

Lawrence Berkeley National Laboratory

Recent Work

Title

Incorporating model quality information in climate change detection and attribution studies

Permalink

<https://escholarship.org/uc/item/87g5z7x6>

Journal

Proceedings of the National Academy of Sciences of the United States of America, 106(35)

ISSN

0027-8424

Authors

Santer, BD
Taylor, KE
Gleckler, PJ
[et al.](#)

Publication Date

2009-09-01

DOI

10.1073/pnas.0901736106

Peer reviewed

Incorporating model quality information in climate change detection and attribution studies

B. D. Santer^{a,1}, K. E. Taylor^a, P. J. Gleckler^a, C. Bonfils^a, T. P. Barnett^b, D. W. Pierce^b, T. M. L. Wigley^c, C. Mears^d, F. J. Wentz^d, W. Brüggemann^e, N. P. Gillett^f, S. A. Klein^a, S. Solomon^g, P. A. Stott^h, and M. F. Wehnerⁱ

^aProgram for Climate Model Diagnosis and Intercomparison, Lawrence Livermore National Laboratory, Livermore, CA 94550; ^bScripps Institution of Oceanography, La Jolla, CA 92037; ^cNational Center for Atmospheric Research, Boulder, CO 80307; ^dRemote Sensing Systems, Santa Rosa, CA 95401; ^eInstitut für Unternehmensforschung, Universität Hamburg, 20146 Hamburg, Germany; ^fCanadian Centre for Climate Modelling and Analysis, University of Victoria, Victoria, BC, Canada V8W 3V6; ^gChemical Sciences Division, National Oceanic and Atmospheric Administration Earth System Research Laboratory, Boulder, CO 80305; ^hHadley Centre, U.K. Meteorological Office, Exeter EX1 3PB, United Kingdom; and ⁱLawrence Berkeley National Laboratory, Berkeley, CA 94720

Edited by Michael E. Mann, Pennsylvania State University, University Park, PA, and accepted by the Editorial Board July 1, 2009 (received for review February 23, 2009)

In a recent multimodel detection and attribution (D&A) study using the pooled results from 22 different climate models, the simulated “fingerprint” pattern of anthropogenically caused changes in water vapor was identifiable with high statistical confidence in satellite data. Each model received equal weight in the D&A analysis, despite large differences in the skill with which they simulate key aspects of observed climate. Here, we examine whether water vapor D&A results are sensitive to model quality. The “top 10” and “bottom 10” models are selected with three different sets of skill measures and two different ranking approaches. The entire D&A analysis is then repeated with each of these different sets of more or less skillful models. Our performance metrics include the ability to simulate the mean state, the annual cycle, and the variability associated with El Niño. We find that estimates of an anthropogenic water vapor fingerprint are insensitive to current model uncertainties, and are governed by basic physical processes that are well-represented in climate models. Because the fingerprint is both robust to current model uncertainties and dissimilar to the dominant noise patterns, our ability to identify an anthropogenic influence on observed multidecadal changes in water vapor is not affected by “screening” based on model quality.

climate modeling | multimodel database | water vapor

Since the mid-1990s, pattern-based “fingerprint” studies have been the primary and most rigorous tool for disentangling the complex causes of recent climate change (1–3). Fingerprinting relies on numerical models of the climate system to provide estimates of both the searched-for fingerprint—the pattern of climate response to a change in one or more forcing mechanisms—and the background “noise” of natural internal climate variability. To date, most formal detection and attribution (D&A) work has used information from only one or two individual models to estimate both the fingerprint and noise (4–6). Relatively few D&A studies have used climate data from three or more models (7–13).

The availability of large, multimodel archives of climate model output has had important implications for D&A research. A prominent example of such an archive is the CMIP-3 (Coupled Model Intercomparison Project) database, which was a key resource for the Fourth Assessment Report of the Intergovernmental Panel on Climate Change (IPCC) (14). The CMIP-3 archive enables D&A practitioners to use information from two dozen of the world’s major climate models and to examine the robustness of D&A results to current uncertainties in model-based estimates of climate-change signals and natural variability noise (10–13).

Multimodel databases offer both scientific opportunities and challenges. One challenge is to determine whether the information from each individual model in the database is equally reliable, and should be given equal “weight” in a multimodel D&A study, or in estimating some “model average” projection of future climate change (15). Previous multimodel D&A investigations with atmo-

spheric water vapor (10) and sea-surface temperatures (SSTs) in hurricane formation regions (13) adopted a “one model, one vote” approach, with no attempt made to weight or screen models based on their performance in simulating aspects of observed climate. An important and hitherto unexplored question, therefore, is whether the findings of such multimodel D&A studies are sensitive to model weighting or screening decisions.

To address this question, objective measures of model performance are required. An obvious difficulty is that model errors are highly complex; they depend on the variable considered, the space and timescale of interest, the statistical metric used to compare modeled and observed climatic fields, the exact property of the fields that is being considered (e.g., mean state, diurnal or annual cycle, amplitude and structure of variability, and evolution of patterns), and uncertainties in the observations themselves (16–22). Recent assessments of the overall performance of CMIP-3 models have relied on a variety of statistical metrics and were primarily focused on how well these models reproduce the observed climatological mean state (23, 24).*

Here, we revisit our multimodel D&A study with atmospheric water vapor over oceans (10). We calculate a number of different “model quality” metrics and demonstrate that use of this information to screen models does not affect our ability to identify an externally forced fingerprint in satellite data.

Observational and Model Water Vapor Data

We rely on observational water vapor data from the satellite-based Special Sensor Microwave Imager (SSM/I). The SSM/I atmospheric moisture retrievals commenced in late 1987 and are based on measurements of microwave emissions from the 22-GHz water vapor absorption line (25–27). Retrievals are unavailable over the highly emissive land surface and sea-ice

Author contributions: B.D.S., K.E.T., P.J.G., C.B., T.P.B., D.W.P., T.M.L.W., C.M., F.J.W., and W.B. designed research; B.D.S., P.J.G., and C.B. performed research; C.M., F.J.W., P.A.S., and M.F.W. contributed new reagents/analytic tools; B.D.S., K.E.T., P.J.G., C.B., T.M.L.W., N.P.G., S.A.K., and S.S. analyzed data; and B.D.S., K.E.T., P.J.G., C.B., T.M.L.W., C.M., W.B., N.P.G., S.A.K., S.S., P.A.S., and M.F.W. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission. M.E.M. is a guest editor invited by the Editorial Board. Freely available online through the PNAS open access option.

¹To whom correspondence should be addressed. Email: santer1@llnl.gov.

This article contains supporting information online at www.pnas.org/cgi/content/full/0901736106/DCSupplemental.

*The processes affecting the gradual response of the climate system to long-term anthropogenic forcing need not be the same as those controlling shorter-timescale phenomena. For example, model inadequacies in simulating the diurnal cycle do not necessarily translate to a deficient simulation of long-term responses.

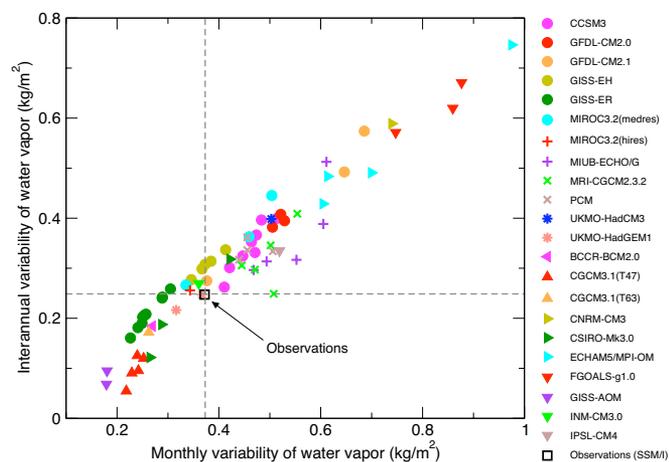


Fig. 1. Comparison of the simulated and observed temporal variability of atmospheric water vapor. Observations are from the SSM/I dataset (25, 26); model data are from 71 realizations of 20th century climate change performed with 22 different models (see *SI Appendix*). All variability calculations rely on monthly mean values of $\langle W \rangle$, the spatial average of total atmospheric moisture over near-global oceans. Model and observational $\langle W \rangle$ data were first expressed as anomalies relative to climatological monthly means over the period 1988–1999 and then linearly detrended. We computed temporal standard deviations from both the unfiltered and filtered anomaly data. The latter were smoothed by using a filter with a half-power point at ≈ 2 years. The raw and filtered standard deviations provide information on monthly and interannual-timescale variability, respectively. All calculations were over the 144-month period from January 1988 to December 1999 (the period of maximum overlap between the SSM/I data and most 20CEN simulations). The dashed gray lines are centered on the observations.

regions. Our focus is therefore on W , the total column water vapor over oceans for a near-global domain.[†]

As noted above, fingerprint studies require estimates of both the climate-change signal in response to external forcing and the noise of internal climate variability. We obtain signal estimates from simulations with historical changes in natural and anthropogenic forcings (“20CEN” runs) and noise information from control integrations with no forcing changes.[‡] We use 20CEN and control integrations from 22 different climate models in the CMIP-3 archive. These are the same models that were used in our original water vapor D&A study (10).

Strategy for Assessment of Model Quality

Fig. 1 illustrates why it may be useful to include model quality information in multimodel D&A studies. The figure shows the simulated and observed temporal standard deviation of $\langle W \rangle$, the spatial average of atmospheric water vapor over near-global oceans.[§] Results are given for monthly and interannual-timescale fluctuations in $\langle W \rangle$. On both timescales, the simulated variability in 20CEN runs ranges from one-third to two-and-a-half times the amplitude of the observed variability.

Are such variability differences between models and observations of practical importance in multimodel D&A studies? Most

[†]Our D&A study area encompasses all oceans between 50°N and 50°S. This domain was chosen to minimize the effect of model-versus-SSM/I water vapor differences associated with inaccurate simulation of the latitudinal extent of ice margins.

[‡]The external forcings imposed in the 20CEN experiments differed between modeling groups. The most comprehensive experiments included changes in both natural external forcings (solar irradiance and volcanic dust loadings in the atmosphere) and in a wide variety of anthropogenic influences (such as well-mixed greenhouse gases, ozone, sulfate and black carbon aerosols, and land surface properties). Details of the models, 20CEN experiments, and control integrations are given in the *SI Appendix*.

[§]Here and subsequently, $\langle \rangle$ denotes a spatial mean.

D&A studies routinely apply some form of statistical test to check the consistency between observed residual variability (after removal of an estimated externally forced signal) and model control run variability (4, 7–13), and many studies compare power spectra of the observed and modeled variables being analyzed (12, 13). Our focus here is not on formal statistical tests or spectral density comparisons; instead, it is calculating metrics that provide more direct information regarding the fidelity with which models simulate the amplitude and structure of key modes of natural internal variability.

Although our D&A study involves water vapor only, we compute performance metrics for water vapor and SST. We examine SST data because observed SST datasets are 130–150 yr in length and therefore provide a better constraint on model-based estimates of decadal variability than the short (21-yr) SSM/I record. Information on low-frequency variability is crucial for D&A applications, because it constitutes the background noise against which we attempt to identify a slowly evolving anthropogenic signal. All SST-based model quality metrics were calculated using observations from the NOAA Extended Reconstructed SST (ERSST) dataset (28).

We evaluate model performance in simulating W and SST in five different regions. The first is the 50°N–50°S ocean domain used in our previous water vapor D&A work. The next three regions were chosen because they provide information on model errors in simulating three characteristic modes of natural climate variability: the El Niño/Southern Oscillation (ENSO), the Pacific Decadal Oscillation (PDO), and the Atlantic Multidecadal Oscillation (AMO).[¶] The final region comprises tropical oceans (30°N–30°S) and is of interest because of claims that modeled and observed atmospheric temperature changes differ significantly in the tropics.

We analyze model performance in simulating the mean state, annual cycle, and amplitude and structure of variability.^{||} There are 10 mean state diagnostics (two variables \times five regions). Each mean state metric is simply a measure of the absolute value of the climatological annual-mean model bias. The 10 annual cycle diagnostics involve the correlations between the simulated and observed climatological mean annual cycle patterns. The 50 variability metrics** are measures of model skill in simulating the amplitude and pattern of observed variability on monthly, interannual, and decadal timescales. The rationale for examining model performance on different timescales is that model variability errors are complex and frequency-dependent (29).

All 70 metrics are normalized by the intermodel standard deviation of the statistical property being considered. This allows us to combine information from the mean, annual cycle, and variability metrics as well as from different climate variables and geographical regions. Details regarding the definition and calculation of our model performance metrics are given in the [supporting information \(SI Appendix\)](#).

Results from Model Quality Assessment

Results for 40 of the 70 individual metrics are shown in Fig. 2. To illustrate the complexity of model errors, we use the example of the

[¶]ENSO variability can be characterized in a number of different ways. We analyze water vapor and SST changes over the Niño 3.4 region (5°N–5°S; 170°W–120°W). The PDO and AMO regions used here are 20°N–60°N; 115°W–115°E and 20°N–60°N; 75°W–0°, respectively.

^{||}We do not calculate metrics that gauge model performance in simulating observed water vapor and SST trends. Results could be biased toward identification of an anthropogenic fingerprint by first selecting a subset of models with greater skill in replicating observed trends and then using the same subset in a D&A analysis that compares modeled and observed trend behavior.

**For the higher-frequency variability comparisons, there are a total of 40 metrics: two variables (SST and W) \times five regions (oceans 50°N–50°S, ENSO, PDO, and AMO regions, and tropical oceans) \times two statistical attributes (variability amplitude and pattern) \times two timescales (monthly and interannual). For comparisons of decadal variability, there are only 10 diagnostics, because these are meaningful to compute for SST only (see the text). All variability pattern metrics are centered correlations, with removal of the spatial means of the two fields being compared.

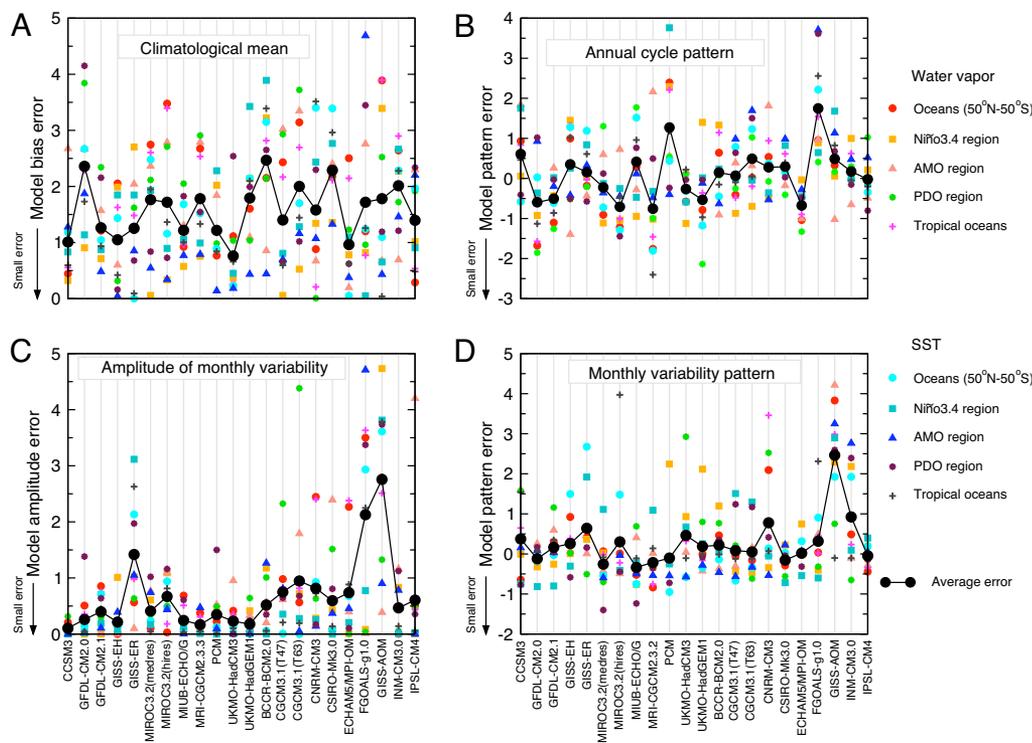


Fig. 2. Results for four different sets of metrics used in the ranking of model performance. The statistics are measures of how well 22 of the models in the CMIP-3 database reproduce key features of observed water vapor and SST behavior in five different geographical regions. The metrics shown here are a subset of the full suite of metrics that we applied for model ranking, and are for the mean state (A), annual cycle pattern (B), amplitude of monthly variability (C), and pattern of monthly variability (D). For models with multiple 20CEN realizations, values of metrics are averaged over realizations. The black dots labeled “average error” represent the arithmetic average (for each model) of the 10 metric values (2 variables \times 5 regions). (A and C) Small values of the normalized metrics indicate greater skill in simulating the mean state and the amplitude of monthly variability. (B and D) Negative values of the normalized pattern correlation metrics denote greater skill in simulating the annual cycle and monthly variability patterns.

UKMO-HadCM3 model (developed at the UK Meteorological Office Hadley Centre). Consider first the results for the absolute bias in the climatological mean state (Fig. 2A). HadCM3 has relatively small bias values for both water vapor and SST, except for SSTs in the PDO region. When models are ranked parametrically on the basis of the “average error” results in Fig. 2A, HadCM3 has the lowest bias values and is therefore ranked first.

In terms of its simulation of the climatological annual cycle pattern (Fig. 2B) and the amplitude of monthly variability (Fig. 2C), HadCM3 also performs well relative to its peers and is ranked seventh and fifth, respectively. For the monthly variability pattern, however, HadCM3 has a large error for water vapor in the PDO region (Fig. 2D). This one component has a marked influence on HadCM3’s low overall ranking (18th) for the monthly variability pattern. For interannual and decadal variability (not shown), HadCM3 ranks 10th and first in terms of its variability amplitude and 15th and 14th in terms of its variability pattern. As is clear from the HadCM3 example and the other model results in Fig. 2,

assessments of the relative skill of the CMIP-3 models are sensitive to a variety of analyst choices.

This message is reinforced in Fig. 3, which shows that for our selected variables, regions, and diagnostics, there are no statistically significant relationships between model skill in simulating the climatological mean state and model skill in capturing either the observed annual cycle or the amplitude and pattern of monthly variability. Similar findings have been obtained in related studies (19, 20, 22, 23). One possible interpretation of this result is that the spatial averages of observed climatological annual means provide a relatively weak constraint on overall model performance. Modeling groups attempt to reduce biases in these large-scale climatological averages by adjusting poorly known physical parameters (and by flux correction, which still is used in several of the CMIP-3 models). Observed annual cycle and variability patterns offer more stringent tests of model performance. Reliable reproduction of these more challenging observational targets is difficult to achieve through tuning alone—accurate representation of the underlying physics is of greater importance.

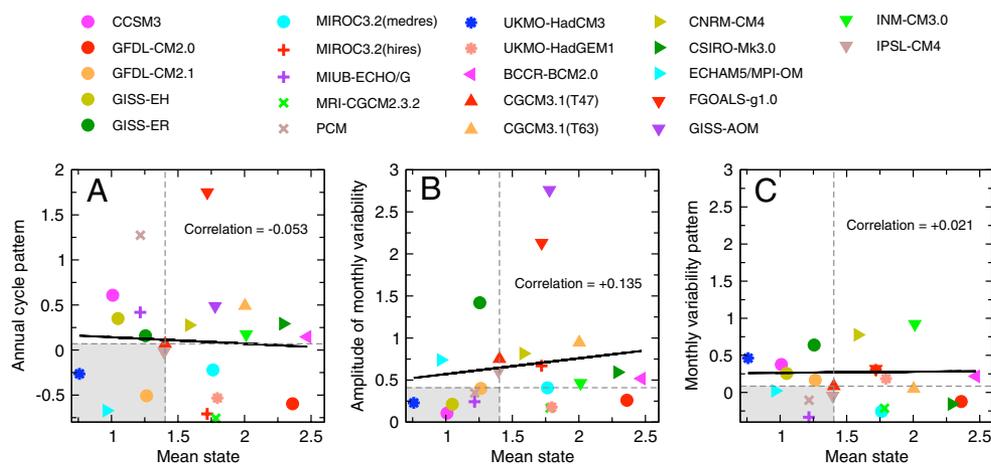


Fig. 3. Relationship between model skill in simulating the mean state and skill in simulating the annual cycle pattern (A), amplitude of monthly variability (B), and monthly variability pattern (C). Results plotted are the “average errors” shown and described in Fig. 2. The black lines are the fitted least-squares regression lines. Models to the left of the vertical dashed gray line are ranked in the top 10 based on values of the mean state metric $\hat{\alpha}$. Models below the horizontal dashed gray line are ranked in the top 10 based on values of the annual cycle pattern metric $\hat{\beta}$ (A), the variability amplitude metric $\hat{\phi}$ (B), and the variability pattern metric $\hat{\psi}$ (C). The gray shaded region indicates the intersection of the two sets of top 10 models plotted in each graph.

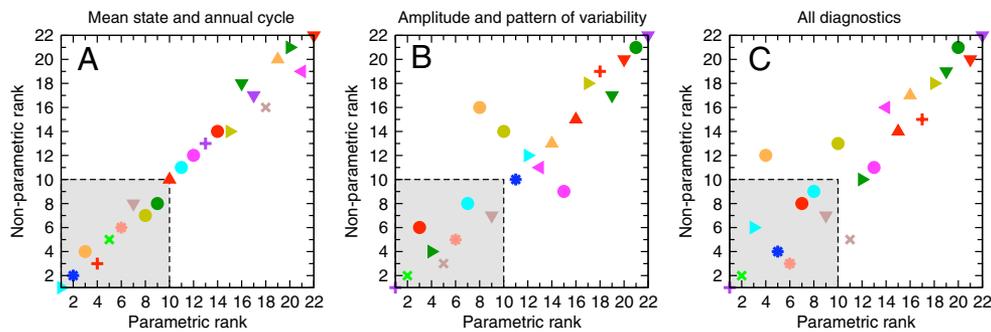


Fig. 4. Parametric and nonparametric ranking of 22 CMIP-3 models. The parametric ranking is based on the \hat{Q}_1 , \hat{Q}_2 , and \hat{Q}_3 statistics, which are (respectively) measures of model skill in simulating the observed mean state and annual cycle (A), the amplitude and pattern of variability (B), and the combined mean state, annual cycle, and variability properties (C). The \hat{Q}_1 , \hat{Q}_2 , and \hat{Q}_3 statistics are averages of the normalized values of 20 mean state and annual cycle metrics (M+AC), 50 variability amplitude and variability pattern metrics (VA+VP), and 70 combined metrics (ALL). In the nonparametric ranking procedure, models are ranked from 1 to 22 for each of the 70 metrics, and the individual ranks are then averaged in each of the three groups of metrics (M+AC, VA+VP, and ALL). Full details of the statistics and ranking procedures are given in the *SI Appendix*. The gray shaded boxes indicate the intersection of the two sets of top 10 models. See Fig. 3 for key.

The final stage in our model quality assessment is to combine information from different performance metrics, which we accomplish in three different ways. The three combinations involve the 10 mean state and 10 annual cycle diagnostics (M+AC), the 25 variability amplitude and 25 variability pattern metrics (VA+VP), and the 70 mean state, annual cycle, and variability diagnostics (ALL). Individual values of these metrics are averaged, yielding the \hat{Q}_1 , \hat{Q}_2 , and \hat{Q}_3 statistics, which are used for the parametric ranking of the CMIP-3 models (see *SI Appendix*). The nonparametric rank is simply the average of the individual ranks rather than the average of individual metric values.

The overall ranking results are shown in Fig. 4. A number of interesting features are evident. First, only three models (MRI-CGCM2.3.2, UKMO-HadGEM1, and IPSL-CM4) are consistently ranked within the top 10 CMIP-3 models based on both ranking approaches and all three sets of performance criteria (M+AC, VA+VP, and ALL). None of the top four models determined with the M+AC metrics (Fig. 4A) is also in the top four based on the VA+VP metrics (Fig. 4B). These results support our previous finding that assessments of model quality are sensitive to the choice of statistical properties used in model evaluation.

Second, there is also some sensitivity to the choice of ranking procedure, particularly for the VA+VP and ALL statistics (Fig. 4B and C). In each of these two cases, the nonparametric and parametric ranking approaches identify slightly different sets of “top 10” models. Only 8 models are in the intersection of these sets.

Third, higher horizontal resolution does not invariably lead to improved model performance. The CMIP-3 archive contains two models (the Canadian Climate Centre’s CGCM3.1 and the Japanese MIROC3.2) that were run in both higher- and lower-resolution configurations. The lower-resolution version of CGCM3.1 outperforms the higher-resolution version in terms of the M+AC diagnostics but not for the VA+VP metrics. The reverse applies to the MIROC3.2 model. The lack of a consistent benefit of higher resolution is partly due to our focus on temperature and moisture changes over oceans. The performance improvement related to higher resolution is more evident over land areas with complex topography (30).

Detection and Attribution Analysis

We now apply the same multimodel D&A method used by Santer et al. (10). Instead of employing all 22 CMIP-3 models in the D&A analysis, we restrict our attention to 10-member subsets of the 22 models. These subsets are determined by ranking models on the basis of the three different sets of metrics (M+AC, VA+VP, and ALL) and two different ranking approaches (parametric and

nonparametric). From each of these six ranking sets, we select the top 10 and bottom 10 models, yielding 12 groups of 10 models.

Fingerprints are calculated in the following way. For each set of 10 models, we determine the multimodel average of the atmospheric moisture changes over the period 1900–1999.^{††} The fingerprint is simply the first empirical orthogonal function (EOF) of the multimodel average changes in water vapor.

Because 10 modeling groups used anthropogenic forcings only, whereas the other 12 applied a combination of anthropogenic and natural external forcings (see *SI Appendix*), we expect the multimodel fingerprint to down-weight the contribution of natural external forcing to the fingerprint. However, previous work has found that the fingerprints estimated from combined historical changes in anthropogenic and natural external forcing are very similar to those obtained from “anthropogenic only” forcing (10). We infer from this that anthropogenic forcing is the dominant influence on the changes in atmospheric moisture over the 20th century and that the multimodel fingerprint patterns are not distorted by the absence of solar and volcanic forcing in 10 of the 22 models analyzed here.^{‡‡}

There is pronounced similarity between the fingerprint patterns estimated from the 12 subsets of CMIP-3 models (Fig. 5). All 12 patterns show spatially coherent water vapor increases, with the largest increases over the warmest ocean areas. There are no systematic differences between the fingerprints estimated from different sets of metrics, different ranking procedures, or from the top 10 or bottom 10 models. This indicates that the structure of the water vapor fingerprint is primarily dictated by the zero-order physics governing the relationship between surface temperature and column-integrated water vapor (25, 31).

For each of our 12 subsets of CMIP-3 models, estimates of natural internal variability are obtained by concatenating the 10 individual control runs of that subset, after first removing residual drift from each control (Fig. S1 in *SI Appendix*). The leading EOF patterns estimated from the concatenated control runs are remarkably similar. Each displays the horseshoe-shaped pattern characteristic of the effects of ENSO variability on atmospheric moisture

^{††}This calculation involves averaging the ensemble mean water vapor changes of each model—i.e., averaging the 20CEN realizations of an individual model before averaging over models (see *SI Appendix*). Note that use of water vapor data for the entire 20th century (rather than simply the period of overlap with SSM/I) provides a less noisy estimate of the true water vapor response to slowly varying external forcings and a response that is more similar across models.

^{‡‡}Because volcanic effects on climate have pronounced structure in space and time, they can and have been identified in D&A studies which include information on the spatiotemporal evolution of signal and noise (12).

fingerprint is governed by very basic physics and is highly similar in all 12 of our sensitivity tests (Fig. 5). Second, the fingerprint is characterized by spatially coherent water vapor increases, whereas the dominant noise modes in the model control runs are ENSO-like in structure and do not show coherent water vapor increases over the entire global ocean (Fig. S2 in *SI Appendix*). Although the structural details of the dominant noise mode differ from model to model (Fig. S3 in *SI Appendix*), the dissimilarity of the water vapor fingerprint and the leading noise patterns does not. This dissimilarity is the main explanation for the robustness of our D&A results.

The water vapor feedback mechanism is of primary importance in determining the sensitivity of the climate system to external forcing (31, 32). Because our fingerprint estimates are robust across models and relatively insensitive to the model quality metrics calculated here, the contribution of water vapor feedback to projected future climate changes may be similarly insensitive to model skill.^{§§}

Our study also demonstrates that it is not easy to make an unambiguous identification of “superior” models, even for a very specific application. Model performance assessments are sensitive to the choice of climate variables, analysis regions and timescales, the physical properties of the fields being compared, the comparison metrics, the way in which individual metrics are normalized and combined, and the ranking approaches (see *SI Appendix*). There is considerable subjectivity in all of these choices. Different sets of choices would yield different model rankings.

In our analysis of water vapor and SST data, we find that model performance in simulating the mean state is virtually uncorrelated with model performance in reproducing the observed annual cycle or the observed amplitude or pattern of variability. This result has

implications for attempts to use model performance metrics to weight projections of future climate change. To date, most of these attempts have relied on mean state metrics. Our findings imply that different projection weights would be obtained with annual cycle and variability metrics. Whether different weighting approaches lead to important differences in climate-change projections is currently unclear and may depend on the region, climate variable, and timescale of interest (20, 22). Identification of the best models for making projections of future climate change will likely require metrics that can better constrain current uncertainties in feedback mechanisms (33).

Although we find that incorporating model quality information has little impact on our ability to identify an externally forced water vapor fingerprint, this does not mean that model quality assessment will be of limited value in D&A studies with other variables (8, 11). In the case of water vapor, S/N ratios are invariably above stipulated significance thresholds. If S/N ratios are closer to these thresholds, it may become more important to screen or down-weight models that are deficient in their simulation of the amplitude and structure of natural variability. As we show here, such variability errors can systematically bias D&A results.

In summary, future multimodel D&A studies must deal with the fundamental challenge of how to make appropriate use of the information from a large collection of models of varying complexity and performance levels. Inevitably, model quality assessment will be an integral component of multimodel D&A studies. Although a democratic “one model, one vote” approach was successful for the water vapor D&A problem, this approach may not be adequate in all cases.

ACKNOWLEDGMENTS. We thank Gabi Hegerl (University of Edinburgh) and an anonymous reviewer for constructive comments on the paper, the modeling groups for providing simulation output for analysis, the Program for Climate Model Diagnosis and Intercomparison for collecting and archiving these data, and the World Climate Research Program’s Working Group on Coupled Modeling for organizing the model data analysis activity. The CMIP-3 multimodel dataset was supported by the Office of Science, U.S. Department of Energy. National Oceanic and Atmospheric Administration ERSST data were provided by Dick Reynolds at the National Climatic Data Center. P.A.S. was supported by the joint Department of Energy and Climate Change/Department for Environment, Food and Rural Affairs (GA01101) and Ministry of Defense Integrated Climate (CBC/2B/0417_Annex C) Program.

^{§§}We note, however, that upper tropospheric water vapor is a key component of the water vapor feedback. Our skill measures address only total column water vapor, which is dominated by water vapor in the lower troposphere. Metrics focusing on model performance in simulating the present-day vertical distribution of water vapor may yield stronger relationships between model skill and the component of climate change projections arising from water vapor feedback.

- Santer BD, et al. (1996) A search for human influences on the thermal structure of the atmosphere. *Nature* 382:39–46.
- Tett SFB, Mitchell JFB, Parker DE, Allen MR (1996) Human influence on the atmospheric vertical temperature structure: Detection and observations. *Science* 274:1170–1173.
- Hegerl GC, et al. (1996) Detecting greenhouse-gas-induced climate change with an optimal fingerprint method. *J Clim* 9:2281–2306.
- Stott PA, et al. (2000) External control of 20th century temperature by natural and anthropogenic forcings. *Science* 290:2133–2137.
- Santer BD, et al. (2003) Contributions of anthropogenic and natural forcing to recent tropopause height changes. *Science* 301:479–483.
- Barnett TP, et al. (2005) Penetration of human-induced warming into the world’s oceans. *Science* 309:284–287.
- Gillett NP, et al. (2002) Detecting anthropogenic influence with a multi-model ensemble. *Geophys Res Lett* 29, 1970, doi:10.1029/2002GL015836.
- Gillett NP, Zwiers FW, Weaver AJ, Stott PA (2003) Detection of human influence on sea level pressure. *Nature* 422:292–294.
- Huntingford C, Stott PA, Allen MR, Lambert FH (2006) Incorporating model uncertainty into attribution of observed temperature change. *Geophys Res Lett* 33:L05710, 10.1029/2005GL024831.
- Santer BD, et al. (2007) Identification of human-induced changes in atmospheric moisture. *Proc Natl Acad Sci USA* 104:15248–15253.
- Zhang X, et al. (2007) Detection of human influence on twentieth-century precipitation trends. *Nature* 448:461–465.
- Hegerl GC, et al. (2007) Understanding and attributing climate change. in *Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change*, Solomon S, Qin D, Manning M, Chen Z, Marquis M, Averyt KB, Tignor M, Miller HL, eds (Cambridge Univ Press, Cambridge, U.K.).
- Gillett NP, Stott PA, Santer BD (2008) Attribution of cyclogenesis region sea surface temperature change to anthropogenic influence. *Geophys Res Lett* 35:L09707, 10.1029/2008GL033670.
- Intergovernmental Panel on Climate Change (2007) Summary for Policymakers. in *Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change*, Solomon S, Qin D, Manning M, Chen Z, Marquis M, Averyt KB, Tignor M, Miller HL, eds (Cambridge Univ Press, Cambridge, U.K.).
- Giorgi F, Mearns LO (2002) Calculation of average, uncertainty range and reliability of regional climate changes from AOGCM simulations via the reliability ensemble averaging (REA) method. *J Clim* 15:1141–1158.
- Preisendorfer RW, Barnett TP (1983) Numerical model-reality intercomparison tests using small-sample statistics. *J Atmos Sci* 40:1884–1896.
- Wigley TML, Santer BD (1990) Statistical comparison of spatial fields in model validation, perturbation, and predictability experiments. *J Geophys Res* 95:851–865.
- Taylor KE (2001) Summarizing multiple aspects of model performance in a single diagram. *J Geophys Res* 106:7183–7192.
- Diffenbaugh NS (2005) Response of large-scale eastern boundary current forcing in the 21st century. *Geophys Res Lett* 32:L19718, 10.1029/2005GL023905.
- Brekke LD, Dettlinger MD, Maurer EP, Anderson M (2008) Significance of model credibility in estimating climate projection distributions for regional hydroclimatological risk assessments. *Clim Change* 89:371–394.
- Waugh DW, Eyring V (2008) Quantitative performance metrics for stratospheric-resolving chemistry-climate models. *Atmos Chem Phys* 8:5699–5713.
- Pierce DW, Barnett TP, Santer BD, Gleckler PJ (2009) Selecting global climate models for regional climate change studies. *Proc Nat Acad Sci USA* 106:8441–8446.
- Gleckler PJ, Taylor KE, Doutriaux C (2008) Performance metrics for climate models. *J Geophys Res* 113:D06104, 10.1029/2007JD008972.
- Reichler T, Kim J (2008) How well do coupled models simulate today’s climate? *Bull Amer Met Soc*, 10.1175/BAMS-89-3-303.
- Wentz FJ, Schabel M (2000) Precise climate monitoring using complementary data sets. *Nature* 403:414–416.
- Mears CA, Wentz FJ, Santer BD, Taylor KE, Wehner MF (2007) Relationship between temperature and precipitable water changes over tropical oceans. *Geophys Res Lett* 34:L24709, 10.1029/2007GL031936.
- Trenberth KE, Fasullo J, Smith L (2005) Trends and variability in column-integrated atmospheric water vapor. *Clim Dyn* 24:741–758.
- Smith TM, Reynolds RW, Peterson TC, Lawrimore J (2008) Improvements to NOAA’s historical merged land-ocean surface temperature analysis (1880–2006). *J Clim* 21:2283–2296.
- AchutaRao K, Sperber KR (2006) ENSO simulation in coupled atmosphere-ocean models: are the current models better? *Clim Dyn* 27:1–15.
- Duffy PB, Govindasamy B, Milovich J, Taylor KE, Thompson S (2003) High resolution simulations of global climate, part 1: Present climate. *Clim Dyn* 21:371–390.
- Soden BJ, Held IM (2006) An assessment of climate feedbacks in coupled ocean-atmosphere models. *J Clim* 19:3354–3360.
- Held IM, Soden BJ (2006) Robust responses of the hydrological cycle to global warming. *J Clim* 19:5686–5699.
- Hall A, Qu X (2006) Using the current seasonal cycle to constrain snow albedo feedback in future climate change. *Geophys Res Lett* 33:L03502, 10.1029/2005GL025127.