

## Improved Combination of Multiple Atmospheric GCM Ensembles for Seasonal Prediction

ANDREW W. ROBERTSON, UPMANU LALL, STEPHEN E. ZEBIAK, AND LISA GODDARD

*International Research Institute for Climate Prediction, The Earth Institute at Columbia University, Palisades, New York*

(Manuscript received 11 December 2003, in final form 24 May 2004)

### ABSTRACT

An improved Bayesian optimal weighting scheme is developed and used to combine six atmospheric general circulation model (GCM) seasonal hindcast ensembles. The approach is based on the prior belief that the forecast probabilities of tercile-category precipitation and near-surface temperature are equal to the climatological ones. The six GCMs are integrated over the 1950–97 period with observed monthly SST prescribed at the lower boundary, with 9–24 ensemble members. The weights of the individual models are determined by maximizing the log likelihood of the combination by season over the integration period. A key ingredient of the scheme is the climatological equal-odds forecast, which is included as one of the “models” in the multimodel combination. Simulation skill is quantified in terms of the cross-validated ranked probability skill score (RPSS) for the three-category probabilistic hindcasts. The individual GCM ensembles, simple poolings of three and six models, and the optimally combined multimodel ensemble are compared.

The Bayesian optimal weighting scheme outperforms the pooled ensemble, which in turn outperforms the individual models. In the extratropics, its main benefit is to bring much of the large area of negative-precipitation RPSS values up to near-zero values. The skill of the optimal combination is almost always increased (in the large spatial averages considered) when the number of models in the combination is increased from three to six, regardless of which models are included in the three-model combination.

Improvements are made to the original Bayesian scheme of Rajagopalan et al. by reducing the dimensionality of the numerical optimization, averaging across data subsamples, and including spatial smoothing of the likelihood function. These modifications are shown to yield increases in cross-validated RPSS skills. The revised scheme appears to be better suited to combining larger sets of models, and, in the future, it should be possible to include statistical models into the weighted ensemble without fundamental difficulty.

### 1. Introduction

Atmospheric general circulation models (GCMs) are now used routinely at several centers as part of a two-tier system for making seasonal climate forecasts up to several seasons in advance (e.g., Goddard et al. 2003). The sea surface temperatures (SSTs) are predicted first, and these are then used as boundary conditions for ensembles of predictions with GCMs. The latter simulate precipitation and temperature and other atmospheric variables, with a resolution of about 300 km across the globe. The two-tier approach approximates the coupled ocean–atmosphere system in which much of the seasonal predictability stems from ocean memory. Two issues confront this system: 1) the difficulty in tier 1 of predicting the SST boundary conditions for use in tier 2, and 2) the optimal use of atmospheric models to simulate seasonal climate. This paper addresses the second issue.

The second tier of the two-tier approach is based on

harnessing atmospheric predictability of the “second kind” (Lorenz 1963), in which the monthly or seasonal-average atmospheric statistical behavior is often sensitive to anomalies in the underlying sea surface and land conditions, with the former being much the stronger effect. In general, atmospheric chaos prevents information in the state of the atmosphere at the initial time of the forecast from being useful at lead times greater than about 2 weeks. Thus, in order to deal with the statistical nature of the problem, forecasts need to be made from ensembles of GCM simulations (typically 10–20 members; Kumar et al. 2001), generated through small perturbations of the initial conditions. These ensembles often differ significantly between one GCM and another because of differences in physical parameterizations between the models. Different GCMs may perform better in different geographical locations, and a combination of models has been shown to outperform a single model globally (Doblas-Reyes et al. 2000).

Several methods exist for combining together the ensemble simulations from multiple GCMs. The simulations or predictions are commonly expressed in terms of three-category probabilities: “below normal,” “near

---

*Corresponding author address:* Andrew W. Robertson, IRI, Monell 230, 61 Route 9W, Palisades, NY 10964.  
E-mail: awr@iri.columbia.edu

TABLE 1. The six GCMs used in the combinations.

Model	ECHAM4.5 <sup>a</sup>	NCEP-MRF9 <sup>b</sup>	NSIPP1 <sup>c</sup>	COLA <sup>d</sup>	CCM3.2 <sup>e</sup>	ECPC <sup>f</sup>
Ensemble size	24	10	9	10	10	10
Horizontal resolution	T42	T40	2° × 2.5°	T63	T42	T62
No. of levels	19	18	34	18	18	28

<sup>a</sup> ECHAM: Max Planck Institute for Meteorology, Hamburg, Germany (Roeckner et al. 1996); <http://www.mpimet.mpg.de/en/extra/models/echam/index.php>.

<sup>b</sup> NCEP-MRF: National Centers for Environmental Prediction Medium-Range Forecast model (Kumar et al. 1996).

<sup>c</sup> NSIPP: National Aeronautics Space Administration's Seasonal-to-Interannual Prediction Project at Goddard Space Flight Center; [http://nsipp.gsfc.nasa.gov/research/atmos\\_descr.html](http://nsipp.gsfc.nasa.gov/research/atmos_descr.html).

<sup>d</sup> COLA: Center for Ocean–Land–Atmosphere studies; [http://www.pcmdi.llnl.gov/modeldoc/amip1/14cola\\_ToC.html](http://www.pcmdi.llnl.gov/modeldoc/amip1/14cola_ToC.html).

<sup>e</sup> CCM: National Center for Atmospheric Research (NCAR) Community Climate Model (Hack et al. 1998); <http://www.cgd.ucar.edu/cms/ccm3>.

<sup>f</sup> ECPC: Experimental Climate Prediction Center at Scripps Institution of Oceanography: a revised version of the GCM earlier implemented at the National Oceanic and Atmospheric Administration (NOAA)/NCEP (Kanamitsu et al. 2002), with some changes to the physics as described in Kanamitsu and Mo (2003).

normal,” and “above normal,” with the terciles computed from a climatological period. The simplest method is to simply “pool” the ensembles of the different models together to form a large superensemble, giving each member equal weight (Hagedorn 2001). To go beyond this, each GCM can be given a weight according to its historical skill. Rajagopalan et al. (2002, hereafter RLZ) introduced a Bayesian methodology to determine the optimal weights by using the equiprobable climatological forecast probabilities as a prior. This method is based on the supposition that seasonal climate predictability is marginal in many areas—even assuming that the SST can be predicted in tier 1 of the forecast—so that a reasonable forecast prior is the climatological three-category probabilities of  $\frac{1}{3}$  for each category. To the extent that a model's historical skill exceeds that of this climatological forecast, its forecast is weighted preferentially in the multimodel combination forecast. Thus our prior belief that the best seasonal forecast is the climatological one is updated by the GCM forecasts according to their skill over the historical record.

The Bayesian scheme was implemented by RLZ for each of the models' land grid boxes independently for precipitation and 2-m temperature separately. The spatial maps of model weights that result often exhibit small-scale variability that may not be physical. No cross-validation was used in that study, and the weights may be sensitive to sampling over the relatively short (41 yr) training period. In addition, the dimensionality of the likelihood optimization in RLZ scales linearly with the number of models. This may be adequate for the three models combined by RLZ but becomes problematic when combining many models together, because of an insufficient length/amount of training data. Despite the success of RLZ's Bayesian scheme, questions remain regarding the usefulness of combining together many models, and whether a simple pooled ensemble might suffice for a larger multimodel ensemble.

The aim of this paper is to use historical ensembles made with six atmospheric GCMs to investigate the skill of multimodel precipitation and near-surface tempera-

ture simulations. The GCM simulations were made with historical analyses of SSTs, and we do not address the issue of the seasonal predictability of SST. We compare a simple pooled ensemble with an improved Bayesian weighting scheme and examine changes in simulation skill when the number of GCMs is increased from three to six.

The set of GCM simulations is described in section 2, along with the observational datasets, the probabilistic (three category) forecast methodology, and the skill measure used for validation. Section 3 describes the optimal weighting methodology and the improvements that are made to the RLZ Bayesian scheme. The skill of the revised optimal combination is presented in section 4 and compared against simply pooling all the GCM ensembles together, as well as the RLZ scheme. The paper's conclusions are presented and discussed in section 5.

## 2. Preliminaries

### a. The general circulation models

This study is based on six GCMs run in ensemble mode (9–24 members) over the period 1950–97, with only the initial conditions differing between ensemble members. The same monthly observational SST dataset was prescribed globally in each case, consisting of the Reynolds (1988) dataset, up until the early 1980s, and the Reynolds and Smith (1994) dataset thereafter. A key to the six GCMs is provided in Table 1 and includes the model resolution and ensemble size.

The results presented below focus on the January–February–March (JFM) and particularly the July–August–September (JAS) seasonal averages of precipitation and 2-m temperature, interpolated (if necessary) to a T42 Gaussian grid (approximately 2.8° in latitude and longitude). Only grid boxes that contain land are considered, yielding 2829 in all.

The observational verification data for both precipitation and near-surface (2 m) air temperature comes

from the New et al. (1999, 2000) 0.5° dataset, compiled by the Climate Research Unit of the University of East Anglia, United Kingdom. The observational datasets were aggregated onto the T42 Gaussian grid of the models.

### b. Probabilistic forecasts and the pooled ensemble

For simplicity, we often refer to the model simulations as “forecasts,” keeping in mind that the observed SSTs were prescribed. Thus we are simulating precipitation and temperature over land, given knowledge of the contemporaneous global distribution of SST. The forecasts are expressed probabilistically by counting how many of the ensemble members fall into the below-normal, near-normal, and above-normal categories. The probabilistic GCM forecast for category  $k$  at time  $t$  is thus expressed as

$$P_{kt}(y) = m_{kt}/m, \quad (1)$$

where  $m$  is the total number of GCM simulations in the ensemble,  $m_{kt}$  is the number of ensemble members falling into category  $k$  at time  $t$ , and  $y$  stands for either seasonal-mean precipitation or temperature.

The terciles are determined for each model (and the observations) separately using the 1968–97 30-yr period as the climate normal. In this way, we remove any overall bias of each model, as expressed in the respective categorical values.

The simplest multimodel ensemble is formed by pooling the ensembles from each model together to produce one large superensemble (after having removed each model’s bias individually as described in the previous paragraph). The forecast probability of an ensemble of  $J$  models is given by

$$P_{kt}^{\text{pool}}(y) = \frac{1}{m_p} \sum_{j=1}^J m_{jkt}, \quad (2)$$

where  $m_p = \sum_{j=1}^J m_j$  is the total number of ensemble members. This is referred to as the pooled ensemble in this paper.

To verify a forecast, we compare  $P_{kt}(y)$  for  $k = 1, 2, 3$  against the category  $k^*$  that was observed to occur at time  $t$ . The ranked probability skill score (RPSS) is used to quantify the skill of the forecasts (Epstein 1969; Wilks 1995). The RPSS is a distance-sensitive measure of the skill of probability forecasts, defined in terms of the squared differences between the *cumulative* probabilities in the forecast and observation vectors. For a single three-category forecast:

$$\text{RPS} = \sum_{l=1}^3 \left( \sum_{k=1}^l y_k - \sum_{k=1}^l o_k \right)^2, \quad (3)$$

where  $o_k = 1$  if category  $k$  was observed to occur ( $o_k = 0$  otherwise), and  $y_k$  are the forecast probabilities. Jointly evaluating a set of  $n$  forecasts (e.g., averaging

over time, space, or both) and expressing the result relative to the climatological probabilities yields

$$\text{RPSS} = 1 - \frac{\sum_{i=1}^n \text{RPS}_i}{\sum_{i=1}^n \text{RPS}_i^{\text{Clim}}}. \quad (4)$$

The RPSS is positive if—on average—the forecast skill exceeds that of the climatological probabilities. Random probability forecasts score worse than climatology and can yield large negative RPSS values, because confident but incorrect forecasts are penalized acutely in (3) (e.g., Goddard et al. 2003; Mason 2004). In this study, the RPSS is computed for the years 1953–95. Cross-validation is used when computing the RPSS by withholding six contiguous years at a time, determining the weights from the remaining 42 yr, and calculating the RPSS for year 4 of the omitted set. This was done so as to leave a training-set length divisible by 3, which is useful if the tercile values are recomputed for each cross-validation sample.

## 3. Optimal model weights

### a. Combining a single model with the climatological prior

The method used here is conceptually Bayesian (e.g., Gelman et al. 1995) and is described fully in RLZ. It is based on the fact that seasonal climate predictability is often marginal, so that a reasonable forecast prior would consist of the climatological probabilities of  $\frac{1}{3}$  for each of three categories. Only to the extent that a particular model or a combination of models shows skill at predicting the quantity of interest over the historical record at a particular location (hindcast skill) do we desire our predictions to deviate from equal odds.

Using the Dirichlet distribution as a conjugate distribution for the multinomial process that is relevant to tercile categories (or quantiles in general), the posterior distribution resulting from the combination of two sources of information (i.e., the climatological forecast plus a single GCM ensemble forecast), with parameters  $\mathbf{a}$  and  $\mathbf{b}$ , is also Dirichlet with parameter  $(\mathbf{a} + \mathbf{b})$ . Here, we consider a weighted combination of the climatological probabilities  $\mathbf{P}_t(x) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ , and the GCM forecast probabilities  $\mathbf{P}_t(y)$  with components  $P_{kt}(y) = m_{kt}/m$ . The posterior distribution of forecast probabilities for year  $t$  can thus be expressed (see RLZ) as the sum

$$f[\mathbf{Q}_t | \mathbf{P}_t(y)] = D(\mathbf{a} + \mathbf{b}), \quad (5)$$

where  $\mathbf{Q}_t$  is a vector of posterior probabilities for each of the categories for year  $t$ . To proceed, we consider only the first moment of the two Dirichlet distributions (i.e., the means), whose sum yields

$$E(Q_{kt}) = \frac{nP_{kt}(x) + wmP_{kt}(y)}{n + wm},$$

$$= \frac{(n/3) + wm_{kt}}{n + wm} \quad \text{for tercile categories,} \quad (6)$$

and  $w$  is the weight to be optimized, which is constrained to be nonnegative.

The uncertainty inherent in estimating the two sample means that are combined in Eq. (6) can be expressed through their respective sample sizes. For the climatological probabilities, there is uncertainty in the estimation of the tercile values (i.e., the break points between the three categories) that depends on the number of years  $n$  in the climatological record (typically  $n = 30$ ). For the GCM probabilities, the sample size is the number of ensemble members  $m$ . This is the reason why  $m$  and  $n$  appear in Eq. (6). The effective sample size of the combined forecast is the weighted sum ( $n + wm$ ) of the sample size of the climatology and the GCM ensemble size (RLZ).

The selection of  $w$  constitutes an optimization problem, the result depending on the choice of skill measure that is to be optimized. Given the Bayesian framework, a natural choice is the posterior likelihood function, defined over the  $N$ -yr common available record of historical data and model simulations at a particular grid location. This has the form

$$L(w) = \prod_{t=1}^N E(Q_{k^*t}), \quad (7)$$

where  $k^*$  represents the category actually observed to occur at each time  $t$ . Thus  $L(w)$  simply reflects the product over all times (years) of the forecast probabilities assigned to the correct category. It represents an integration of the model's performance over a run of events. In practice, it is computationally more accurate to sum over log likelihoods, rather than computing the product of likelihoods over all times. Similar results are obtained by minimizing the sum of squared errors, or by maximizing the RPSS.

#### b. Combining several models

The above scheme was generalized by RLZ to construct a posterior probability forecast through a combination of forecasts from  $J$  different models plus a climatological forecast. The mean of the posterior categorical probability forecast is then defined as

$$E(Q_{kt}) = \frac{\sum_{j=1}^{J+1} w_j m_j P_{jkt}(y)}{\sum_{j=1}^{J+1} w_j m_j}, \quad (8)$$

where  $m_j$  is the size of the ensemble for model  $j$  ( $m_j = n$  for climatology), and  $w_j$  is the weight given to model

$j$ . The weights are determined by maximizing the posterior likelihood function as before [Eq. (7)].

This scheme has been used successfully at the International Research Institute for Climate Prediction (IRI) to make routine seasonal climate forecasts using three–six GCMs (Barnston et al. 2003). However, estimation difficulties started to arise when more models became available and were added to the mix. The resulting weight maps became more noisy and speckled in appearance (cf. RLZ). Upon closer inspection, it was found (not shown) that the weights often become exactly zero for all except one model (or climatology), so that the scheme [Eq. (8)] tends to “choose” one (or two) particular model(s), with large variability in this choice between neighboring model grid boxes. This problem appears to be associated with the high dimension of the optimization space, given the short length of the time series: the likelihood in (7) has to be maximized over a  $J$ -dimensional space of model weights, using only  $N$  years of data.

To circumvent the problem of the increasingly high dimensionality of the optimization space with increasing  $J$ , we now introduce a two-stage optimization procedure, wherein the model combination is always limited to a *single* model plus climatology, as given by Eq. (6).

In stage 1, each model is combined with the climatological forecast individually by performing  $J$  separate optimizations using (6) and (7). This yields a set of (nonnegative) model weights  $w_j^{(1)}$  ( $j = 1 \dots J$ ) that express each model's performance compared to an  $n$ -yr climatology.

In stage 2, we combine the forecast probabilities of the  $J$  models together according to the normalized values  $w_j^{(1)}$  (i.e., the weights from stage 1) to form a new set of GCM forecast probabilities:

$$P_{kt}^{(2)}(y) = \frac{1}{\sum_j w_j^{(1)}} \sum_{j=1}^J w_j^{(1)} \frac{m_{jkt}}{m_j}. \quad (9)$$

Equation (6) is then solved for  $w^{(2)}$ , by substituting  $P_{kt}^{(2)}(y)$  and  $m^{(2)} = \sum_{j=1}^J m_j$ , and then using Eq. (7) as before. The final weights of the individual models are then disaggregated according to their values in stage 1:

$$w'_j = \frac{1}{\sum_j w_j^{(1)}} w_j^{(1)} w^{(2)}. \quad (10)$$

The weight maps (not shown) produced using Eqs. (9) and (10) are much more evenly weighted between models than those from Eq. (8), but they continue to exhibit noise at the grid-box scale. The short length of the training dataset used to derive the weights (48 yr) suggests that sampling variability is still potentially a problem. To help alleviate this, a repeated subsampling procedure akin to the bootstrap was performed, by repeating the entire two-stage procedure multiple times. Each time, a contiguous block of 6 yr was withheld from the dataset, and the optimal weights computed. The resulting 43



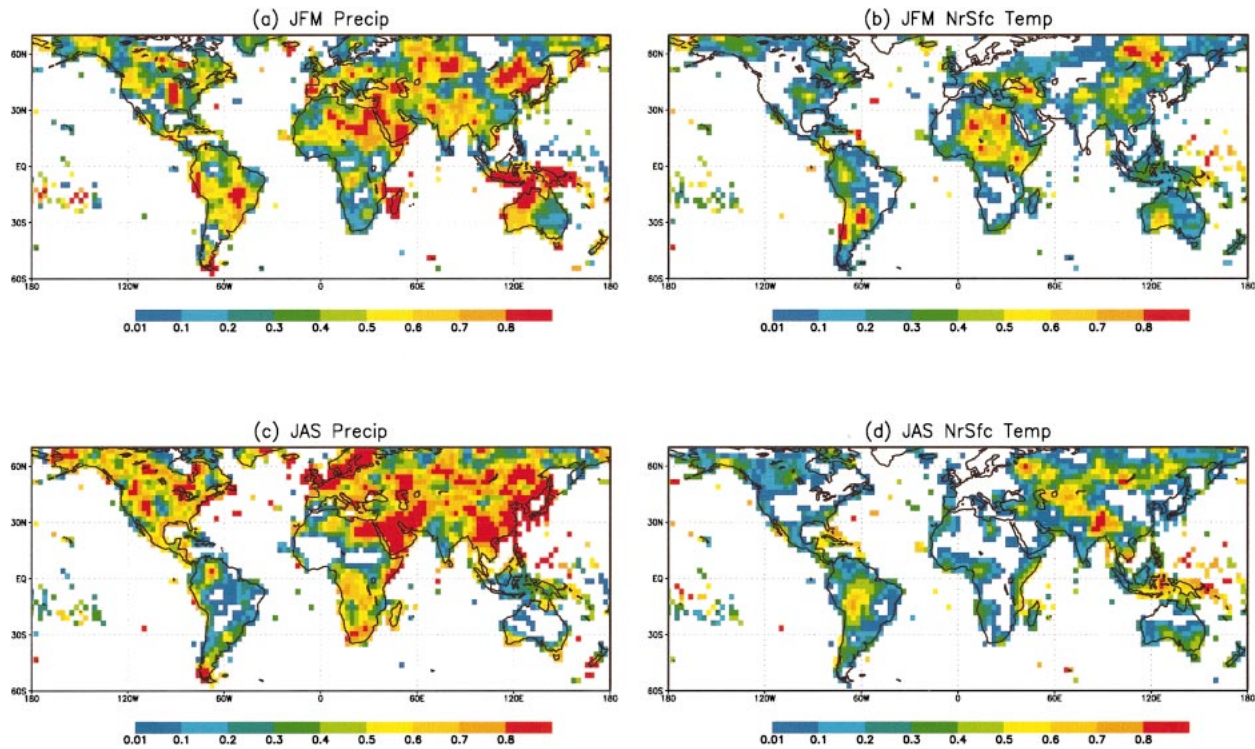


FIG. 1. The weight values assigned to the climatological forecast by the revised six-model optimal combination scheme: (a) JFM precipitation, (b) JFM temperature, (c) JAS precipitation, and (d) JAS temperature. Weights  $< 0.01$  are denoted by white.

estimates of the optimal weights were then simply averaged together. Averaging over multiple subsamples is designed to reduce the effects of sampling variations on the optimization: it has the effect of largely removing the white areas on the weight maps where weights are zero, replacing them with small values ( $< 0.1$ ) (not shown). The observed and model tercile values were kept fixed at their 1968–97 values. Little sensitivity was found to recomputing them from the 42-yr subsample each time. Six-year blocks were chosen here as an ad hoc way to take into account serial correlation, considering that the “low frequency” component of the El Niño–Southern Oscillation (ENSO) has a period of 4–5 yr.

Up until this point, we have computed the weights independently at each of the 2829 land grid boxes of the models. While the resolution of a GCM is nominally at the grid-box scale, it is not expected to be as skillful at this scale as at more aggregated scales (Gong et al. 2003), and much of the variability in model weights between adjacent grid boxes must be regarded as sampling variability. To this end, we introduce a nine-point binomial spatial smoother into the two-stage algorithm. In this case we maximize the likelihood:

$$L(w) = \prod_{i=1}^9 \prod_{t=1}^N E(Q_{itk*}), \quad (11)$$

where the  $i$  subscript sums over adjacent grid points.

The central point is counted twice (to give a binomial smoother), and grid boxes that fall over ocean areas are excluded (for which there is no observational verification data).

The weights assigned to the climatological equal-odds forecast are plotted in Fig. 1, computed using the revised two-stage scheme, including both spatial averaging of the likelihood function and averaging the weights across data subsamples. Here we plot the normalized climatological weights  $w_{\text{Clim}} = n/[n + m^{(2)}w^{(2)}]$ , so that the climatological and model weights sum to 1 at each grid box. The optimal climatological weights are smaller for temperature than for precipitation, consistent with the expected higher skill of near-surface temperature (given prescribed observed SST) compared to precipitation, which is a complex derived variable in GCMs. The  $w_{\text{Clim}}$  exhibits considerable spatial and seasonal variation. The red shading ( $w_{\text{Clim}} > 0.8$ ) occurs at many locations in the precipitation-weight maps, denoting areas where the multimodel ensemble hindcasts lack skill; forecasts issued for these regions will largely resort to the “prior” climatological forecast probabilities, with near-zero RPSS. It is seen later in Tables 2 and 3 that this largely removes negative values in spatial averages of RPSS.

The optimal model weights for JAS precipitation are shown in Fig. 2, with the normalization  $w_{\text{Model}} = w^{(2)}m^{(2)}/[n + m^{(2)}w^{(2)}]$ . The weights of the individual models often tend to be in the range 0.1–0.3, for the six-model

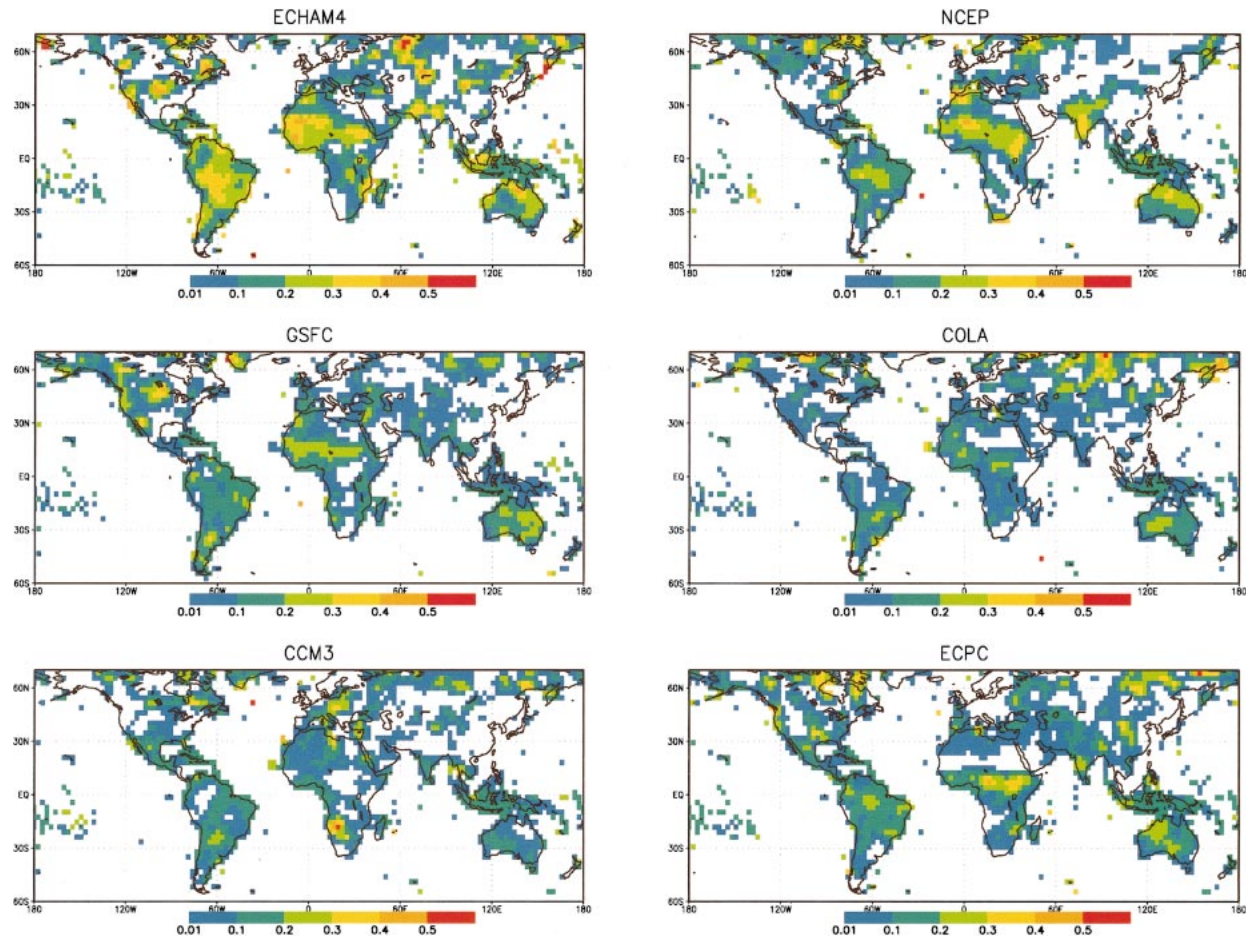


FIG. 2. The weight values assigned to each model simulation by the revised six-model optimal combination scheme, for JAS precipitation. Weights  $< 0.01$  are denoted by white.

combination. In many areas, the revised scheme tends to weight the models fairly evenly. Nonetheless, closer inspection reveals important intermodel differences that may be informative to model developers. Figure 3 shows the optimal model weights for JAS temperature, which are, not surprisingly, generally somewhat larger than for precipitation. The 24-member ECHAM4.5 model receives higher weights in both precipitation and temperature than the other five models, which have only 9–10 members (Table 1). The effect of ECHAM4.5 ensemble size is investigated in section 4. In general, the weight maps in Figs. 1–3 are less noisy than those of RLZ and better reflect the spatial scale on which GCMs are expected to be more skillful.

#### 4. Combined model skill

##### a. Time-average RPSS maps

Figure 4 shows maps of time-average RPSS values for precipitation and temperature during JAS, for both the simple pooled multimodel ensemble and the optimal combination, together with the difference between them.

In all cases the RPSS is cross-validated as described in section 2, so that the weights and RPSS are not computed from the same data. The models' skill varies considerably by geographical location and by variable. Indeed, the JAS precipitation skill is highly regional and is largely confined to most of South America, equatorial Africa, South Asia, and Australasia; this skill originates from the sensitivity of the tropical atmosphere to SST anomalies and to ENSO in particular (Ropelewski and Halpert 1987; Barnston and Smith 1996).

The precipitation skill of the simple pooled ensemble is largely negative in the extratropics. The optimal multimodel combination replaces a large fraction of the negative-precipitation RPSS values with near-zero values. From Fig. 1, this can be seen to be due to the high weighting given to the climatological forecast in many of these areas, clearly demonstrating the impact of including the climatology in the multimodel combination. The impact is smaller in the more-skillful temperature hindcasts, although the negative RPSS values over Amazonia and Indonesia are much reduced in the optimal combination. The skill scores of temperature are



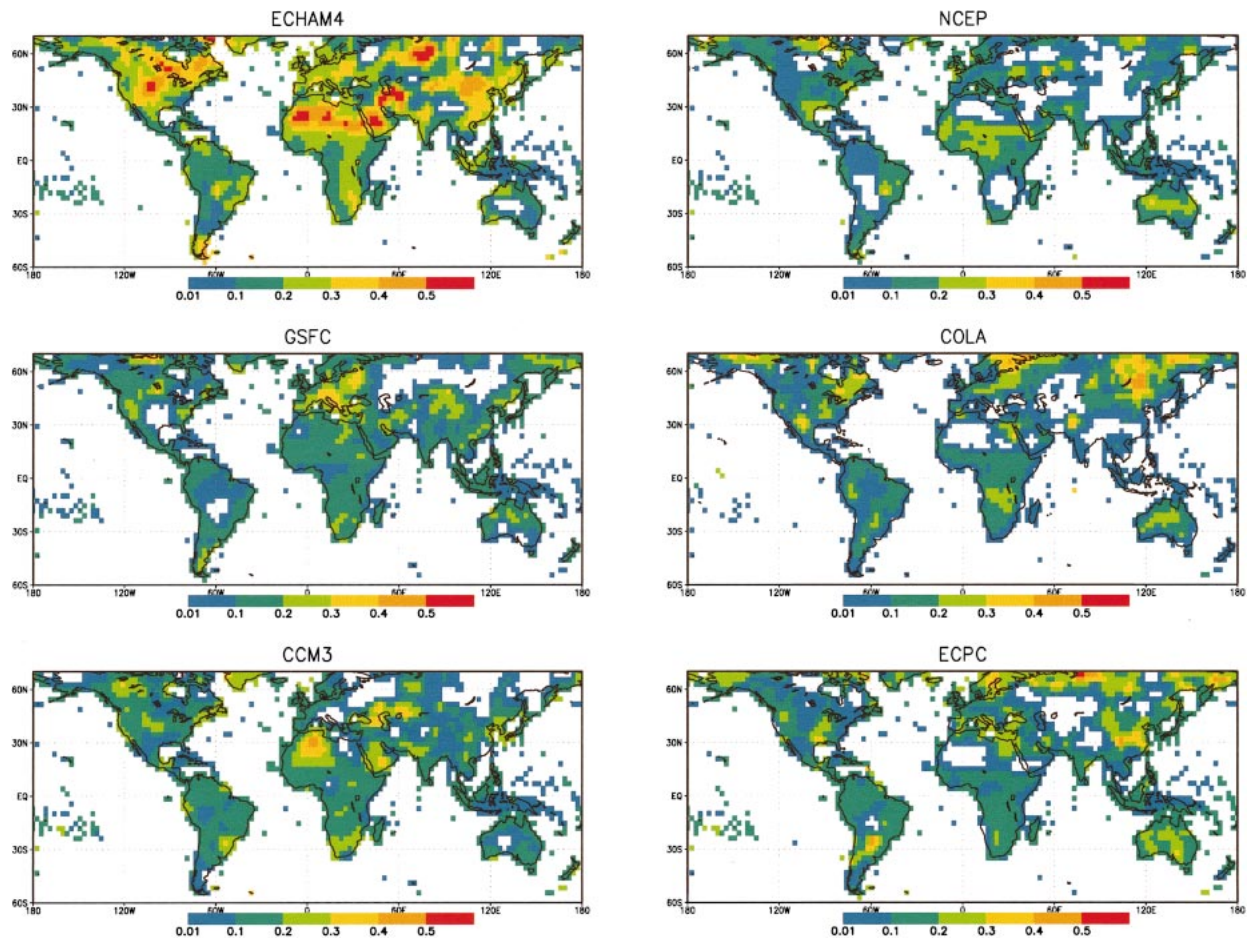


FIG. 3. The optimal model weights from the revised six-model combination for JAS near-surface temperature. Weights  $< 0.01$  are denoted by white.

low over central Asia may be associated with remoteness from the prescribed SSTs. The difference maps (Figs. 4e and 4f) demonstrate that the optimal combination is generally considerably more skillful than the simple pooling for both precipitation and temperature. There is often improvement in skill in areas that already have skill in the simple pool, so that the improvement is not just limited to replacing low-skill areas with climatology. There are, nonetheless, also a few regions of decreasing skill, particularly in temperature over North America and Siberia.

#### b. RPSS of individual models

The RPSS of the individual models and various multimodel ensembles are summarized in Tables 2 and 3, in terms of spatiotemporal averages over the land areas of the Tropics and extratropics (divided at  $30^\circ$  latitude). For precipitation (Table 2), all the *individual* models have negative skills (i.e., worse than climatology) in both domains. This is largely the case for temperature as well, except when all 24 members are included in

the ECHAM4 ensemble (Table 3). Increasing the ensemble size here has a clear benefit on both temperature and precipitation RPSS averages. The pooled ensembles perform much better than the individual models, but the average RPSS values are still near-zero for precipitation over these large domains.

#### c. Combinations of three versus six models

The sensitivity of the RPSS to the number of models included in the ensemble is shown in Tables 2 and 3 and Figs. 5 and 6, for both the pooled ensemble and optimal model combination. Here we compare the full six models, against all possible (i.e., 20) subsets of three models. To construct Figs. 5 and 6, we first identified the best, middle, and worst three-model subsets by ranking the 20 *time-averaged* RPSS values. We then plot the time series of these three particular three-model subsets, together with the full six-model combination. In Tables 2 and 3, we simply give the range of RPSS over all 20 possible subsets.

Combining six models instead of three almost always

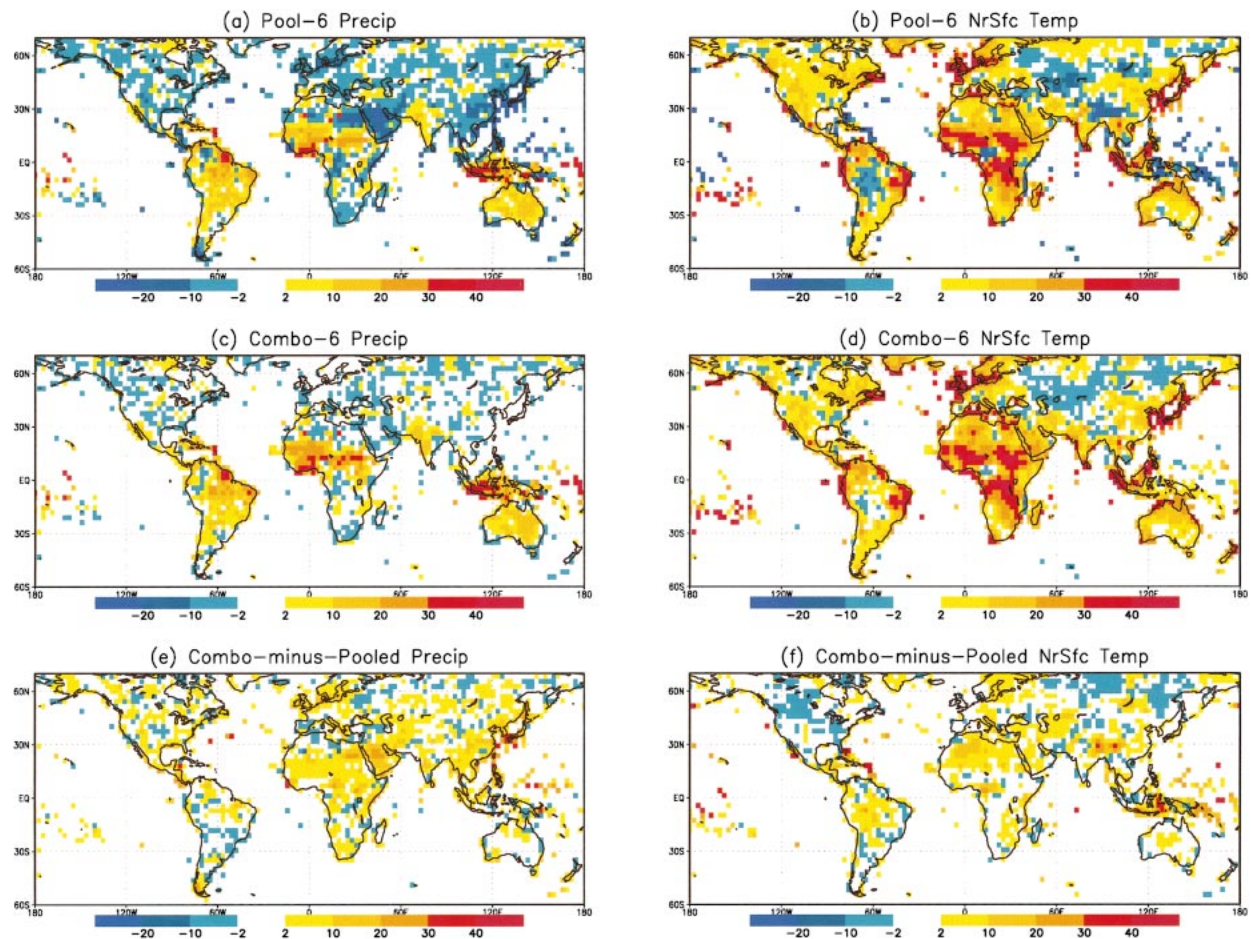


FIG. 4. The RPSS of six-model ensembles of (left) JAS precipitation and (right) JAS near-surface temperature. (a), (b) Simple pooled ensembles; (c), (d) optimal combinations; and (e), (f) differences between pooled and optimal combinations. Blue denotes negative RPSS values, near-zero values (i.e., the climatological forecast) are white, and positive RPSS values are denoted by yellow and red. The pooled ensemble comprises all six models, with 83 members in total. All values computed for 1953–95. Absolute RPSS differences  $<2\%$  are denoted by white in (e) and (f).

TABLE 2. Spatially averaged RPSS values for precipitation, over the Tropics ( $30^{\circ}\text{S}$ – $30^{\circ}\text{N}$ ) and extratropics (poleward of  $30^{\circ}$ ), for the individual models and various multimodel ensembles. Pool: pooled ensemble; Cmbo: revised two-stage Bayesian combination with spatial smoothing of objective function; RLZ: Bayesian combination of Rajagopalan et al. (2002) (with cross-validation). The  $-n$  suffix denotes the number of models in the ensemble. The three-model combination is given as the range of all 20 possible such combinations. The ECHAM4+, Cmbo-6+, and RLZ-6+ entries use the extended 24-member ensemble; all other entries use a 10-member ECHAM4 ensemble. All results are for the 1953–95 period.

	JFM		JAS	
	Tropics	Extratropics	Tropics	Extratropics
ECHAM4	−12.35	−9.41	−7.71	−11.33
ECHAM4+	−8.04	−4.04	−3.22	−5.85
NCEP	−18.72	−11.27	−14.30	−12.81
NSIPP1	−20.40	−11.75	−19.76	−13.83
COLA	−22.83	−13.68	−23.59	−13.67
CCM3	−13.54	−8.69	−15.54	−12.81
ECPC	−14.07	−14.79	−15.10	−13.76
Pool-3	−5.82 to −1.87	−4.16 to −1.90	−3.20 to 0.68	−5.06 to −4.33
Pool-6	−0.35	−0.24	3.08	−2.21
Cmbo-3	2.33 to 3.10	−0.06 to 0.24	4.55 to 6.01	−0.39 to −0.15
Cmbo-6	3.30	0.19	6.85	−0.53
Cmbo-6+	3.39	0.44	6.97	−0.55
RLZ-6+	1.42	−1.12	5.26	−2.37



TABLE 3. Spatially averaged RPSS values for near-surface temperature. See Table 2 for details.

	Jan–Mar (JFM)		Jul–Sep (JAS)	
	Tropics	Extratropics	Tropics	Extratropics
ECHAM4	−1.03	−4.82	−1.35	−4.36
ECHAM4+	2.76	−0.14	2.24	1.01
NCEP	−7.99	−14.40	−16.35	−11.05
NSIPP1	−7.28	−11.36	−15.23	−11.66
COLA	−19.62	−17.90	−19.84	−12.59
CCM3	1.18	−4.21	−4.60	−5.15
ECPC	−2.00	−6.63	−10.41	−7.81
Pool-3	8.40 to 11.22	−0.68 to 4.17	4.71 to 9.62	2.0 to 5.56
Pool-6	13.41	4.95	11.79	7.38
Cmbo-3	12.50 to 14.71	2.69 to 4.49	12.44 to 14.31	4.11 to 7.51
Cmbo-6	15.68	5.07	15.78	7.99
Cmbo-6+	15.75	5.27	16.01	8.16
RLZ-6+	14.79	4.01	15.48	7.35

leads to increases in skill. The payoff is larger for the simple pool than for the optimal combination. If we know a priori which three models to pick, the increase in skill of adding the remaining three models is often quite modest.

The RPSS of the six-model optimal combination with the extended 24-member ECHAM4 ensemble is denoted in Tables 2 and 3 as Cmbo-6+. Even in the optimal six-model combination, including an additional 14 ECHAM4 members does yield increases in overall skill.

#### d. Interannual skill variations

Time series of spatially averaged RPSS values are plotted in Figs. 5 and 6. In general, interannual variations in skill are larger in the Tropics than the extratropics for both precipitation and temperature. This reflects the fact that interannual anomalies in tropical SST such as El Niño produce large-scale responses in the Tropics, but much less so in the extratropics. The peaks in skill are largely consistent with the timing of ENSO events (Goddard 2005). Note that the spatial averages are not weighted by area and are thus biased toward higher latitudes, because of the convergences of the meridians.

Figure 5 clearly illustrates how the optimal weighting boosts the tropical skill of precipitation forecasts in years in which it is relatively low, reducing the amount of interannual and interdecadal skill variability. In the extratropics, substantial spurious interannual variation in the pooled-model skill are largely eliminated in the optimal combination, to yield near-zero RPSS in all years. Similar comments apply to temperature (Fig. 6), although interannual variations in skill are larger. In the extratropics there appears to be a trend toward increasing temperature skill; this may be an artifact associated with recent upward trends in temperature, together with the use of a fixed climatological normal.

#### e. Comparison with the RLZ scheme

The skill of the revised multimodel optimal combination is compared to the original RLZ scheme in Tables

2 and 3. The revised Bayesian scheme of Eqs. (9)–(11) is found to be more skillful on average than the RLZ scheme [Eqs. (7)–(8)], especially for precipitation. This is also clear in RPSS maps similar to Fig. 4 (not shown). Tables 2 and 3 also indicate that for six models, the RLZ scheme is actually *less* skillful than the pooled ensemble in the extratropics for both precipitation and temperature. All our computations with the RLZ scheme were performed with the weights averaged across data subsamples, as described in section 3, so that sampling variability should be reduced compared to the results reported by RLZ.

### 5. Discussion and conclusions

An improved Bayesian weighting scheme is developed and used to combine several atmospheric GCM ensembles forced with observed SSTs. We combine the GCM simulations of precipitation or near-surface temperature at each land grid point, based on the prior belief that the GCM-simulated tercile-category probabilities are equal to climatological probabilities of  $\frac{1}{3}$ . The scheme's skill is compared against the individual model ensembles (with 9–24 members), simple pooled ensembles of three and six models, as well as the original version of the Bayesian weighting scheme devised by Rajagopalan et al. (2002). The ranked probability skill score (RPSS) is used as the skill measure, cross-validated by withholding six contiguous years at a time from the 48-yr 1950–97 time series of model simulations and observed precipitation and temperature.

Our results demonstrate clear gains in skill by simply pooling together the ensemble hindcasts made with individual GCM ensembles, corroborating previous studies (Fraedrich and Smith 1989; Graham et al. 2000; Palmer et al. 2000; Pavan and Doblas-Reyes 2000; Peng et al. 2002). A pooling of six models is found to be almost always superior to a pooling of just three models, although the gain is modest if the three best models (measured over the 46-yr period used to compute RPSS) can be identified a priori. As expected, the precipitation

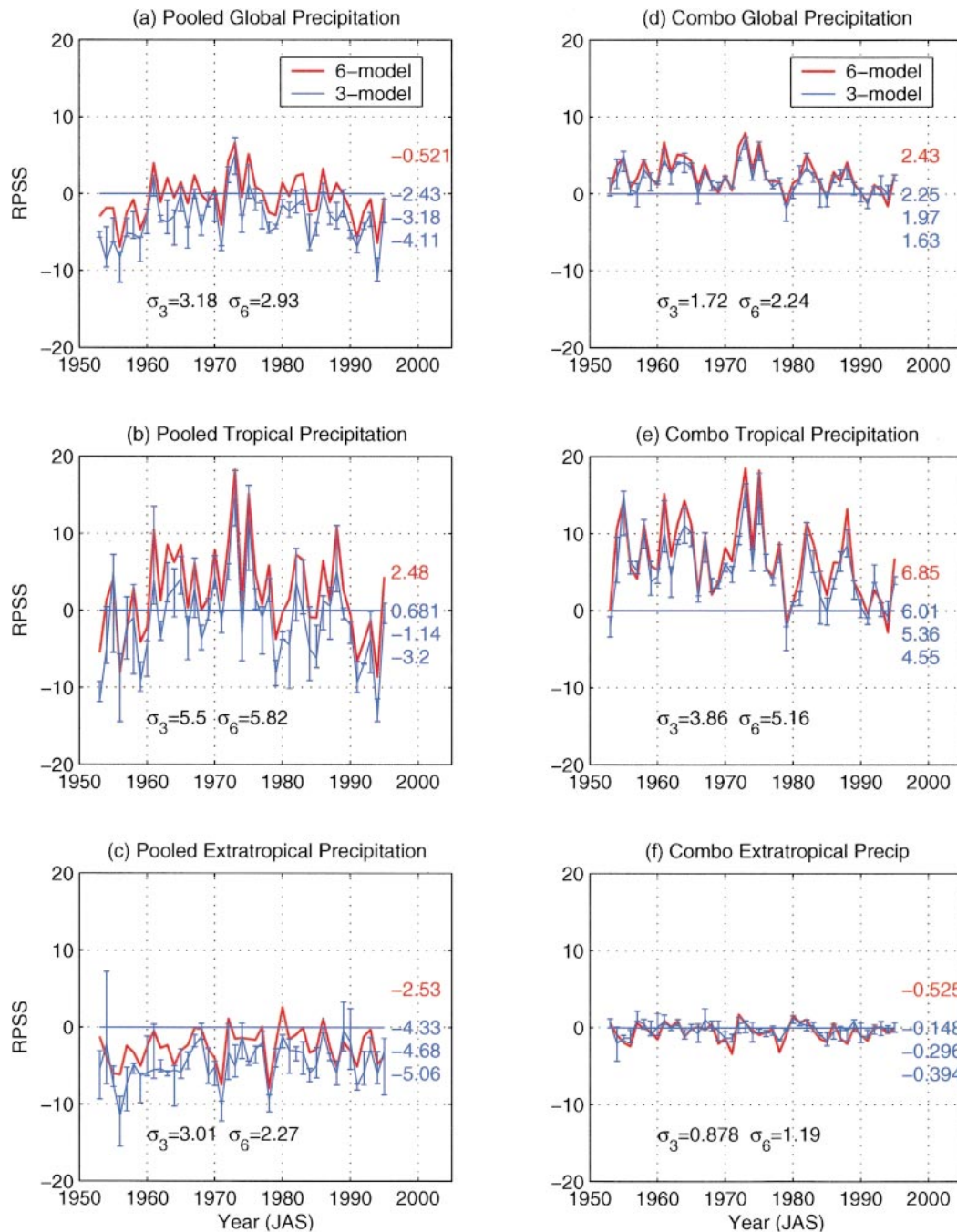


FIG. 5. Spatially averaged RPSS values for JAS precipitation as a function of year. (top) Global average, (middle) Tropics ( $30^{\circ}\text{S}$ – $30^{\circ}\text{N}$ ), and (bottom) extratropics. (left) The results from the pooled ensembles, and (right) the optimally combined ensembles. Each panel shows the six-model combination (red) and the three-model combination with median overall skill (blue). The error bars on the latter show the range of three-model RPSS values given in Table 2. The zero line is indicated in blue. The numbers on the right give the respective time averages. The interannual standard deviations are also indicated, with the three-model value taken from the median-performing pick of three models.

skill is higher within the Tropics than in the extratropics, and the temperature skill is higher than for precipitation.

The revised Bayesian optimal weighting scheme is shown to outperform the pooled ensemble. In the extratropics, its main impact is to bring much of the large

area of negative precipitation RPSS up to near-zero values. Effectively, it progressively replaces the model GCM-simulated probabilities with climatological equal-odds values in these areas by downweighting the model simulations relative to the climatological forecast. There

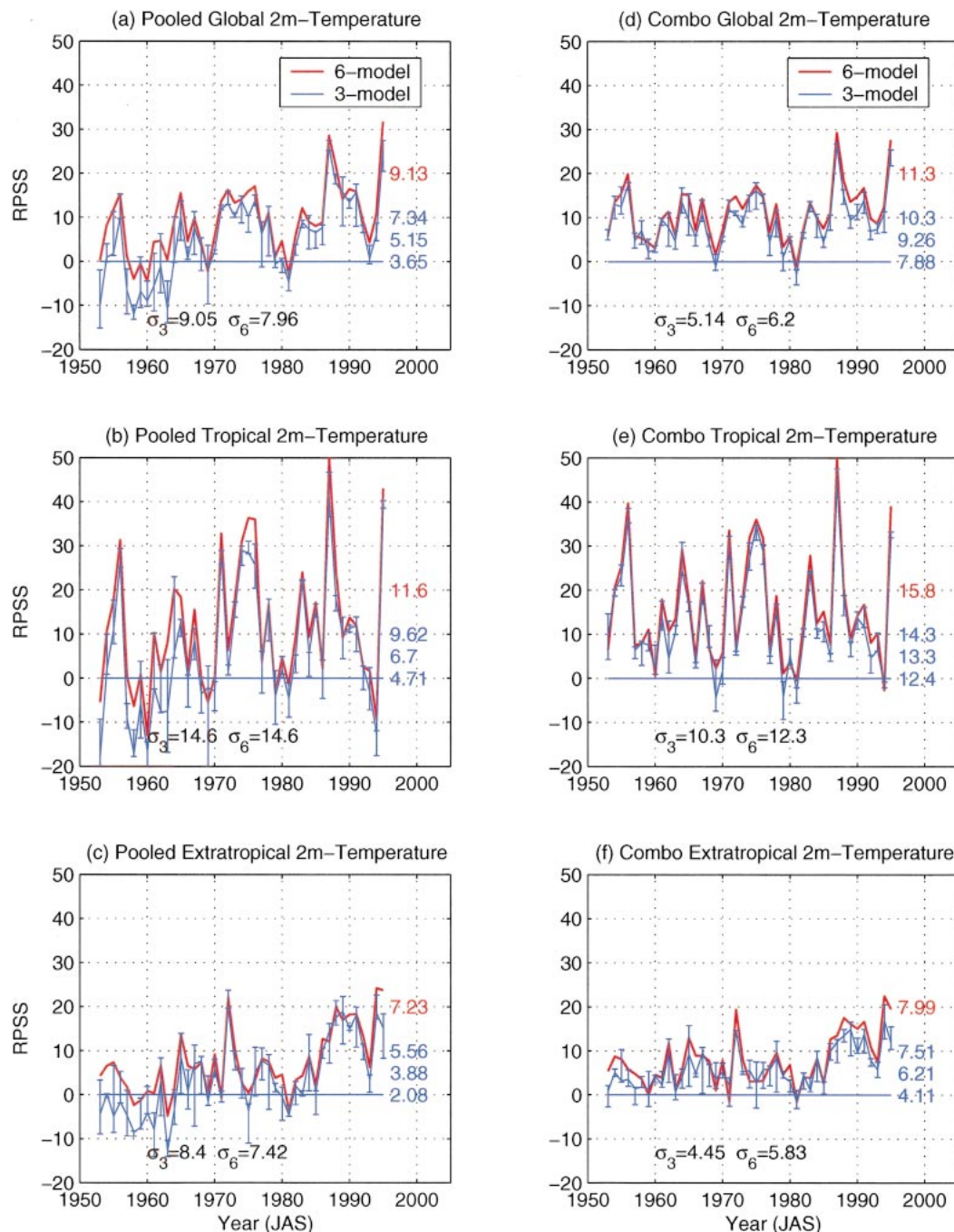


FIG. 6. Same as Fig. 5, but spatially averaged RPSS values for JAS near-surface temperature as a function of year.

are also substantial gains in the average tropical precipitation skill, and even in individual regions of positive skill in the pooled ensemble. Increases in skill are more modest for the temperature simulations, which are more skillful to begin with. However, there are nonetheless regions of negative RPSS for temperature in the pooled ensemble that are much reduced in the optimal combination. Interannual variations in skill are reduced,

especially in extratropical precipitation, where they are largely spurious.

Improvements made to the original Bayesian scheme in the form of reducing the dimensionality of the numerical optimization and including spatial smoothing of the likelihood function are shown to substantially increase the RPSS cross-validated skills. Maps of the model weights are less noisy than in the original scheme,



and the weights are distributed more evenly among the models.

The number of parameters to be estimated increases with the number of models, while the amount of training data remains the same. This translates into higher sampling variability in the optimal weights selected and hence a degradation in the performance of the “best” model combination selected. In the revised scheme, each model is first calibrated against climatology independently, and this potentially leads to a more robust weighting and smoothing of that model’s results toward climatology. The multimodel combination can be interpreted as a way to reduce sampling variance, together with the conditional biases of the individual models. Sampling variance is inherent in estimating the GCM probabilities in Eq. (1), and in determining the terciles of the 30-yr observational climatology. Increasing each model’s ensemble size will decrease the variance of the estimate, and this is seen to increase the skill of the ECHAM model (Tables 2 and 3). Model conditional biases are more difficult to alleviate, although model output statistics (MOS) corrections can often remove spatial conditional biases. Averaging over six models in Eq. (9) should lead to a decrease in the conditional bias of the combined forecast, provided the individual models are suitably weighted. There is evidence of this in regions where both (a) the weighted ensemble skill beats the pooled ensemble, and (b) the latter is itself skillful, so that the increase in skill is not just coming from including the climatology forecast. The JAS precipitation simulations over West Africa and India are examples. However, it is less clear in general whether the benefit of increasing the number of models is likely to be greater than the mere increase in the number of ensemble members of selected models (Pavan and Doblas-Reyes 2000).

Multiple colinearity is often a concern when combining together several predictors using multiple linear regression. In the context of the multimodel ensemble, consider the case of two identical GCMs run with the same number of ensemble members, but from a different set of initial conditions. In the one-stage scheme of RLZ [Eq. (8)], there will be nonuniqueness in the weights assigned to the two models, since *any* combination of them will yield a similar log-likelihood score. However, a forecast made with the multimodel ensemble will not be impacted, and this only presents a problem if we wish to use the “optimal” weights to attribute skill to either model. In contrast, the revised two-stage scheme does not suffer from this nonuniqueness in the weights. Each model is calibrated independently against climatology, so two near-identical GCMs will receive similar weight, whose magnitude depends upon skill against climatology; any colinearity of errors will not be reflected in the weights. Nonetheless, it is worth pointing out that even the revised scheme will not be able to distinguish between models with similar skill, but which achieve that skill through different mechanisms.

Maps of the optimal model weights, such as Figs. 2 and 3, provide a useful byproduct of the optimal weighting exercise. The maps provide an additional metric of model skill and intercomparison that may be of value to GCM developers.

One weakness of the Bayesian scheme that persists despite the improvements to the algorithm is an occasional tendency toward high GCM precipitation weights in some high-latitude regions (see Fig. 2). We would not expect the GCMs’ precipitation simulations to be skillful in many of these regions. If the GCM probabilities—or the combined second-stage model probabilities in Eq. (9)—beat the climatological ones over the training period, even by a slight amount, then the optimal combination can heavily favor the model. In effect, the likelihood optimization is not sensitive to distance. Averaging the weights across data subsamples and spatial averaging of the likelihood function both alleviate the problem to some extent because they reduce sampling variability, which is the root of the spurious model skill in question. No account was taken of the convergence of the meridians toward the poles. In future work, the spatial smoothing could be performed over a fixed area, rather than a fixed number of grid points. However, sampling variability can have large spatial scale and will never be completely eliminated given the relatively short records available.

The revised scheme appears to be well suited to combining larger sets of models, and, in the future, it should be possible to include statistical models into the weighted ensemble without fundamental difficulty. The skill of the optimal combination is always increased (at least in the large spatial averages considered) when the number of models in the combination is increased from three to six, regardless of which models are included in the three-model combination. With the exception of the 24-member ECHAM4 ensemble, the number of ensemble members for each model was limited to about 10. Increasing the size of the ECHAM4 model ensemble from 10 to 24 members increases this individual model’s RPSS substantially and even has a positive impact on the six-model combination. Thus, there is a potential payoff to be achieved by increasing the size of the model ensembles.

Finally, it should be remembered that the RPSS values reported in this paper apply to the case of prescribed monthly mean SST. These skills decrease substantially in retrospective forecasts in which predicted SST is used to force the atmospheric GCMs (Goddard and Mason 2002). On the other hand, some increased skill can be expected from initializing the models with observed estimates of soil moisture, snow cover, and atmospheric initial data. In any case, optimally weighted multimodel ensembles form a valuable component of a seasonal climate forecasting system.

*Acknowledgments.* We are grateful to Tony Barnston, Simon Mason, and Balaji Rajagopalan for helpful dis-

cussions, and Tony Barnston for his valuable comments on an earlier version of the manuscript. We especially wish to thank the six GCM modeling centers whose model runs formed the basis for our study. The comments of two anonymous reviewers improved the manuscript. This work was supported by the International Research Institute for Climate Prediction and a National Oceanic and Atmospheric Administration Grant.

## REFERENCES

- Barnston, A. G., and T. M. Smith, 1996: Specification and prediction of global surface temperature and precipitation from global SST using CCA. *J. Climate*, **9**, 2660–2697.
- , S. J. Mason, L. Goddard, D. G. DeWitt, and S. E. Zebiak, 2003: Multi-model ensembling in seasonal climate forecasting at IRI. *Bull. Amer. Meteor. Soc.*, **84**, 1783–1796.
- Doblas-Reyes, F. J., M. Deque, and J.-P. Piedelievre, 2000: Multi-model spread and probabilistic seasonal forecasts in PROVOST. *Quart. J. Roy. Meteor. Soc.*, **126**, 2069–2088.
- Epstein, E. S., 1969: A scoring system for probability forecasts of ranked categories. *J. Appl. Meteor.*, **8**, 985–987.
- Fraedrich, K., and N. R. Smith, 1989: Combining predictive schemes in long range forecasting. *J. Climate*, **2**, 291–294.
- Gelman, A., J. B. Carlin, H. S. Stern, and D. B. Rubin, 1995: *Bayesian Data Analysis*. Chapman and Hall, 526 pp.
- Goddard, L., 2005: El Niño: Catastrophe or opportunity? *J. Climate*, in press.
- , and S. J. Mason, 2002: Sensitivity of seasonal climate forecasts to persisted SST anomalies. *Climate Dyn.*, **19**, 619–631.
- , A. G. Barnston, and S. J. Mason, 2003: Evaluation of the IRI's "net assessment" seasonal climate forecasts: 1997–2001. *Bull. Amer. Meteor. Soc.*, **84**, 1761–1781.
- Gong, X., A. G. Barnston, and M. N. Ward, 2003: The effect of spatial aggregation on the skill of seasonal precipitation forecasts. *J. Climate*, **16**, 3059–3071.
- Graham, R. J., A. D. L. Evans, K. R. Mylne, M. S. J. Harrison, and K. B. Robertson, 2000: An assessment of seasonal predictability using atmospheric general circulation models. *Quart. J. Roy. Meteor. Soc.*, **126**, 2211–2240.
- Hack, J. J., J. T. Kiehl, and J. W. Hurrell, 1998: The hydrological and thermodynamic characteristics of the NCAR CCM3. *J. Climate*, **11**, 1179–1206.
- Hagedorn, R., 2001: Development of a multi-model ensemble system for seasonal to interannual prediction. *Proc. XXVI General Assembly of the EGS*, Nice, France, European Geophysical Society.
- Kanamitsu, M., and K. C. Mo, 2003: Dynamical effect of land surface processes on summer precipitation over the southwestern United States. *J. Climate*, **16**, 496–509.
- , and Coauthors, 2002: NCEP dynamical seasonal forecast system 2000. *Bull. Amer. Meteor. Soc.*, **83**, 1019–1037.
- Kumar, A., M. P. Hoerling, M. Ji, A. Leetmaa, and P. Sardeshmukh, 1996: Assessing a GCM's suitability for making seasonal predictions. *J. Climate*, **9**, 115–129.
- , A. G. Barnston, and M. P. Hoerling, 2001: Seasonal predictions, probabilistic verifications, and ensemble size. *J. Climate*, **14**, 1671–1676.
- Lorenz, E. N., 1963: Deterministic nonperiodic flow. *J. Atmos. Sci.*, **20**, 130–148.
- Mason, S., 2004: On using "climatology" as a reference strategy in the Brier and ranked probability skill scores. *Mon. Wea. Rev.*, **132**, 1891–1895.
- New, M., M. Hulme, and P. D. Jones, 1999: Representing twentieth-century space–time climate variability. Part I: Development of a 1961–90 mean monthly terrestrial climatology. *J. Climate*, **12**, 829–856.
- , —, and —, 2000: Representing twentieth-century space–time climate variability. Part II: Development of a 1961–90 monthly grid of terrestrial surface climate. *J. Climate*, **13**, 2217–2238.
- Palmer, T. N., C. Brankovic, and D. S. Richardson, 2000: A probability and decision-model analysis of PROVOST seasonal multi-model ensemble integrations. *Quart. J. Roy. Meteor. Soc.*, **126**, 2013–2034.
- Pavan, V., and F. J. Doblas-Reyes, 2000: Multi-model seasonal hindcasts over the Euro-Atlantic: Skill scores and dynamic features. *Climate Dyn.*, **16**, 611–625.
- Peng, P., A. Kumar, H. Van den Dool, and A. G. Barnston, 2002: An analysis of multimodel ensemble predictions for seasonal climate anomalies. *J. Geophys. Res.*, **107**, 4710, doi:10.1029/2002JD002712.
- Rajagopalan, B., U. Lall, and S. E. Zebiak, 2002: Categorical climate forecasts through regularization and optimal combination of multiple GCM ensembles. *Mon. Wea. Rev.*, **130**, 1792–1811.
- Reynolds, R. W., 1988: A real-time global sea surface temperature analysis. *J. Climate*, **1**, 75–87.
- , and T. M. Smith, 1994: Improved global sea surface temperature analyses using optimum interpolation. *J. Climate*, **7**, 929–948.
- Roeckner, E., and Coauthors, 1996: The atmospheric general circulation model ECHAM4: Model description and simulation of present-day climate. Max-Planck-Institut für Meteorologie Rep. 218, 90 pp.
- Ropelewski, C. F., and M. S. Halpert, 1987: Global and regional scale precipitation patterns associated with the El Niño/Southern Oscillation. *Mon. Wea. Rev.*, **115**, 1606–1626.
- Wilks, D. S., 1995: *Statistical Methods in the Atmospheric Sciences*. International Geophysical Series, Vol. 59, Academic Press, 464 pp.