physical relationships between variables in the model output data (e.g., between adjacent grid cells or between surface temperature and outgoing longwave radiation), then any correlation between these would be represented as a single mode in the EOF analysis. Thus, any linear relationships that exist in the original data are effectively

preserved in the random ensemble also. However, a strong nonlinear relationship between two variables in the CMIP5 archive could not be represented in a single EOF mode and might be represented in two or more modes. In this case, then there would be some of the space that should be physically off limits. Hence, by using normally distributed data to define the random ensembles and their associated length scale for interpoint distances, we make the assumption that multimodel variability can be appropriately described by a linear basis set. Although one could potentially consider designing a random sample that fitted a high-dimensional distribution to the existing ensemble to account for nonlinear relationships between modes, the increase in complexity, the lack of samples in the original ensemble, and the necessary subjective parameterization of such a distribution means this is impractical for the present study.

d. Weighting for uniqueness

In this section, we seek to use the relationships derived in section 2b to define a weighting scheme that would effectively downweight closely related model pairs the ensemble, which we can assess using the expectation values for near-neighbor distances in the random ensembles proposed in section 2c. Our scheme should also provide the capability to downweight models that exhibit low fidelity in a desirable metric.

The limiting cases of such a scheme are easy to define. We consider the models, as before, to be represented as points in a space defined by the loadings of the model in an ensemble-wide EOF analysis. In the extreme case, if the distance between two models is exactly zero then the models are considered identical and each member of the pair should be given half the weight that they would otherwise have (equivalently, a statement that adding an identical model to an existing ensemble member should not change the results).

We propose a simple functional form for model similarity that satisfies the requirements for a given model pair (i, j) , separated by a distance δ_{ij}^{20c} or δ_{ij}^{21c} ,

$$S(\delta_{ij}^{20c}) = \exp \left[- \left(\frac{\delta_{ij}^{20c}}{D_u} \right)^2 \right] \quad \text{and} \quad (10)$$

$$S(\delta_{ij}^{21c}) = \exp \left[- \left(\frac{\delta_{ij}^{21c}}{D_u} \right)^2 \right], \quad (11)$$

where D_u is a free parameter, a “radius of similarity,” such that model pairs separated by less than this value are considered similar. The distance is squared so that the metric tends to unity for values $\ll D_u$. The smallest

reasonable value for D_u would be the expected distance between two identical models exhibiting different realizations of internal model variability, given this represents a case where the model structure is identical. As D_u is increased from this value, increasingly distant pairs of models are considered similar. In the extreme case, as D_u approaches the largest interpoint distances (i.e., the largest values of δ_{ij}^{20c} or δ_{ij}^{21c}) in the ensemble, then only the models with the largest biases would exhibit a value of S of close to unity and all other members would be downweighted.

In section 2c, we derived D_u empirically by considering the nearest neighbors one would expect to find by chance in a t -dimensional normal distribution of equal population, variance, and dimensionality as U . This is achieved in practice by considering the randomly generated distributions from the section 2a. We define D_u to be the 50th percentile of nearest-neighbor distances in the 10^5 randomly generated ensembles.

One can thus obtain a value for the effective repetition of model i in the ensemble,

$$R_u(i)^{20c} = 1 + \sum_{j \neq i}^m S(\delta_{ij}^{20c}) \quad \text{and} \quad (12)$$

$$R_u(i)^{21c} = 1 + \sum_{j \neq i}^m S(\delta_{ij}^{21c}), \quad (13)$$

for the past and future cases, respectively, where m is the total number of models. We then propose a uniqueness weighting for model i by taking the inverse of the number of models similar to i ,

$$w_u(i)^{20c} = [R_u(i)^{20c}]^{-1} \quad \text{and} \quad (14)$$

$$w_u(i)^{21c} = [R_u(i)^{21c}]^{-1}. \quad (15)$$

for the past and future cases, respectively. If desired, a weighting scheme could also consider model quality; a model should be given increasingly less weight the farther that model lies from the point representing the observations in the EOF space. In the limiting case, the model weight should tend to zero as the distance of the model to the observations tends to infinity. These attributes are satisfied by the following construction for w_q , the model quality weighting:

$$w_q(i) = \exp \left[- \left(\frac{\delta_{i(\text{obs})}^{20c}}{D_q} \right)^2 \right], \quad (16)$$

where $\delta_{i(\text{obs})}^{20c}$ is the Euclidean distance between the EOF loading for model i and the loading of the observed climatology projected onto the same EOF basis set. This

is only calculated for the historical data where observations are available. The parameter D_q is a “radius of model quality” and is a free parameter in the weighting scheme. As $D_q \rightarrow +\infty$, then $w_q \rightarrow 1$ for all models, and the quality weighting has no distinguishing effect. As the value of D_q is reduced, models closer to the observations are increasingly upweighted. The smallest reasonable value for D_q would be the smallest observational bias seen in the ensemble [i.e., $\min(\delta_{i(\text{obs})})$]. In the extreme case as $D_q \rightarrow 0$, the majority of the weight is placed on the single best performing model.

To explore the sensitivity to this parameter, we consider two values for D_q : a “wide” choice where D_q is equal to the mean intermodel distance in the CMIP5 ensemble and a “narrow” choice that is half of this value. Expressing D_q in terms of the CMIP variance has the disadvantage that the variance itself can be influenced by both model quality and reproduction, but this decision is a matter of practicality. We present the values of D_q as subjective, effectively as a statement that relative skill, rather than any absolute measure, should define whether we accept or reject a model. In effect, the wide case describes a situation where only the models with the largest biases in the ensemble are down-weighted, while in the narrow case a distinction is made between the “average” and “best” performers. It might be desirable to let internal or natural variability define D_q , but, as we show in section 3a, this would lead to a situation where $\delta_{i(\text{obs})}^{20c} \gg D_q$ for all i , which, given Eq. (16), would place the majority of the weight on the model with the lowest value of $\delta_{i(\text{obs})}^{20c}$.

e. Eliminating interdependent models

If the researcher’s goal is simply to produce a multi-model average that is less susceptible to bias by model replication, then simply weighting each model by the appropriate value of w_u would suffice. This approach could be used directly for calculating a central estimate of combined multimodel projections.

However, some issues associated with model co-dependence cannot be solved by weighting alone. For example, the potential bias associated with regression-based predictions of unknown climate parameters can only be addressed by removing the interdependent models. This can be achieved in a pure statistical fashion (see Caldwell et al. 2014), but the interpretation of such constructions is not always intuitive.

We propose here a less formal approach that should be readily reproducible for a variety of purposes where it is desired to remove the most blatant model co-dependencies. Our method is a stepwise model elimination, where the models with the highest codependencies are removed first.

The simplest approach here would be to recursively remove a member of the closest near-neighbor pair until the remaining ensemble conforms to a plausible random distribution in the n -dimensional EOF space. Since better models are replicated more, however, such an approach preferentially eliminates the models clustering closer to observations while models with large biases would be preserved. This has a significant detrimental effect on the mean performance of the remaining ensemble. Instead, we propose a strategy that considers both model performance and model independence when creating an ensemble subset.

First, we introduce a bulk quantity that describes the ensemble characteristics, the “independent ensemble quality score,”

$$S_m^{20c} = \sum_i^m w_u^{20c}(i) w_q(i) \quad \text{and} \quad (17)$$

$$S_m^{21c} = \sum_i^m w_u^{21c}(i) w_q(i), \quad (18)$$

for historical and future cases, where w_u^{20c} , w_u^{21c} , and w_q are described in section 2d as the individual model weights corresponding to model i . Using the product of the two weights is a subjective decision, and other functional forms could potentially be explored. However, as we now demonstrate, this simple combination of the uniqueness and quality weights addresses our goals to remove the influence of exactly replicated models and of very poor models.

This can be illustrated as follows for the historical simulation: If an independent model is added to the ensemble, $w_u^{20c}(i)$ equals 1 for model i and so S_m will increase by the model quality score $w_q(i)$. The increase is large for a high performing model and approaches zero for a very poor model. However, if two identical models i and j are added to the ensemble together, $w_u^{20c}(i)$ and $w_u^{20c}(j)$ each equal 0.5, and so S_N will still only increase by $w_q(i)$.

If we start with an N -member ensemble, we eliminate a single member by considering the maximum possible ensemble quality score for each combination of $N - 1$ members. The excluded model j is removed from the ensemble and the process is repeated until an appropriate stopping criterion has been reached. We can assess the effective number of models remaining at any point by considering the “number of effective models,” for both historical and future cases,

$$n_{\text{eff}}^{20c} = \sum_i^m w_u^{20c}(i) \quad \text{and} \quad (19)$$

$$n_{\text{eff}}^{21c} = \sum_i^m w_u^{21c}(i), \quad (20)$$

with each representing the sum of the uniqueness weights for the remaining models in the ensemble.

The approach outlined here is quantitative but subjective, with a number of free parameters. To demonstrate its utility, we consider a case study of the CMIP5 ensemble, where we can objectively demonstrate that we can use the algorithm to produce a subset of CMIP5 models that provides comparable model diversity, improved mean model performance, and reduced model replication in comparison to the original model sample.

3. Results

a. CMIP5 ensemble properties

The initial dataset from which we draw our conclusion is the matrix of pairwise distances between models in the CMIP5 archive, δ^{20c} and δ^{21c} , which are calculated from \mathbf{U}^{20c} and \mathbf{U}^{21c} matrices. This matrix is represented graphically in Fig. 1 for the all-variable case using both present-day climatological fields calculated from 1970 to 2000 in historical simulations and the anomalies from those fields in the RCP8.5 simulation between 2070 and 2100. In both cases, recognizable structure relating to model genealogy is visible in the intermodel distance field.

We can compare, in a bulk sense, the distribution of distances in the matrices to that one might expect from a purely random distribution. The distributions for the CMIP5 derived matrix and the random distributions are plotted in Figs. 2a,b for a number of different variable choices.

The random distributions have the same variance as the original CMIP5 distributions by design because each dimension of the random pseudo ensembles is normally distributed with the same variance as the original CMIP5 case in each dimension of \mathbf{U}^{20c} and \mathbf{U}^{21c} . Because we consider a large number of pseudorandom normally distributed ensembles, we can produce best estimates and confidence intervals for the distribution of intermodel distances one would expect if the models were normally distributed in the space defined by \mathbf{U}^{20c} and \mathbf{U}^{21c} . If the CMIP5 distribution falls outside of this range, this implies that the models in CMIP5 are distributed in a nonnormal fashion in the space.

We find there are some significant deviations in the CMIP5 distribution from what one would expect in a purely random case. First, there are a number of model pairs that lie closer to each other in the EOF space than ever occurs by chance in the random samples (less than 50% of the expected mean interpoint distance for the random case). However, there is also an absence of

models at intermediate distances (between 50 and 90% of the mean interpoint distance), relative to the random distributions. This indicates that the distribution of CMIP5 models in the EOF space has a rather heterogeneous, clustered distribution, with families of closely related models lying close together but with significant voids in-between model clusters. These features are especially clear in the future case, where the distances are measured in terms of (2070–2100) anomalies from the (1970–2000) climate mean state. We also show the histogram of intermodel distances in initial condition CCSM4 ensemble, demonstrating that intermodel distances resulting from internal model variability alone are an order of magnitude smaller than the mean intermodel distances seen in the CMIP5 archive.

The responsible model pairs can be explicitly plotted. Figure 3a shows model pairs that are closer together than the expected nearest-neighbor distances in the random distributions, using all variables. Many of these samples correspond to identical models from the same institution submitted at a different resolution (e.g., IPSL-CM5A-MR/LR and MPI-ESM-LR/MR). Other model pairs relate to changes in model configuration that do not influence the set of atmospheric diagnostics considered here (e.g., HadGEM2-AO and HadGEM2-ES share the same atmosphere, ocean, and ice models, but the former lacks treatment of the carbon cycle, which has little effect in these concentration-driven simulations). Finally, there are some cases where models from two institutions share a large fraction of code base, and this is reflected in their proximity in EOF space (e.g., HadGEM2-AO and ACCESS1.0 or FIO-ESM and BNU-ESM). Several other model pairs are plotted with dotted lines. These, to a lesser degree, still occur closer together than one might expect by chance (for the models joined by a black line, one such pair would be expected by chance in a 36-member ensemble). These connections can also be related to common model components (e.g., NorESM and CCSM4 share atmosphere and land surface, and MPI-ESM and CMCC-CSM5 share atmospheric code). We also include the observational point in the same analysis in Fig. 3a, which shows that none of the models in the CMIP5 archive is considered closer to the observations than would be expected by chance. In the later part of the study, where we prune similar models from the archive, this gives us some confidence that similar models are not being removed because they are all converging on the “true” climate. We can repeat the analysis for future changes in the same variables (Fig. 3b), which show a similar close relationships to present-day case. Using specific fields produces similar (but nonidentical) relationships (Figs. 3c–e). The all-variables case shows

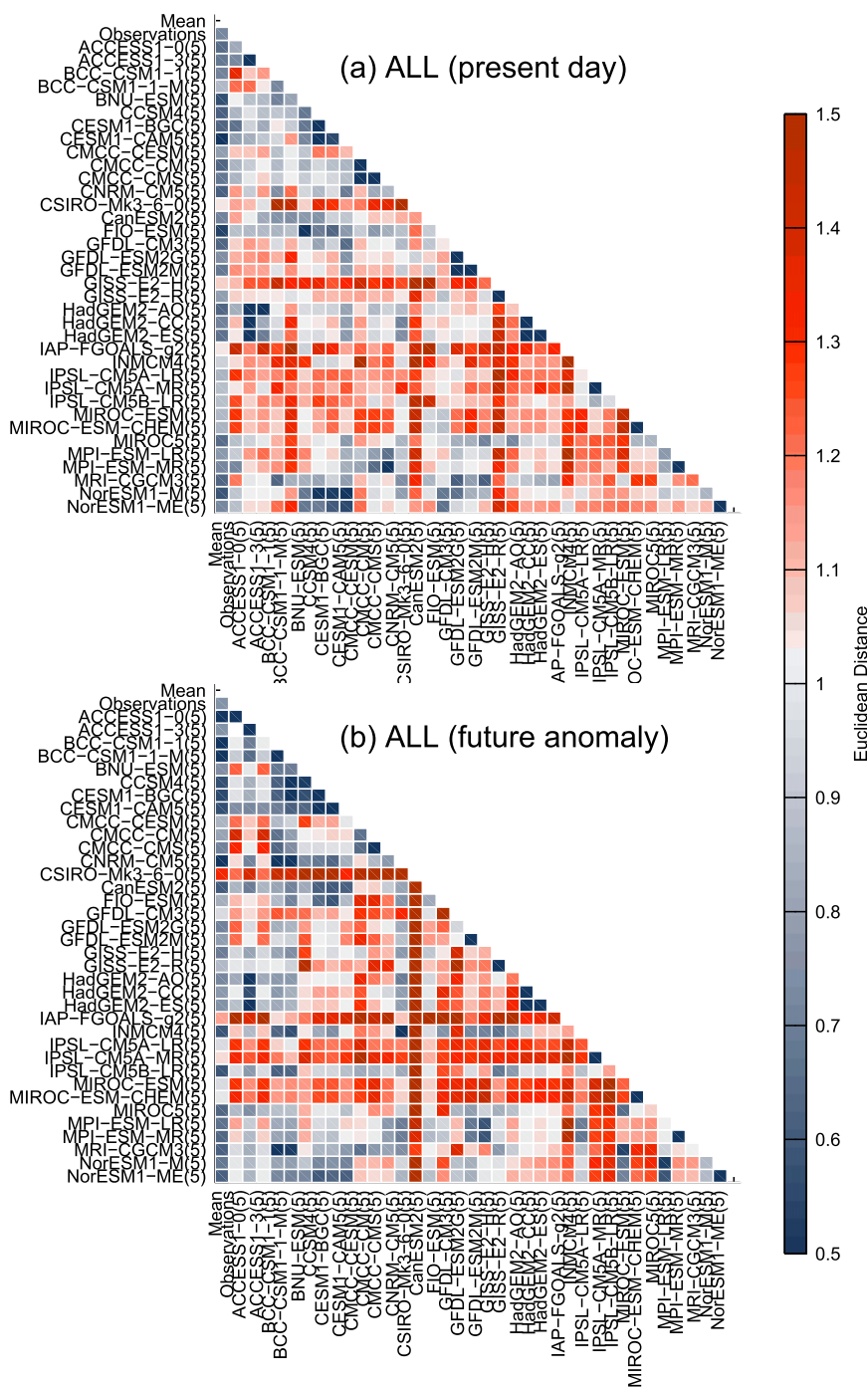


FIG. 1. A graphical representation of the intermodel distance matrix for CMIP5 calculated for ALL using (a) 1970–2000 monthly mean climatological fields (as defined in Table 2) and (b) changes in the aforementioned fields between 1970–2000 and 2070–2100. Each row and column of (a),(b) represents a single climate model (or observation). Each box represents a pairwise combination, where warm colors indicate a greater distance. Distances are measured as a fraction of the mean intermodel distance in CMIP5.

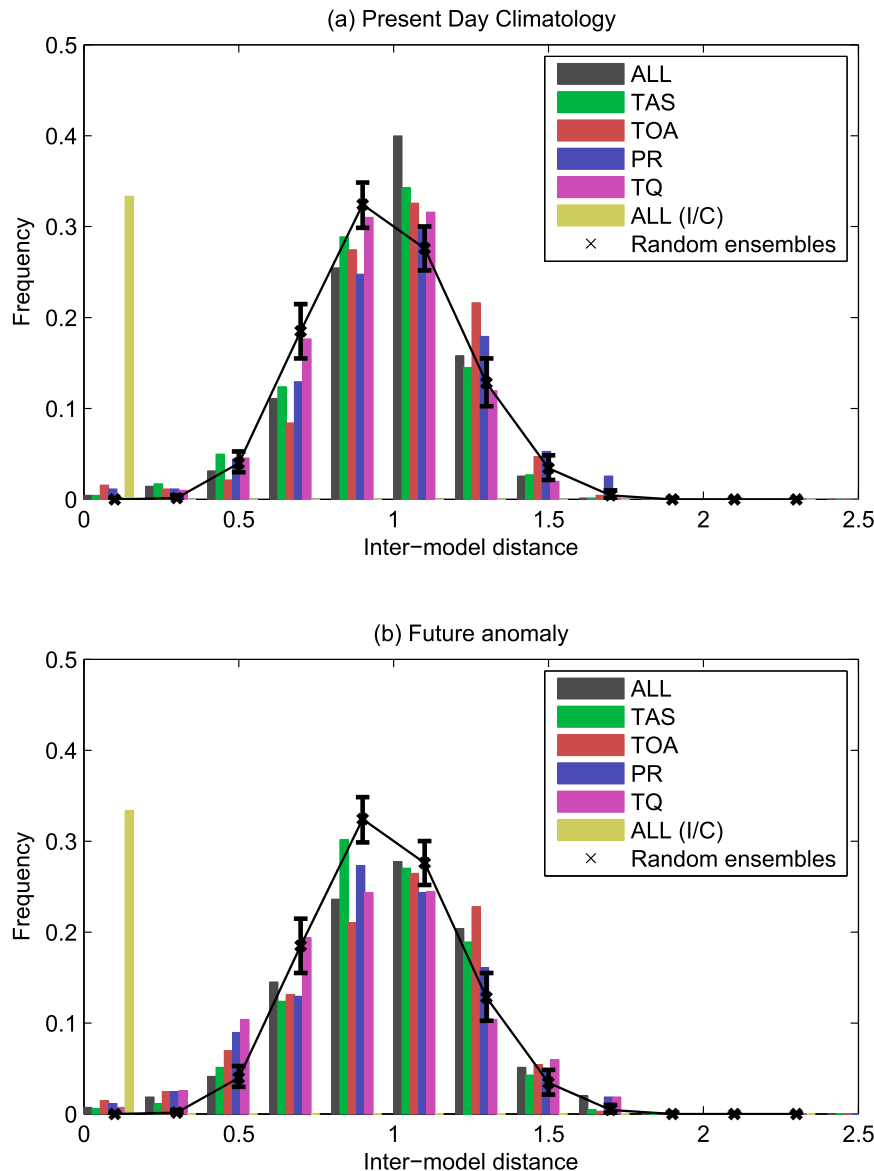


FIG. 2. Histograms of CMIP5 intermodel Euclidean distances in the EOF loading space derived from (a) 1970–2000 monthly mean climatological fields (as defined in Table 2) and (b) changes in the aforementioned fields between 1970–2000 and 2070–2100, as compared to a sample of 10^5 histograms calculated from randomly sampled distributions. Gray bars show the histogram of intermodel distances in the CMIP5 ensemble in an EOF space constructed with all available variables, while other colors show distances constructed with only a subset of variables: TAS, TOA shortwave and longwave fluxes, PR, and TQ. The yellow bars indicate the distribution using all variables from the CCSM4 initial condition ensemble. The box-and-whisker plots show the range of bin values observed in the random distributions showing the 10th, 50th, and 90th percentiles of the distribution.

that all close relationships would be expected from a genealogical perspective. However, when one uses single variables (PR especially), there are some unexpected results (e.g., MIROC and CAM5 are considered close). We attribute this to the difficulty of representing intermodel precipitation variability in a low-dimensional

basis set (although models from different centers may in some cases share parameterizations).

b. Stepwise model elimination

There are various arguments to support the hypothesis that the CMIP5 ensemble is biased by the inclusion

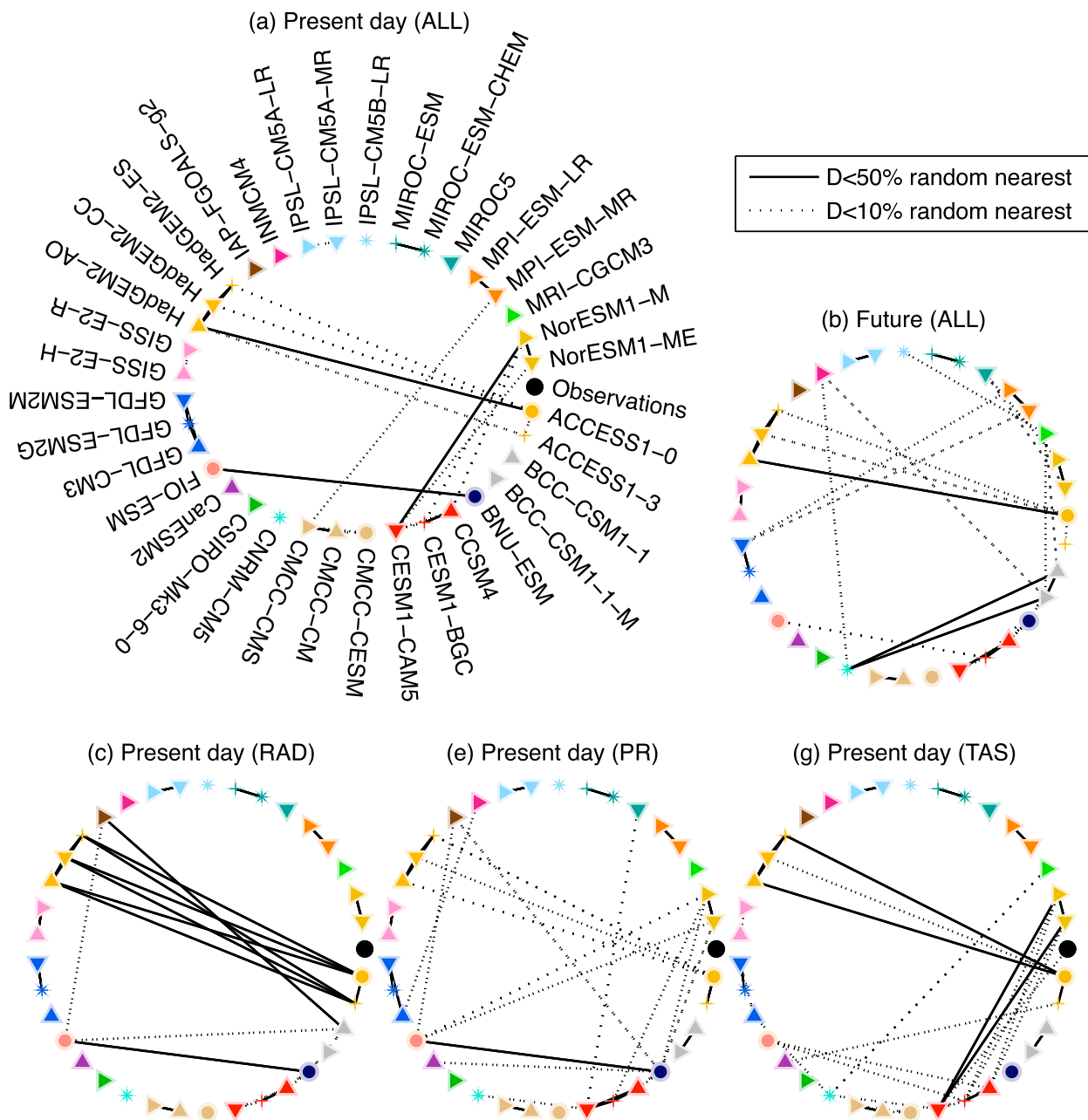


FIG. 3. An illustration of intermodel and observation–model distances in an EOF space defined by (a) 1970–2000 simulated climatology for ALL and (b) the anomaly between 1970–2000 and 2070–2100 under the RCP8.5 scenario for ALL. Plots are repeated for individual variables: (c) TOA shortwave and longwave fluxes, (d) precipitation, and (e) surface air temperature. Intermodel lines illustrate where the intermodel distance is less than 50% (dotted) or 90% (solid) of nearest interpoint distances in a randomly generated distribution of with the same dimensionality, variance, and population.

of common components: some of which are featured more frequently than others. One can make this argument from a consideration of the models themselves (see section 1 and Table 1) or by examining the spatial distribution of models in orthogonal dimensions derived from model output. We have proposed a method of model removal that maximizes a metric reflecting both

model diversity and fidelity. The iterative model elimination process is illustrated for the CMIP5 ensemble in Fig. 4.

The plot shows the consecutive removal of models from the set of 36 considered in this study until a single model remains. The process is demonstrated by eliminating interdependent models as judged by the

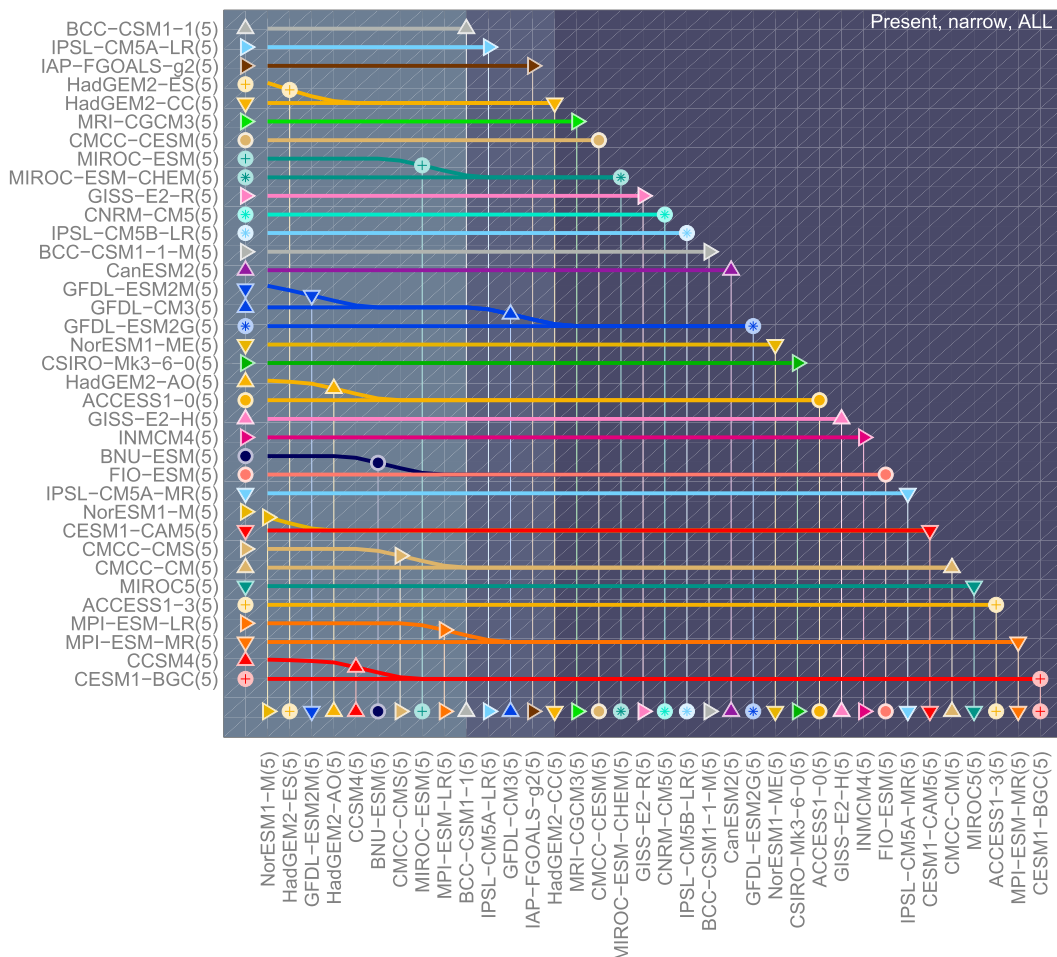


FIG. 4. An illustration of the stepwise model elimination procedure outlined in section 2e as applied to the 36 models from the CMIP5 ensemble, using model similarity information from the present-day (1970–2000) climatology for ALL and the wide quality radius. The full set of models is shown on the left axis, and the order of model removal is shown along the bottom axis, with the leftmost model removed first. If the number of effective models n_{eff} decreases by less than 0.5, then the removed model is shown merging with its nearest neighbor in EOF space. If the number of effective models decreases by more than 0.5, the line ends, indicating the removal of that model family from the ensemble. Background shading indicates whether the smallest interpoint distance in EOF space using the remaining archive is less than 90% (light gray), 50% (mid gray), or 10% (dark gray) of purely random distributions of the same population, variance, and dimensionality.

simulation of present-day climatology. The model quality weights w_q are obtained using the mean state climatology from the models as compared to the observations. Model uniqueness is calculated as in section 2e after each iteration.

We demonstrate the sequence of model removal in Fig. 4 (for present-day similarities, all variables and a wide quality radius). The figures show the order in which models are removed from the archive to achieve the maximum independent ensemble quality. If the removed model is closer than D_u (a function of the number of models remaining) to any other remaining model, then that model is shown to merge with its nearest

neighbor. However, if the model is further than D_u from any other model, the model branch is shown as terminating in the diagram.

We have not yet fully discussed an appropriate point to stop trimming models. This question is ultimately subjective, and the conclusion is somewhat dependent on the specific needs of the researcher. However, Fig. 5 shows some changing characteristics of the remaining ensemble as the ensemble size is decreased, and these can be used to recommend ensemble subsets for different scenarios. In essence, a first phase of eliminating models just removes redundant data, and a second phase improves the characteristics of the ensemble by

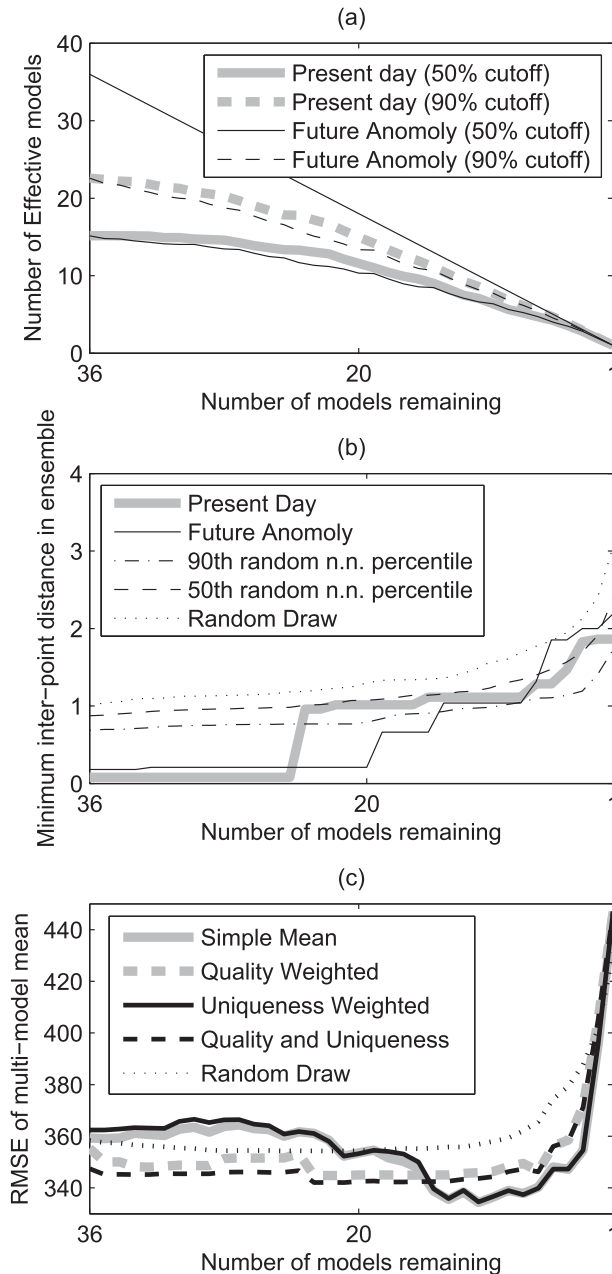


FIG. 5. Plots illustrating the stepwise model elimination following the procedure in section 2e. Calculations are conducted using model similarity metrics derived from both present-day climatology and from future climate change under RCP8.5. (a) The number of effective models as a function of the number of actual models remaining in the ensemble. The percentile cutoff is the fraction of nearest-neighbor distances seen in purely random ensembles used to define the radius of similarity D_u in Eq. (10). (b) The nearest-neighbor distance as a function of the number of models remaining. For comparison, the 10th, 50th, and 90th percentiles of nearest-neighbor distances in purely random ensembles of the same dimensionality and variance are shown. (c) RMSE of weighted and unweighted multimodel means as a function of remaining models.

removing poor models and partly redundant ones. Going beyond that potentially worsens the ensemble mean bias representation.

Figure 5a shows how n_{eff} varies as models are removed from the archive as described in section 2e. The actual number is dependent on the choice of D_u , the radius of similarity. Two choices of D_u are illustrated, using either the 50th percentile of nearest-neighbor distances in the set of 10^5 random ensembles (as was used in section 2d) or, for comparison, the 90th percentile. Using all the models in the archive, n_{eff} is 15.5 using the larger value for D_u or 22.5 using the smaller value (using present-day climatology metrics of similarity). The removal of the first 10 models has little effect on n_{eff} (especially using the larger value of D_u). The removal of the remaining models results in a monotonic decrease in n_{eff} .

As was indicated by Fig. 5a, most of the early model eliminations have little effect on n_{eff} . Figure 4 shows that many of the initial removals represent models (from CCSM4 to CESM1(BGC), from HadGEM2-ES to HadGEM2-AO, and from GFDL-ESM2M to GFDL-ESM2G) that are largely structurally identical, at least in terms of their long-term atmospheric climatology, differing only in the presence of an active carbon cycle that would not influence the diagnostics used in this study. It is thus largely random which member of the pair is eliminated. In this regime, there is a strong inverse relationship between model quality weights (w_q) and uniqueness weights (w_u), as shown in Fig. 6a.

The second broad class of eliminations is models with strong connections, often from the same institutions but with some differing components. In these cases, the model with the higher value quality weighting (w_q) is generally preserved (e.g., GISS-E2-H and GISS-E2-R, which differ in their ocean components). In this regime, the inverse relationship between the model quality weight and uniqueness weights is weaker (Fig. 6b), as the clear duplicates have already been removed. Note that the uniqueness weights now refer to uniqueness within the remaining subset and not within the full CMIP5 archive.

The final stages of removal (approximately the final 20 models) do result in a reduction in the number of effective models, illustrated by the termination of the model path. As shown in Fig. 5b, in this regime, the distribution of intermodel distances are now consistent with what one might expect from a purely random sample. Each family of closely related models is now represented, to a large extent, by its own “champion.” Figure 6c shows that when only 10 models remain, the relationship between w_u and w_q is rather weak, with all remaining models having comparable uniqueness weights.

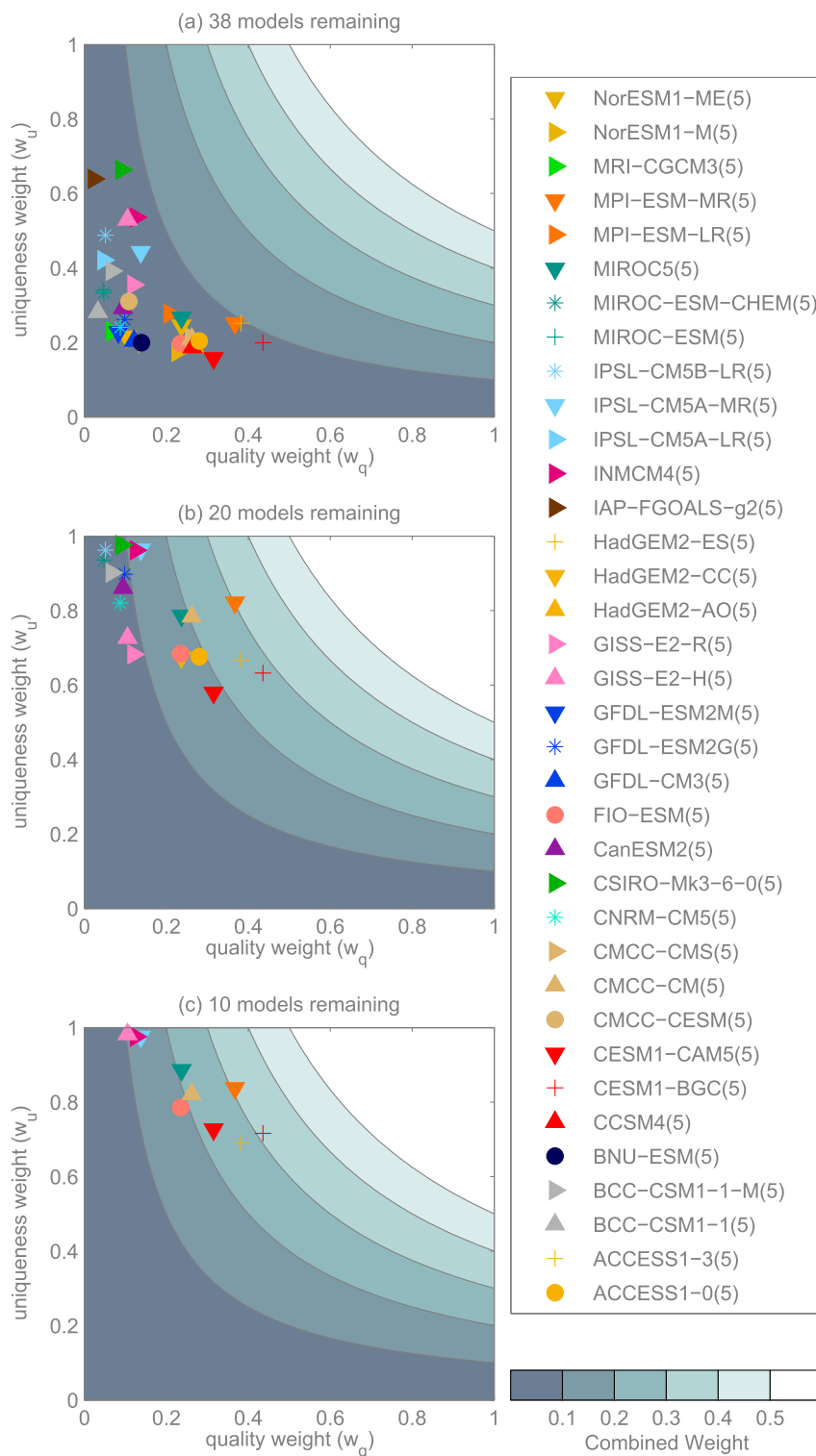


FIG. 6. A plot demonstrating how model uniqueness weights and model quality weights change as models are eliminated in the sequence shown in Fig. 4, for (a) 36, (b) 20, and (c) 10 models remaining.

Our value judgment for an appropriate stopping criterion is thus dependent on the application. If one wishes to only remove near-identical models, one should stop trimming when the number of effective models n_{eff} begins to significantly decrease. However, if one wishes to produce the best performing ensemble mean simulation of the mean state, it is more logical to also remove the worst performing models such that the RMSE error of the subensemble mean is minimized.

c. Sensitivity to initial choices

The algorithm as described in section 2e requires several assumptions and we explore the sensitivity of the results to those choices in this section. Figure 7 shows the models that are retained in the analysis with a range of different initial variable and parameter choices. In each case, the analysis is repeated and there is a stepwise removal of models based on maximizing the ensemble quality score. On each line of the plot, we show which models remain when the smallest interpoint distance in the remaining archive is first greater than 50% (unfilled symbols) or 10% (filled symbols) of purely random distributions of the same population, variance, and dimensionality (regions marked by mid-gray and dark gray shading in Fig. 4). Thus, we can explore the sensitivity of the retained models to our initial assumptions.

First, there is the choice of which variables are used to derive the intermodel distance matrix. To address this, we repeat the analysis with a variety of individual fields, as well as the multivariate example discussed in the previous section. The analysis is repeated for zonal mean temperature and humidity (TQ), gridded PR, gridded top-of-atmosphere (TOA) shortwave and longwave fluxes, gridded TAS, and all variables combined (ALL). Second, we explore the radius of model quality D_q introduced in Eq. (16). The analysis is repeated for two values, a wide value where D_q is equal to the mean intermodel distance in the CMIP5 ensemble and a narrow choice that is half of this value. The latter narrow case effectively increases the role of the model quality metric, such that models with a low quality score are removed earlier in the algorithm, unlike in the wide case, where highly interdependent models are removed first. Finally, we construct the model uniqueness weightings w_u using the intermodel distances derived from the 30-yr mean 1970–2000 present-day data in the “present” case but use the anomaly between 2070–2100 and 1970–2000 for the “future” case.

We find that variable choice has little impact on the final choice of model subsets. Although in some cases the choice of model from a given institution can change, the overall number of models retained is similar for each of the variable choices. The use of the narrow radius of

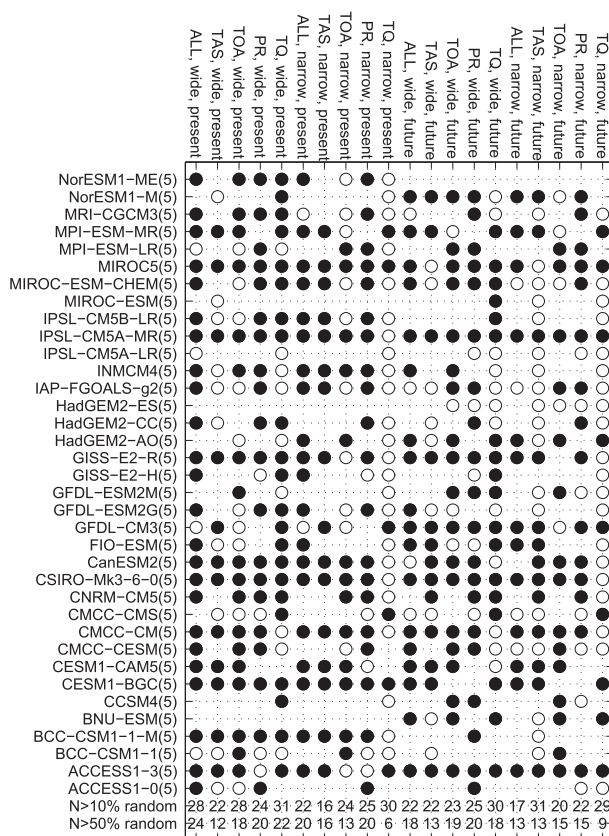


FIG. 7. A plot showing suggested subsets of CMIP5 given model quality scores and codependencies derived in a number of ways. Each line in the figure repeats the analysis leading to Fig. 4 with different assumptions. Plotted are the remaining models where the smallest interpoint distance in EOF space using the remaining archive is greater than 10% (unfilled symbols) or 50% (filled symbols) of purely random distributions of the same population, variance, and dimensionality (regions marked by mid-gray and dark gray shading in Fig. 4). The analysis is conducted with TQ, gridded PR, gridded TOA shortwave and longwave fluxes, gridded TAS, and all variables combined. The parameter D_q , the radius of model quality, is set to wide or narrow (the latter increasing the role of model quality metrics in model elimination). The variable w_u , the model uniqueness weighting, is shown calculated with the future RCP8.5 data or the present-day data. Numbers at the bottom of the plot indicate the number of retained models for the two conditions where the minimum remaining intermodel distance is greater than the 10th or 50th percentile of the random smallest intermodel distances.

model quality, however, significantly decreases the number of retained models with respect to the wide value. This can be explained by considering that the narrow setting increases the ratio of the model quality weighting for models lying close to the observations and those far away. In the narrow regime, the ensemble quality score is best maximized by removing the poorly performing models earlier in the analysis; thus, after the interdependent remaining models have been removed,

the number of remaining unique models is smaller than in the wide case.

EOF TRUNCATION CHOICES

Some subjective decisions are required in the interpretation and subsequent usage of the PCA conducted in section 2a, and we discuss these at greater length here. In previous studies like Masson and Knutti (2011), the intermodel distances were calculated without the PCA stage, simply calculating distances in the space defined by the anomaly matrices $\Delta\mathbf{X}^{20c}$ and $\Delta\mathbf{X}^{21c}$. For the purposes of this study, and its companion studies (Sanderson et al. 2015), it is necessary to decrease the dimensionality (and codependence) of the data in order to establish prior expectations of near-neighbor distances.

In this study, as in Sanderson et al. (2015), the intermodel distances are calculated with the truncated set of nine modes. The resulting intermodel distance matrix calculated with \mathbf{U}^{20c} truncated to nine modes has a 0.93 correlation with the matrix one would calculate using the full-field matrix $\Delta\mathbf{X}^{20c}$, but using the orthogonal basis set allows us to form random matrices with which to compare the results (Fig. 2).

For smaller values of t , only the leading patterns of model difference are retained, which results in large intermodel distances between different model families (e.g., CESM1 and GFDL models) and very small distances between models in the same family [e.g., CESM1(CAM5) and CESM1(CAM4)]. With such few degrees of freedom, very small intermodel distances cannot be ruled out by chance in the random ensembles, and so no models can be excluded from the ensemble (see Fig. 8 for truncation values of 3 or less). The analysis produces very similar results, and the minimum number of retained models, for values of t between 8 and 12 (see Fig. 8), with relatively little sensitivity to variable choice (not shown). For values of t of 15 or greater, the higher-order modes increasingly represent subtle and often noisy differences between models in the archive, which inflates the distance between the near neighbors in the ensemble. Hence, once again we see fewer models ruled out.

To test the sensitivity of the intermodel distance matrix to variable choice, we also repeat the EOF analysis with a number of different subsets of diagnostic variables. The resulting correlation depends significantly on which exact variable is retained. The intermodel distances calculated using gridded TAS only are highly correlated with the multivariate case ($R = 0.95$ untruncated). Top-of-atmosphere radiative fluxes (RAD; $R = 0.85$ untruncated), PR ($R = 0.66$ untruncated), and zonally averaged vertical temperature and humidity

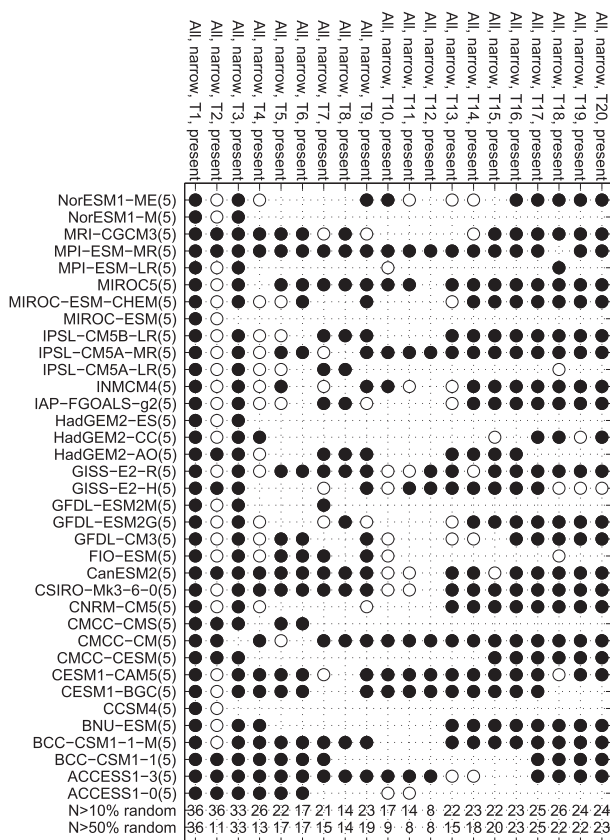


FIG. 8. As in Fig. 7, but showing suggested subsets of CMIP5 with different truncation lengths for the EOF analysis. Plotted are the remaining models where the smallest interpoint distance in EOF space using the remaining archive is greater than 50% (unfilled symbols) or 10% (filled symbols) of purely random distributions of the same population, variance, and dimensionality (regions marked by mid-gray and dark gray shading in Fig. 4).

(QT; $R = 0.42$ untruncated) are increasingly poorly correlated with the full-field multivariate case. This implies that some fields, such as surface temperature have sufficient information to render a multivariate approach unnecessary.

With a truncation length of 9, which we used for the bulk of this study, the resulting distance matrix remains highly correlated to the full-field distance matrix, but the influence of covariant fields and models is reduced [see Caldwell et al. (2014) for an extensive discussion of these issues].

d. Ensemble mean performance

The results of section 3b suggest that eliminating the strongest interdependent models to leave a plausibly random distribution would leave between 10 and 25 of the 36 CMIP5 models considered here (depending on variable and parameter choices). Trimming the ensemble to its more independent subset does not worsen the

fidelity of the climatological mean result, and removing the poorer performing outliers (models with large biases) can actually improve it, as we show in this section.

We can first examine how the multimodel mean of present-day climatology compares against observations. Figure 5c considers the root-mean-square errors (RMSEs) of various weighted and unweighted multimodel means calculated using the same multivariate climate state vectors described in section 2a and the observations listed in Table 2. We illustrate this using the ALL case, with the wide radius of model quality and present-day derived intermodel distances. We also compare with the average RMSE seen when a completely random sample (without replacement) of the same size is taken, as compared to the detailed technique outlined in section 3b.

If one considers only the far left of the plot, where all 36 models are retained, weighting the models by uniqueness actually increases the RMSE. This is largely to be expected—as we have seen in Fig. 6a that the best performing models have the lowest uniqueness weights. It also suggests that a mean of the CMIP5 ensemble is already weakly weighted toward the better performing models. If we explicitly weight the model mean toward models that lie closer to the observations in the EOF space, the RMSE can be reduced significantly.

As the first 10 (highly interdependent) models are removed from the archive, the simple mean RMSE increases slightly while the random draw RMSE remains constant, likely because the high-performing models have less representation when the duplicates have been pruned. The uniqueness weighted mean also becomes more similar to the simple mean case (u_w is now more consistent across the ensemble). Between 28 and 12 models remaining, the simple RMSE decreases significantly; when 20 models remain, the subset outperforms the RMSE of the random sample. The lowest RMSE values occur with between 12 and 5 models remaining. Removing any further models increases the RMSE of the simple multimodel mean. With 5 or fewer models remaining, all models have a high value of both w_u and w_q , so weighting by uniqueness or quality has little effect. In all cases, any further removal of models (below 5) significantly increases the RMSE, a fact that is likely attributable to the Cauchy–Schwarz inequality (Annan and Hargreaves 2011).

4. Discussion and conclusions

The present study considers how one might remove potential biases that might arise from shared components in the CMIP5 archive of climate models and its predecessors. We also propose some simple diagnostics

that might be used to identify interdependent models using model diagnostic output and a possible strategy to choose a model subset to maintain model diversity without replication and to incorporate model quality information into this decision.

This study represents a proof of concept; the choice of diagnostics used in this study is of course arbitrary, to some degree, although the results of which models are interdependent do seem to be relatively resilient to changes in variable and time period (see Fig. 3; Pennell and Reichler 2011; Knutti et al. 2013). However, we do assume that a model's mean state climatology can be used to assess both its skill and independence. Clearly, if our final goal is to assess the plausibility of a model's future simulations, then the mean state simulation is not a perfect assessment of model skill, although it could be argued that it is a necessary condition and as such a weighting strategy based on present-day climatology can be justified in the absence of any additional information.

Certainly, which model exhibits the highest quality score is very much dependent on the specific metrics in which the researcher might be interested (Santer et al. 2009), and it is far beyond the scope of this study to conduct an exhaustive comparison of possible model metrics. In this study, we have focused primarily on diagnostic output from the atmospheric model, and our results are thus liable to be most sensitive to common component in that model. As such, the results of this study should be interpreted as illustrative of a potential method for reducing the effects model interdependency and not as a prescriptive list of models that should be used for future studies. Most studies based on CMIP5 could easily use such a framework, but the value judgements of future researchers should be embedded into the choice of metric used to assess model similarity and quality.

We assess the likelihood of near-neighbor models occurring by chance using a large number of random distributions of the same dimensionality as the truncated orthogonal set of EOF loadings we derive from the original ensemble. The random sample is not a proxy for the space that might be attainable by the real climate; rather, it is a proxy for the distribution of models represented in an orthogonal basis set defined by multimodel variability. As such, we are making the assumption that, if there are physical relationships between variables in the model output data (e.g., between surface temperature and outgoing longwave radiation), then any correlation between these would be represented as a single mode in the EOF analysis. However, if there is a strong nonlinear relationship between two variables in the CMIP5 archive, then this relationship could not be represented in a single EOF mode and

might be represented in two or more modes. In this case, then the distribution of models in the space could be more complex than a simple Gaussian. One could imagine designing a random sample that fitted a high-dimensional distribution to the CMIP5 ensemble to account for such nonlinearities, but the increase in complexity, the lack of samples in the original ensemble, and the necessary parameterization of such a distribution means this is impractical.

We also assume, by drawing random samples using the variance defined by the original ensemble, that none of the CMIP5 members can be ruled out a priori. One could imagine a situation where an arbitrarily poor model was included in the ensemble that would increase the variance represented in each mode such that any realistic models would look self-similar and would be downweighted by the uniqueness weighting. Therefore, the method only makes sense if there is some level of base confidence that none of the models in the archive is completely unrepresentative of the true system. However, we would argue that this is true of any analysis that uses the CMIP5 archive and that even a simple multi-model mean is subject to a sanity check of the participating models.

Caveats aside, this study illustrates some interesting characteristics of the CMIP5 archive and potential issues that might arise from treating this archive as a random sample of possible climate models. There is extensive replication of model code in the archive, primarily within institutions but also in some cases between institutions (see Table 1). This should come as little surprise: a quick examination of AOGCM makeup in the CMIP5 models indicates that some individual components are used by over 25% of the archive. However, we show in this study (e.g., Masson and Knutti 2011) that many of those similarities can be identified also through a simple analysis of model output. A more detailed discussion of shared model components is given in the supplementary material of Knutti et al. (2013).

Similarities in diagnostic output are not always predictable from a consideration of model construction alone. One can find examples of cases with significant changes in code base but with minor changes in diagnostic similarity. For example, CCSM4 and CESM1(CAM5) have significantly different aerosol schemes, dynamics, and cloud microphysics, and yet our results show the two models as very strongly related when considering the distribution of intermodel distances. This indicates that tuning strategies and nonatmospheric components may play a significant role in diagnostic model similarity, even when primarily atmospheric output is used to assess intermodel distance. This implies that, although the diagnostic output is a useful indicator of

model similarities, those similarities may not be a function of shared code alone. The climateprediction.net (Stainforth et al. 2005) and Quantifying Uncertainty in Model Predictions (QUMP; Murphy et al. 2007) experiments, for example, show that considerable diversity in model behavior is achievable through parameter perturbation alone with an identical codebase.

There are several possible additional factors that might influence diagnostic similarity. First, the tendency for various generations of models from a single institution to exhibit strong similarities in spite of extensive model component changes [see Fig. 2 in Sanderson and Knutti (2012) with reference to NSF–DOE–NCAR CESM, GFDL, or Hadley Centre models] indicates that some elements of model calibration tend to cluster models from a given modeling center. The reasons for this clustering have multiple possible candidates that could lie in institutional policy or regional focus (institutions might be more concerned with their model's performance in the region's climate). Standard metrics used to judge model performance during the model development process or preferred observational datasets may also vary from institution to institution. Second, models rarely change all components at the same time, so we would posit that evaluating when a model is “new” is a subjective matter. Finally, the CMIP5 protocol allows for some flexibility in the way that models implement external forcings, so different groups, even with identical models, can choose to represent the historical and future boundary conditions in different ways to produce differences in the simulated climate. Knutti et al. (2013) see similar relationships in control simulations, but one cannot exclude the possibility that the control simulations themselves might also include common assumptions on boundary conditions.

In summary, we confirm earlier arguments that models are not independent, some are essentially duplicates, and the effective number of independent models based on this method is less than half of the actual number of models, consistent with earlier studies (Jun et al. 2008; Annan and Hargreaves 2011; Sanderson and Knutti 2012). Some models are closer to observations than others (Gleckler et al. 2008; Knutti and Sedáček 2013). We believe that our method and results do not strongly hinge on the way in which one interprets the ensemble as “truth centered” (Knutti 2010), “indistinguishable from truth” (Annan and Hargreaves 2011; Rougier et al. 2013), or neither (Sanderson and Knutti 2012; Bishop and Abramowitz 2013). One could imagine a hypothetical ensemble following any of these frameworks; by duplicating some of its members, bias would be introduced in the ensemble distribution. By evaluating our ensemble subset performance in terms of

ensemble mean performance, we do not necessarily advocate a truth centered ensemble, as the ensemble mean would also be the best estimate of future change in the indistinguishable case.

There are of course different ways to account for model performance and interdependence. In the companion paper (Sanderson et al. 2015), we proposed a method to produce probabilistic estimates that are largely insensitive to model duplicates and can consider model performance. However, when high-dimensional data and/or spatially and temporally consistent fields are required (e.g., for impact models), a fully probabilistic method becomes unwieldy and might even hinder the development of tractable impact analyses (Dessai and Hulme 2004). Bishop and Abramowitz (2013) also propose an alternative technique where models in the archive are subject to a linear transformation, where the weighted mean of transformed models is calculated to be optimally close to an observed climate. This transformation and weighting can then be extrapolated for future projections. This method has the advantage that the resulting transformed models have independent errors and weight future projections by climatological skill. However, the transformed models are not themselves physically self-consistent and there is a potential for simulations to be overfitted to historical data in a manner that could potentially result in overconfident future projections. In comparison, the method we present here preserves a subset of self-consistent physical models (for both present-day and future projections); although they might not be independent in the strict sense of orthogonality, this subset can be simply used for almost any application or analysis.

We thus propose that there is significant utility in spanning the potential uncertainty in future climate by representing spread with an appropriate subset of models. This study introduces weights that assess model uniqueness and model climatology fidelity. We find that the two were inversely related such that the models with the best simulations of the present-day climate were also least unique. A part of this is possibly due to the fact that models have been calibrated by the observations and will thus appear to cluster around those observations (and each other). However, a closer examination reveals that a large fraction of the high-scoring models' lack of uniqueness can be explained by other models that have duplicated some or all of their code. When these duplicates are removed, this strong inverse relationship is weakened (but not entirely eliminated).

This property of the ensemble is clearly, to some extent, contingent on the choice of metrics used, but it does raise a potentially interesting property of the ensemble; the best performing models might also be the most

promiscuous. This situation implies that the ensemble as a whole is already strongly weighted toward the better performing models. We show that, if the models are weighted to reward their uniqueness, then the RMSE of the ensemble mean is increased. Thus, through a mechanism of quasi-natural selection, the climate community has created an ensemble of models that has already upweighted its climatologically best performing members. In other words, relying on model democracy is to some degree upweighting skilled model structures without deliberately thinking about it or discussing it, by the mechanism of duplication of well-proven code.

This could be seen as an argument in support of keeping the entire ensemble when performing an analysis and at least some justification that the multimodel mean result is a defensible best estimate. However, it is at best an accidental property that is not guaranteed to remain in future ensembles and may not at all be visible for more specific questions or metrics. Whether a model is extensively duplicated is not a pure function of its quality or fidelity. A submodel with open source code and few restrictions on its use is more likely to be utilized by another group than another model with a closed-source policy. However, a model that is jointly used by a large number of groups also has a large development pool invested in improving that model. Duplication within institutions depends also on funding and the available computing resources. One could make the argument that the CMIP5 ensemble distribution and the social and intellectual landscape of the climate community are surely related but certainly not in any simple fashion.

A question also remains of whether the original CMIP5 ensemble is sufficient to assess systematic uncertainty in future climate change. This question could easily form a study in itself, but our results are somewhat informative in this matter. First, the number of truly independent models in the archive is significantly less than the number of submitted models, when gauged by model output. Hence, adding another model to the existing archive has most value if the developers introduce novel components and assumptions. It is true that exploring different configurations of existing components through submodel exchange or parameter perturbation can certainly modify model behavior, and we would argue that such experiments should continue in order to fully explore the inherent uncertainties in the existing model set.

However, this uncertainty is conditional on the number of independent models available to us and establishing whether the current set is sufficient is a question that might not be a useful, because there is not a convenient space in which systematic model assumptions

can be defined. For example, the current CMIP5 ensemble might have n fundamentally different convection schemes, each with its own advantages and biases, but nobody would argue that this constituted a “full set.” Where there is approximation and parameterization, there are potentially limitless ways to address this. Because nobody can know the behavior of the $(n + 1)$ th model, the question of ensemble adequacy cannot be answered in a strict sense. Within the ensemble we have, we can tractably experiment with subsetting to assess how many models are required to have confidence in the distribution of future climate change formed by the full set, but we can never know if the $(n + 1)$ th model will adopt different assumptions or resolve a new process to place its projection outside of the existing distribution.

We argue that a joint consideration of model similarity and quality metrics allows the researcher to make use of a more quantitatively defensible sample of simulations available in the CMIP archives, either through weighting or by model elimination (in itself, an extreme form of weighting) to produce a best estimate of combined model projections. Our approach for achieving this can be controlled with a small number of subjective but clearly defined parameters, which can potentially mitigate some of the arbitrary sampling issues that arise from relying on model democracy and can be tailored to specific questions by choosing appropriate metrics and datasets.

It should be noted in this discussion that the CMIP5 archive is not a full representation of the uncertainty space for GCM projections. Rather, it is a collection of intended best possible models: the final iterations of their respective tuning processes as model developers calibrate their parameterization choices to best represent the observed climate properties that they find most important, although there may be other acceptable configurations (Mauritsen et al. 2012). Clearly, these choices and targets will vary from model to model, but the fact that there are implicitly a near-infinite number of rejected parameter configurations for each model must be remembered when trying to interpret the significance of the spread of simulations in the archive. In a practical sense, we ignore these rejected configurations because we do not have access to them. In addition, there is some evidence to suggest that the model diversity one can attain by structural changes significantly exceeds that of parameter changes in currently available perturbed parameter ensembles (Yokohata et al. 2013). Nevertheless, it should be remembered that both the CMIP5 ensemble (and by definition our subsets of that ensemble) is already a subset of all possible model configurations that have been chosen by model developers.

There are some cases where we would argue it is essential to eliminate interdependent models, such as when a correlation found in the multimodel ensemble is used as a constraint on a climate parameter [i.e., for climate sensitivity in Fasullo and Trenberth (2012) or for high-latitude surface albedo feedbacks in Hall and Qu (2006)]. The presence of closely related or even identical models in the archive would tend to artificially inflate the significance of any correlation simply because identical models would exhibit similar values for both the predictor and for the unknown quantity (Caldwell et al. 2014). Removing the obvious interdependent models as shown in this study would certainly be better than assessing a correlation based on the entire archive, but a method for achieving this in a strict statistical sense is presented in Caldwell et al. (2014).

There is a danger that, as models improve, the better models have the potential to converge on the true climate state, which might lead to their elimination if interdependent models are removed. We show in Fig. 3 that this is unlikely to be the case for CMIP5, given none of the models lies close enough to the observations to be influenced by the uniqueness weighting. However, one could imagine if a small group of models make a real advance that removes a long-standing systematic bias (e.g., as some models begin to explicitly resolve convection), then it would be necessary to accept a higher level of similarity among the better performing models (i.e., the uniqueness weighting u_w could no longer be independent of the skill weighting u_s).

Proposing a subset of models to consider for a less biased analysis could be seen as overly prescriptive, but our aim is not to focus on the exact set of models that should be used for future studies but rather to establish a framework in which researchers could make their selection based upon metrics that are most relevant to their question. We would argue that, although the collection of models arising from the “ensemble of opportunity” is often seen as sacrosanct, the democratic policy of one model, one vote is no longer a logical one in the increasingly complex family tree of models available to the researcher. A subset of 10–20 models that are reasonably independent and perform well for the criteria that are judged to be relevant is very likely to be more skillful than the full ensemble. Giving equal weight to all models that have completed a simulation of interest is, albeit implicitly, adopting a weighting scheme that rewards model components that are highly replicated. This weighting scheme might fortuitously have the property of rewarding the most skilled components but, we would argue, this property should be demonstrated and the decision how to incorporate it should be made consciously.

Acknowledgments. We acknowledge the World Climate Research Programme's Working Group on Coupled Modeling, which is responsible for CMIP, and we thank the climate modeling groups for producing and making available their model output. For CMIP the U.S. Department of Energy's Program for Climate Model Diagnosis and Inter-Comparison provides coordinating support and led development of software infrastructure in partnership with the Global Organization for Earth System Science Portals. Portions of this study were supported by the Office of Science (BER), U.S. Department of Energy, Cooperative Agreement DE-FC02-97ER62402. We would also like to thank our anonymous reviewers for their extensive and insightful comments.

REFERENCES

- Adler, R., and Coauthors, 2003: The Version-2 Global Precipitation Climatology Project (GPCP) monthly precipitation analysis (1979–present). *J. Hydrometeorol.*, **4**, 1147–1167, doi:10.1175/1525-7541(2003)004<1147:TVGPCP>2.0.CO;2.
- Annan, J., and J. Hargreaves, 2011: Understanding the CMIP3 multimodel ensemble. *J. Climate*, **24**, 4529–4538, doi:10.1175/2011JCLI3873.1.
- Aumann, H. H., and Coauthors, 2003: AIRS/AMSU/HSB on the Aqua mission: Design, science objectives, data products, and processing systems. *IEEE Trans. Geosci. Remote Sens.*, **41**, 253–264, doi:10.1109/TGRS.2002.808356.
- Bishop, C. H., and G. Abramowitz, 2013: Climate model dependence and the replicate Earth paradigm. *Climate Dyn.*, **41** (3–4), 885–900, doi:10.1007/s00382-012-1610-y.
- Brohan, P., J. Kennedy, I. Harris, S. Tett, and P. Jones, 2006: Uncertainty estimates in regional and global observed temperature changes: A new dataset from 1850. *J. Geophys. Res.*, **111**, D12106, doi:10.1029/2005JD006548.
- Caldwell, P. M., C. S. Bretherton, M. D. Zelinka, S. A. Klein, B. D. Santer, and B. M. Sanderson, 2014: Statistical significance of climate sensitivity predictors obtained by data mining. *Geophys. Res. Lett.*, **41**, 1803–1808, doi:10.1002/2014GL059205.
- Dessai, S., and M. Hulme, 2004: Does climate adaptation policy need probabilities? *Climate Policy*, **4**, 107–128, doi:10.1080/14693062.2004.9685515.
- Evans, J. P., F. Ji, G. Abramowitz, and M. Ekström, 2013: Optimally choosing small ensemble members to produce robust climate simulations. *Environ. Res. Lett.*, **8**, 044050, doi:10.1088/1748-9326/8/4/044050.
- Fasullo, J. T., and K. E. Trenberth, 2012: A less cloudy future: The role of subtropical subsidence in climate sensitivity. *Science*, **338**, 792–794, doi:10.1126/science.1227465.
- Gleckler, P., K. Taylor, and C. Doutriaux, 2008: Performance metrics for climate models. *J. Geophys. Res.*, **113**, D06104, doi:10.1029/2007JD008972.
- Hall, A., and X. Qu, 2006: Using the current seasonal cycle to constrain snow albedo feedback in future climate change. *Geophys. Res. Lett.*, **33**, L03502, doi:10.1029/2005GL025127.
- Jeffrey, M., L. Rotstayn, M. Collier, S. Dravitzki, C. Hamalainen, C. Moeseneder, K. Wong, and J. Syktus, 2013: Australia's CMIP5 submission using the CSIRO-Mk3.6 model. *Aust. Meteor. Oceanogr. J.*, **63**, 1–13.
- Jun, M., R. Knutti, and D. W. Nychka, 2008: Local eigenvalue analysis of CMIP3 climate model errors. *Tellus*, **60A**, 992–1000, doi:10.1111/j.1600-0870.2008.00356.x.
- Knutti, R., 2010: The end of model democracy? *Climatic Change*, **102** (3–4), 395–404, doi:10.1007/s10584-010-9800-2.
- , and J. Sedáček, 2013: Robustness and uncertainties in the new CMIP5 climate model projections. *Nat. Climate Change*, **3**, 369–373, doi:10.1038/nclimate1716.
- , D. Masson, and A. Gettelman, 2013: Climate model genealogy: Generation CMIP5 and how we got there. *Geophys. Res. Lett.*, **40**, 1194–1199, doi:10.1002/grl.50256.
- Li, L., and Coauthors, 2013: The Flexible Global Ocean-Atmosphere-Land System model, grid-point version 2: FGOALS-g2. *Adv. Atmos. Sci.*, **30**, 543–560, doi:10.1007/s00376-012-2140-6.
- Masson, D., and R. Knutti, 2011: Climate model genealogy. *Geophys. Res. Lett.*, **38**, L08703, doi:10.1029/2011GL046864.
- , and —, 2013: Predictor screening, calibration, and observational constraints in climate model ensembles: An illustration using climate sensitivity. *J. Climate*, **26**, 887–898, doi:10.1175/JCLI-D-11-00540.1.
- Mauritsen, T., and Coauthors, 2012: Tuning the climate of a global model. *J. Adv. Model. Earth Syst.*, **4**, M00A01, doi:10.1029/2012MS000154.
- Murphy, J., B. Booth, M. Collins, G. Harris, D. Sexton, and M. Webb, 2007: A methodology for probabilistic predictions of regional climate change from perturbed physics ensembles. *Philos. Trans.*, **365A**, 1993–2028, doi:10.1098/rsta.2007.2077.
- NASA, 2011: CERES EBAF datasets, Langley Research Center, accessed 2011. [Available online at <http://ceres.larc.nasa.gov/products.php?product=EBAF-TOA>.]
- Pennell, C., and T. Reichler, 2011: On the effective number of climate models. *J. Climate*, **24**, 2358–2367, doi:10.1175/2010JCLI3814.1.
- Qu, X., and A. Hall, 2013: On the persistent spread in snow-albedo feedback. *Climate Dyn.*, **42**, 69–81, doi:10.1007/s00382-013-1774-0.
- Rougier, J., M. Goldstein, and L. House, 2013: Second-order exchangeability analysis for multimodel ensembles. *J. Amer. Stat. Assoc.*, **108**, 852–863, doi:10.1080/01621459.2013.802963.
- Sanderson, B. M., and R. Knutti, 2012: On the interpretation of constrained climate model ensembles. *Geophys. Res. Lett.*, **39**, L16708, doi:10.1029/2012GL052665.
- , —, and P. Caldwell, 2015: Addressing interdependency in a multi-model ensemble by interpolation of model properties. *J. Climate*, doi:10.1175/JCLI-D-14-00361.1, in press.
- Santer, B., and Coauthors, 2009: Incorporating model quality information in climate change detection and attribution studies. *Proc. Natl. Acad. Sci. USA*, **106**, 14 778–14 783, doi:10.1073/pnas.0901736106.
- Stainforth, D. A., and Coauthors, 2005: Uncertainty in predictions of the climate response to rising levels of greenhouse gases. *Nature*, **433**, 403–406, doi:10.1038/nature03301.
- Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of CMIP5 and the experiment design. *Bull. Amer. Meteor. Soc.*, **93**, 485–498, doi:10.1175/BAMS-D-11-00094.1.
- Tebaldi, C., and R. Knutti, 2007: The use of the multi-model ensemble in probabilistic climate projections. *Philos. Trans. Roy. Soc.*, **365A**, 2053–2075, doi:10.1098/rsta.2007.2076.

- , J. M. Arblaster, and R. Knutti, 2011: Mapping model agreement on future climate projections. *Geophys. Res. Lett.*, **38**, L23701, doi:[10.1029/2011GL049863](https://doi.org/10.1029/2011GL049863).
- Volodin, E. M., N. A. Dianskii, A. V. Gusev, 2010: Simulating present-day climate with the INMCM4.0 coupled model of the atmospheric and oceanic general circulations. *Izv. Atmos. Oceanic Phys.*, **46**, 414–431, doi:[10.1134/S000143381004002X](https://doi.org/10.1134/S000143381004002X).
- Watanabe, M., and Coauthors, 2010: Improved climate simulation by MIROC5: Mean states, variability, and climate sensitivity. *J. Climate*, **23**, 6312–6335, doi:[10.1175/2010JCLI3679.1](https://doi.org/10.1175/2010JCLI3679.1).
- Wu, T., and Coauthors, 2014: An overview of BCC climate system model development and application for climate change studies. *J. Meteor. Res.*, **28**, 34–56, doi:[10.1007/s13351-014-3041-7](https://doi.org/10.1007/s13351-014-3041-7).
- Yang, D., and O. A. Saenko, 2012: Ocean heat transport and its projected change in CanESM2. *J. Climate*, **25**, 8148–8163, doi:[10.1175/JCLI-D-11-00715.1](https://doi.org/10.1175/JCLI-D-11-00715.1).
- Yokohata, T., and Coauthors, 2013: Reliability and importance of structural diversity of climate model ensembles. *Climate Dyn.*, **41** (9–10), 2745–2763, doi:[10.1007/s00382-013-1733-9](https://doi.org/10.1007/s00382-013-1733-9).