

Atmospheric Seasonal Predictability and Estimates of Ensemble Size

ČEDO BRANKOVIĆ AND T. N. PALMER

ECMWF, Shinfield Park, Reading, United Kingdom

(Manuscript received 30 June 1995, in final form 26 January 1996)

ABSTRACT

Results from a set of nine-member ensemble seasonal integrations with a T63L19 version of the European Centre for Medium-Range Weather Forecasts (ECMWF) model are presented. The integrations are made using observed specified sea surface temperature (SST) from the 5-year period 1986–90, which included both warm and cold El Niño–Southern Oscillation (ENSO) events. The distributions of ensemble skill scores and internal ensemble consistency are studied. For years in which ENSO was strong, the model generally exhibits a relative high skill and high consistency in the Tropics. In the northern extratropics, the highest skill and consistency are found for the northern Pacific–North American region in winter, whereas for the northern Atlantic–European region the spring season appears to be both skillful and consistent. For years in which ENSO was weak, the distributions of ensemble skill and consistency are relatively broad and no clear distinction between Tropics and extratropics can be made.

Applying a t test to interannual fluctuations over various tropical and extratropical regions, estimates of a minimum useful ensemble size are made. Explicit calculations are done with ensemble size varying between three and nine members; estimates for larger sizes are made by extrapolating the t values. Based on an analysis of 2-m temperature and precipitation, the use of relatively large (approximately 20 members) ensembles for extratropical predictions is likely to be required; in the Tropics, smaller-sized ensembles may be adequate during years in which ENSO is strong, particularly for regions such as the Sahel.

The role of the SST forcing in a seasonal timescale ensemble is to bias the probability distribution function (PDF) of atmospheric states. Such PDFs can, in addition, be a convenient way of condensing a vast amount of data usually obtained from ensemble predictions. Interannual variability in PDFs of monsoon rainfall and regional geopotential height probabilities is discussed.

1. Introduction

The scientific basis for extended-range atmospheric prediction derives principally from the predictability of the atmosphere's lower boundary conditions, particularly sea surface temperature (SST). However, even if SST could be predicted without error, the associated atmospheric evolution would not be uniquely determined, essentially because of the chaotic nature of atmospheric dynamics. As a result, the SST anomaly should be thought of as having a well-defined impact, not on a specific phase-space trajectory corresponding, say, to the atmosphere's evolution over one season, but on the phase-space geometry of the whole atmospheric climate attractor. This impact can be specified in terms of changes to the atmospheric probability distribution function (PDF) over atmospheric states (Palmer 1993; Kumar and Hoerling 1995).

In practice, estimating the impact of prescribed SST anomalies on such probability distributions can be de-

termined only from ensembles of integrations of a dynamical atmospheric model. Examples of such seasonal ensemble integrations have been discussed recently by Branković et al. (1994), hereafter referred to as BPF; Palmer and Anderson (1994); and Barnett (1995). In all these studies, the issue of what constitutes a reasonable lower bound on ensemble size was raised; this paper in part addresses the same issue. Stern and Miyakoda (1995) also explored the feasibility of seasonal prediction from ensembles of 10-year-long integrations.

In BPF, results from three-member seasonal timescale ensembles were reported. Not surprisingly, it was found that the skill of the ensemble mean fields was higher for the strong ENSO years than for the weak ENSO years. In the northern extratropics, the skill tended to be highest in the spring season. This was consistent with the internal spread of the ensemble, which tended to be smallest in spring. However, the three-member ensembles were inadequate to assess the statistical significance of the extratropical response to the underlying SST anomalies.

The ensemble size for all the experiments in BPF has since been increased to nine members (see section 2). The skill scores of anomaly fields and the distribution of ensemble skill scores and ensemble consistency val-

Corresponding author address: Dr. Čedo Branković, European Centre for Medium-Range Weather Forecasts, Shinfield Park, Reading, RG2 9AX United Kingdom.
E-mail: c.brankovic@ecmwf.int

TABLE 1. Initial dates for ECMWF seasonal simulations.

En-semble member	DJF	MAM	JJA	SON
1	28 October	27 January	27 April	28 July
2	29 October	28 January	28 April	29 July
3	30 October	29 January	29 April	30 July
4	31 October	30 January	30 April	31 July
5	1 November	31 January	1 May	1 August
6	2 November	1 February	2 May	2 August
7	3 November	2 February	3 May	3 August
8	4 November	3 February	4 May	4 August
9	5 November	4 February	5 May	5 August

ues for difference fields are calculated on a regional basis (section 3). Based on statistical tests, we determine the extent to which the increase in ensemble size has increased our confidence in being able to estimate reliably the impact of the imposed SSTs on regional variables of practical interest, specifically precipitation and near-surface temperature (section 4). By extrapolation of the results of the statistical analyses presented, the likely impact of further increases of ensemble size is assessed. An assessment of probability forecasts using the nine-member ensembles is made in section 5. A summary and discussion of the results are given in section 6.

2. Experimental details

As in BPF, all integrations described in the main body of this paper were made with the ECMWF model at the reduced horizontal resolution of T63L19, using the so-

called cycle 36 physics package (Simmons et al. 1988; Miller et al. 1992). The integrations were about 120 days long, depending on season and initial date. They cover all seasons over the 5-yr period, from spring (MAM) 1986 to winter (DJF) 1990/91. The nine initial dates for each calendar season are shown in Table 1. They were chosen around the first day of the month preceding the season of interest. Thus, the range of the BPF initial dates is extended to include four days before and two days after the dates of the original three-member ensembles.

Within a given season, the same SSTs (based on the U.S. National Meteorological Center analyses), were used for each ensemble member and were updated every 5 days throughout the integration. For every calendar season, the interannual variation of the SST anomalies for the period spring 1986 to winter 1990/91 includes both warm and cold El Niño–Southern Oscillation (ENSO) events. As discussed in BPF, we categorize seasons according to whether anomalies in an equatorial Pacific SST index, computed over the tropical Pacific strip (7°N–7°S, 160°E–80°W), were either strong and positive, strong and negative, weak and positive, or weak and negative (see Figs. 1 and 2 of BPF).

In the rest of this paper we shall be discussing differences between pairs of ensemble integrations, each ensemble having been made using SSTs for a specific year. Where results are described as associated with “strong ENSO-index years,” one ensemble was made with SSTs where the Pacific index had strong positive anomalies; the other ensemble was made with SSTs where the index had strong negative anomalies. Similarly, for “weak ENSO-index years,” one ensemble was made using SSTs where the index had weak positive

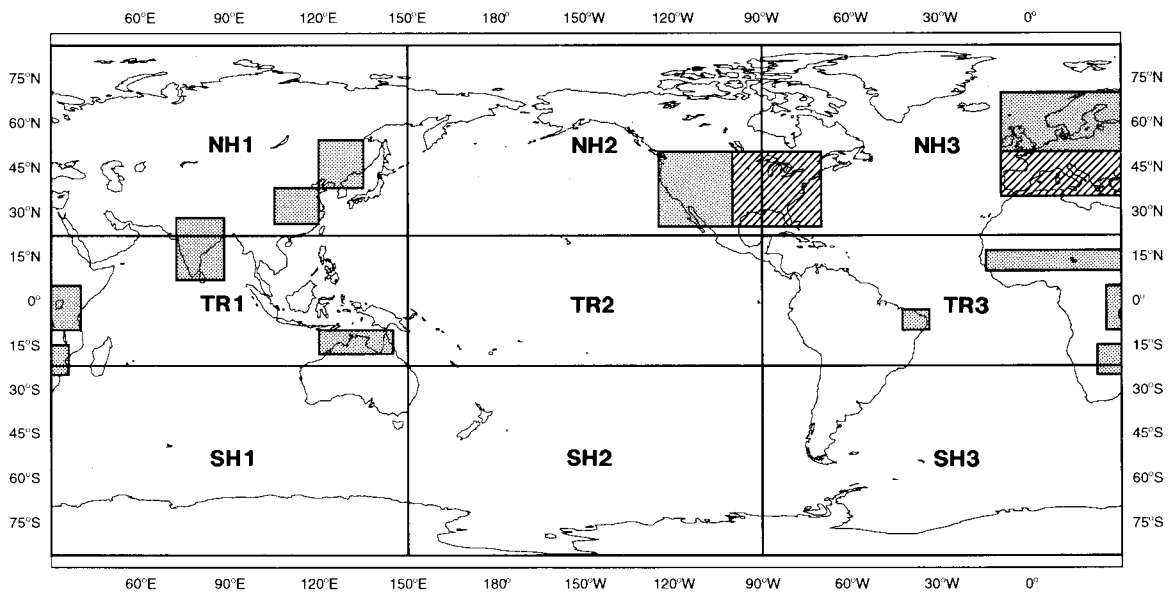


FIG. 1. Nine regions (NH1, NH2, . . .) for which distribution of skill scores and ensemble consistency were computed. Regions for which *t*-test statistic was calculated are shaded. (The eastern Africa and northern Kalahari regions are broken along 30°E.)

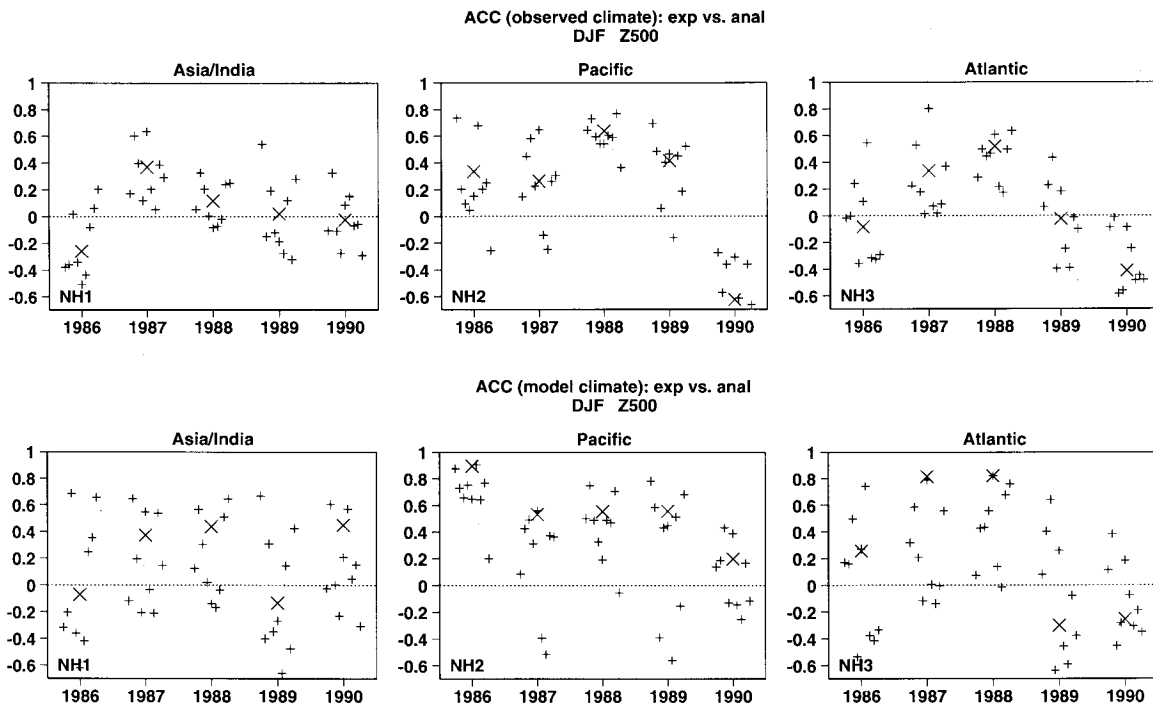


FIG. 2. The DJF 500-mb anomaly correlation coefficients for the three Northern Hemisphere regions (NH1, NH2, NH3): when the “observed climate” is used to calculate model anomalies (top), and when the “model climate” is used (bottom).

anomalies; the other ensemble was associated with an index having weak negative anomalies. In practice, the term “weak ENSO-index years” refers to years in which significant El Niño SST anomalies were absent. Table 2 shows how each season can be characterized with this index.

Throughout the paper we focus our discussion on seasonal averages only. For each experiment, the last three months, corresponding to conventional calendar seasons, were averaged. The seasons are denoted conventionally as: spring—MAM, summer—JJA, autumn—SON, and winter—DJF. For verification purposes, seasonal averages were also computed from ECMWF analysis data.

3. Objective verification of ensembles

a. Skill scores of anomaly fields

In this section, model skill scores are given in terms of anomaly correlation coefficients (ACCs) between seasonally averaged observed anomaly fields and model anomaly fields. For a given season, the observed anomalies have been computed with respect to the 5-yr average (1986–90); we refer to this average as the “observed climate.” For the model, anomalies have been computed with respect to the two different mean fields. First, we used the “observed climate” as above. This methodology is consistent with the ECMWF operational practice for determining the skill of medium-range forecasts. (Of course, in the ECMWF operations the climate

is derived from a much longer period than used here.) In addition to this, model anomalies were also computed with respect to the mean from all integrations in the same 5-yr period, that is, from the total of 45 runs. We refer to this average as the “model climate.” For reason of space, the discussion in this subsection will be restricted to ACCs of the 500-mb heights in the three Northern Hemisphere regions (NH1, NH2, NH3; see Fig. 1) for the DJF season only.

The top row in Fig. 2 shows ACCs when the observed climate is used to calculate model anomalies; the bottom row of Fig. 2 shows ACCs when the model climate is used. The ACCs are shown for all individual integrations within an ensemble (depicted for each year as nine small crosses) and for ensemble averages (larger diagonal crosses). The distribution of small crosses in the vertical is indicative of the intraensemble range of

TABLE 2. Division of the experimental years/seasons according to the index based on the equatorial Pacific SST anomalies.

Amplitude of ENSO	Season	Positive	Negative
Strong	DJF	1986/87	1988/89
	MAM	1987	1989
	JJA	1987	1988
	SON	1987	1988
Weak	DJF	1990/91	1989/90
	MAM	1988	1986
	JJA	1990	1989
	SON	1990	1989

scores, while the small shift between the crosses in the horizontal represents the different initial dates for individual model integrations.

When the observed climate is used, the highest scores are found for the NH2 region (covering the northern Pacific and much of North America) during DJF 1988/89. The range of scores for this winter appears to be the smallest of all winters considered. According to Table 2, the 1988/89 winter was classified as a strong negative ENSO-index season. Skill scores for DJF 1986/87, the strong positive ENSO-index winter, do not differ very much from, for example, skill scores for DJF 1987/88.

On the other hand when the model climate is used (Fig. 2, bottom row), the highest scores for the NH2 region are found during DJF 1986/87. For all but one integration, the ACCs for that season fall between 0.6 and 0.9, and the range of scores is small when compared with the other winters.

The difference between the two sets of skill scores for NH2 during DJF 1986/87 is striking. The differences for DJF 1988/89 and for the weak ENSO-index winters are somewhat smaller, though by no means negligible. Similar differences in skill scores are also seen for the NH3 (Asian) and for NH3 (North Atlantic–European) regions.

Clearly, these inconsistencies are associated with the different reference fields used to calculate model anomalies. Apparent improvement in the model skill in DJF 1986/87, when the model climate is used, may be because of a larger covariance between model and observed anomalies than that obtained when the observed climate is used. Such an increase occurs because the 1986–90 model climate projects more strongly on the La Niña than on the El Niño flow pattern over the northern Pacific–North American region.

These results highlight the dependence of the skill scores, and in particular the interannual variability of the skill scores, on the choice of a reference climate. Even using an observed climate, the scores will be arbitrary to some degree, being dependent on which years are chosen to form the climate fields. Because of this, results are shown for most of the body of this paper, not in terms of the skill of anomalies, but rather in terms of the skill of differences between pairs of chosen years (e.g., between an El Niño year and a La Niña year). For such measures, we do not need to refer to any sample climatology.

b. Skill-score distributions of difference fields

The ensemble distribution of skill scores was estimated by comparing simulated and observed seasonal-mean difference fields for the pairs of years shown in Table 2. Specifically, for a given season let $E^1 = \{e_i^1\}$ and $E^2 = \{e_i^2\}$ denote two ensembles, the first for a year taken from the third column of Table 2 (denoted “positive”), the second corresponding to the year shown on

the same row in the fourth column of Table 2 (denoted “negative”). The differences $(e_i^1 - e_i^2)$ are then correlated with the corresponding observed difference field, $O^1 - O^2$, for all combinations of subscripts i, j , where i and j run from 1 to N , and N denotes the ensemble size. For two nine-member ensembles, there are 81 such difference fields $(e_i^1 - e_j^2)$ with 81 corresponding correlation coefficients. The distribution of relative frequencies of these correlation coefficients is then computed by binning them into categories of equal correlation intervals of 0.2. The distributions have been computed corresponding to all pairs of years from the rows of Table 2, and for nine regions, which together cover the globe (see Fig. 1). For the six extratropical regions (three in the Northern Hemisphere, three in the Southern Hemisphere), the skill scores are shown for the 500-mb height field. For the three tropical regions, they are shown for the 200-mb zonal wind field (geopotential height being rather featureless in the Tropics).

1) STRONG ENSO-INDEX YEARS

Figure 3 shows the distribution of skill scores for strong ENSO-index years for DJF over the nine regions. Categories between -1 and $+1$ are depicted on the x axis; on the y axis the relative frequency of skill scores (%) in each category is shown. Compared with the extratropics, the distributions are strongly peaked in the tropical regions. This is consistent with the relatively chaotic nature of the extratropics (e.g., Charney and Shukla 1981; Palmer 1996). In the region TR2, covering much of the tropical Pacific, the distribution is peaked toward the most skillful category, while in the other two tropical regions (TR1 and TR3), the distributions are strongly peaked toward the second most skillful category. Such a well-defined shift in the skill scores in the regions TR1 and TR3 is presumably associated with model error. Overall, these results are consistent with the fact that the dominant signal is associated with ENSO itself.

For the region NH2, the distribution is clearly peaked toward the two categories with largest correlation values. However, for other Northern Hemisphere regions, farther away from the ENSO signal, the distribution of skill scores is both broader and shifted toward less skillful categories. Nevertheless, the distributions are clearly skewed toward positive values, indicative of an overall level of skill. The distributions are also broad in the Southern Hemisphere, and, as for the Northern Hemisphere, the most skillful region (SH2) lies closest to the ENSO region. The shift toward negative skill values in SH3 may again be indicative of the influence of model error, but it also may be related to the quality of ECMWF verifying analysis for the years in question in these data sparse areas.

The ensemble-mean difference skill scores have also been categorized in the same way as individual members' differences. They are shown as small vertical ar-

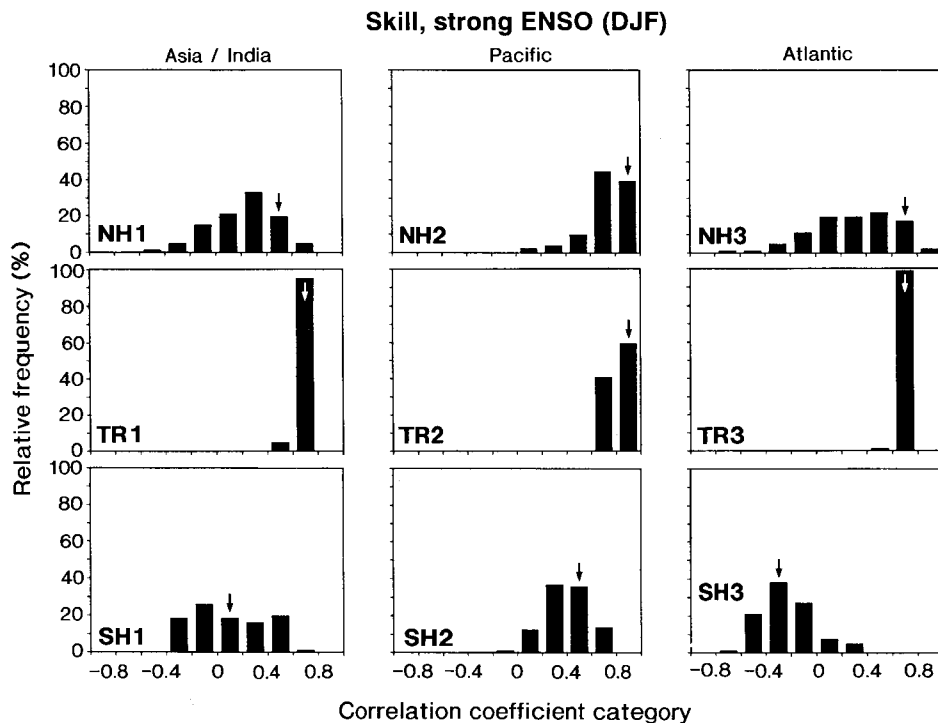


FIG. 3. Distributions of relative frequencies (%) of the difference skill scores over the nine regions of the globe (Fig. 1) for the strong ENSO-index DJF season. In the extratropics (top and bottom panels) correlation coefficients are computed for 500-mb height differences, in the Tropics (center row panels) for 200-mb zonal wind differences. Vertical arrows point to category of ensemble-mean difference skill score.

rows in Fig. 3 (and in subsequent skill score figures), pointing to the category they fall in. Generally, if the distribution is skewed strongly toward positive values, there is a tendency for the ensemble-mean score to either fall within the same category as distribution's peak value or to reside in the adjacent higher category of the peak value. This reflects the fact that the spatial variance of individual runs is larger than the spatial variance of the ensemble mean field and, consequently, ensemble averaging will increase the value of the ACC (Branković et al. 1990). If the distribution is skewed strongly toward negative values (as found for some weak-ENSO cases, see discussion and figures below), the ensemble-mean score may again coincide with distribution's peak value or fall in the adjacent lower category (here, ensemble averaging can make poor scores even worse!).

Figure 4 shows the skill score distributions in the three northern extratropical regions for the strong ENSO-index years, for the three remaining seasons, MAM, JJA, and SON. It is interesting to note that while in NH2 (the northern Pacific), the distributions are skewed to the most skillful categories in DJF (Fig. 3), in NH1 and NH3, the distributions are most skillful in MAM. It is not clear at present whether this is associated with the additional influence of more local (extratropical) lower boundary forcing anomalies in this season. The relatively high levels of skill around the Northern

Hemisphere in spring may have some important practical consequences in the application of seasonal forecasts to agricultural production.

In general, it can be seen that the broadest distributions of skill scores occurs in SON. The distributions for NH1 and NH3 in SON are a particularly striking illustration of the nature of internal chaotic variability in the atmosphere. Within the ensemble difference fields, there are (pairs of) members with skill scores exceeding 0.8, and yet others with skill scores between -0.6 and -0.8 .

2) WEAK ENSO-INDEX YEARS

Figure 5 shows the DJF skill score distributions for the weak ENSO-index years for all nine regions. It can be seen that for these years, comparing with Fig. 3, there is no clear-cut difference in the distributions in the Tropics and extratropics. In general, it is clear that these weak-ENSO years are associated with much weaker levels of skill. In the Tropics, both TR1 and TR3 are strongly shifted toward negative skill scores. Once more, this would appear to be associated with the impact of model error. It is interesting to note that despite a relatively broad distribution, the NH2 region still shows a shift toward positive skill values.

Figure 6 shows the Northern Hemisphere distribu-

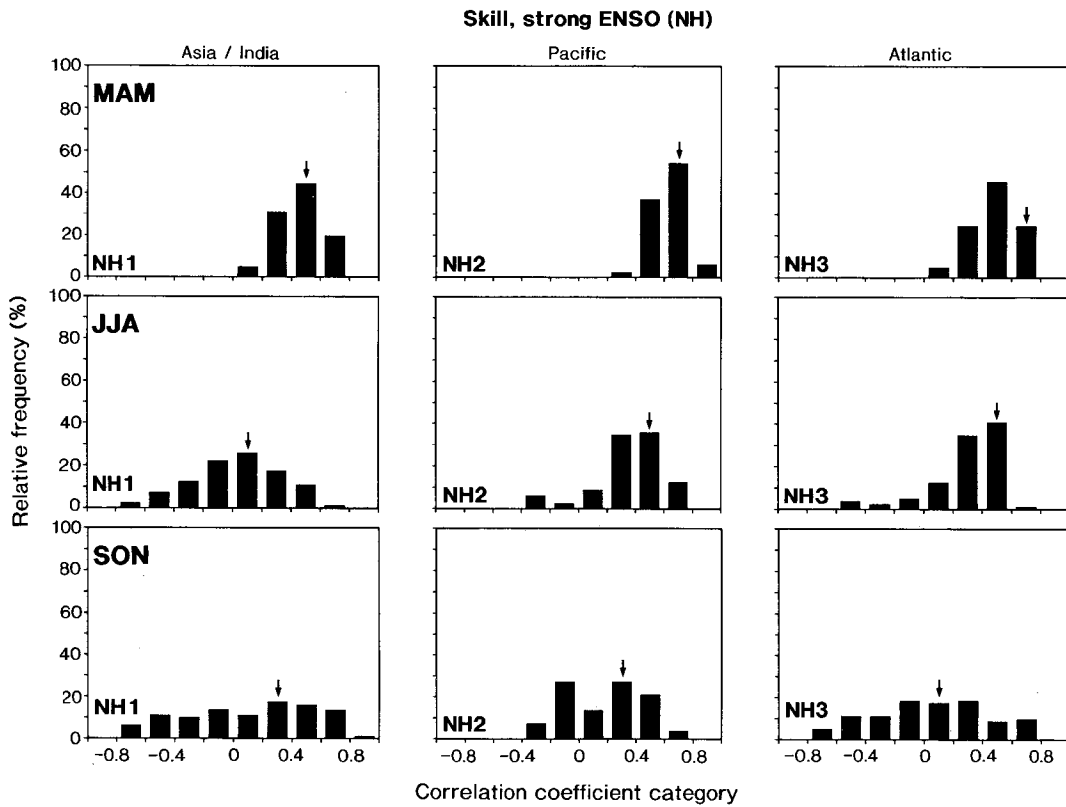


FIG. 4. Same as Fig. 3 but for the strong ENSO-index MAM (top row), JJA (middle), and SON (bottom) seasons over the Northern Hemisphere regions (NH1, NH2, NH3) only.

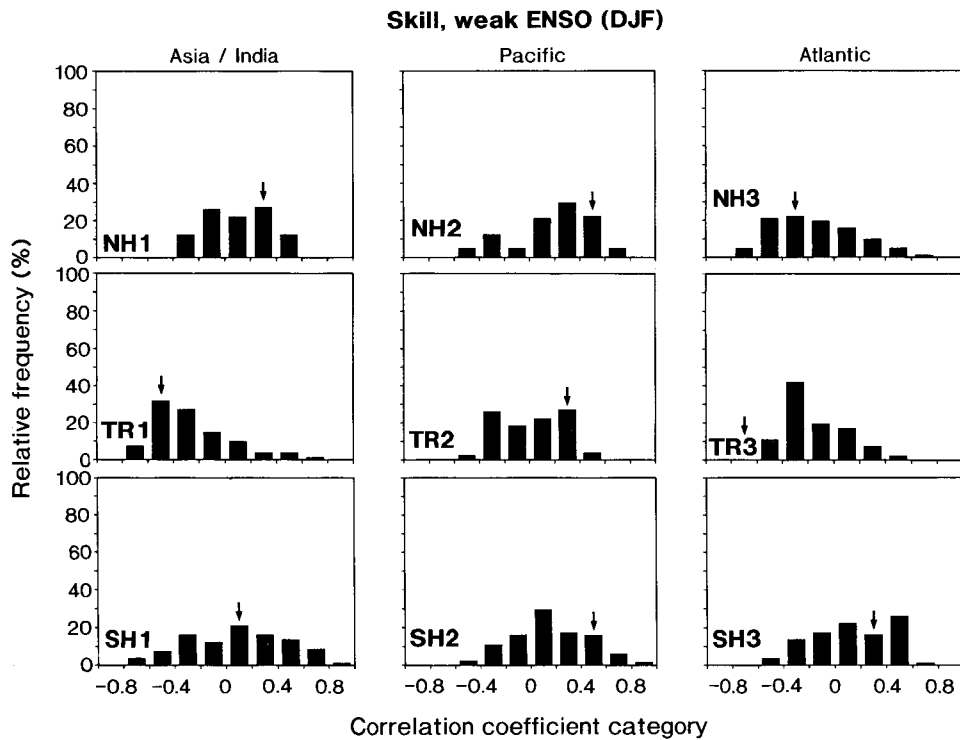


FIG. 5. Same as Fig. 3 but for the weak ENSO-index DJF season.

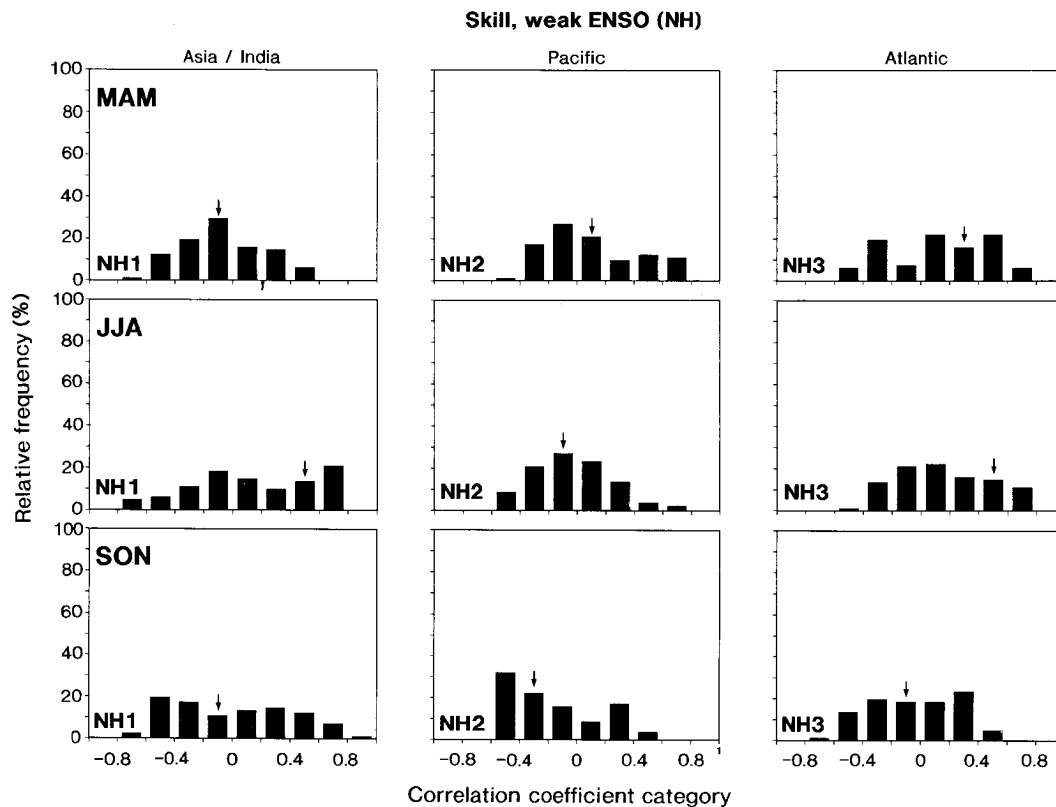


FIG. 6. Same as Fig. 4 but for the weak ENSO-index MAM, JJA, and SON seasons.

tions for the other three seasons for the weak ENSO-index years (these can be compared with Fig. 4). Although the distributions are broad, there is some evidence of a shift toward positive skill values, especially in the Atlantic (NH3) region in MAM and JJA. It is possible [see section 3b(1)] that lower boundary forcing, local to NH3, may be having an influence on the predictability in that region.

c. Consistency of ensembles

As mentioned a number of times above, distributions of skill scores are influenced by model error. To study the impact of SST variations on distributions of “internal” ensemble differences, we again consider two ensembles $\{e_i^1\}$ and $\{e_i^2\}$ from the third and fourth columns (respectively) of Table 2. For a given difference field, $e_i^1 - e_j^2$ (i th element from E^1 and j th element from E^2), we calculate $N \times N - 1$ correlation coefficients between this field and all possible pairs of difference fields, $e_k^1 - e_l^2$, where $k, l = 1, N$. This calculation is performed $N \times N$ times for all different correlations $C(e_i^1 - e_j^2, e_k^1 - e_l^2)$. For nine-member ensembles, we obtain 80×81 correlation coefficients. As for the skill scores, correlation coefficients measuring consistency have been calculated for the 500-mb height in the extratropical regions and for the 200-mb zonal wind in the Tropics. The distribution of relative frequencies of correlation

coefficients is estimated, as with skill scores, by binning into equal categories of 0.2.

Although not shown, in strong ENSO-index years, the consistency between ensembles is very high in the Tropics, regardless of season. Generally speaking, values are strongly peaked in the highest correlation category. In the extratropics, the distributions are much broader. Figure 7 shows the distributions for the Northern Hemisphere regions for all four seasons. (The meaning of the graph axes is the same as in Figs. 3–6.) Interestingly, it can be seen that for all regions MAM has the highest consistency. We noted above that MAM was the most skillful for NH1 and NH3. For the Atlantic sector (NH3), both MAM and JJA are more consistent than DJF.

For the weak ENSO-index years (Fig. 8), there is a clear shift toward positive correlation values; however, the overall distributions are much broader than for the strong ENSO-index years. The DJF season shows a stronger signal than JJA and SON, and in the NH2 area the shift toward positive values is stronger than in the other Northern Hemisphere regions.

If we compare consistency distributions with skill score distributions, it can be seen that there are some clear similarities. For example, both sets of distributions are broader for the weak ENSO-index years. Moreover, the NH2 has the tightest distributions for both consistency and skill. In general, the MAM season is seen as the most

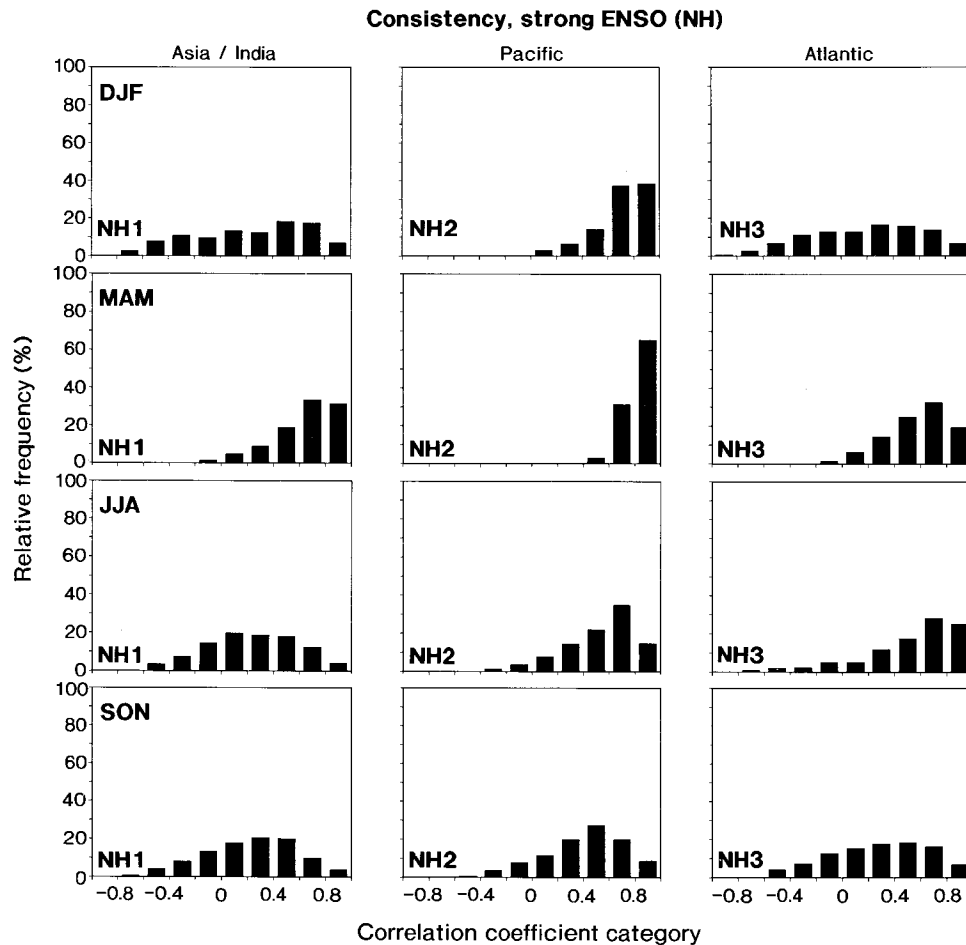


FIG. 7. Distributions of relative frequencies (%) of the 500-mb height consistency correlation coefficients over the Northern Hemisphere regions (NH1, NH2, NH3) for the strong ENSO-index DJF (top row panels), MAM (second row), JJA (third row), and SON (bottom row) seasons.

consistent and most skillful season at least for the strong ENSO-index years. In NH3, both MAM and JJA are consistent and skillful for strong ENSO-index years.

4. Potential predictability and ensemble size

In this section we focus on confidence values associated with the SST forcing of precipitation and near-surface temperature (2-m postprocessed temperature) as a function of ensemble size. These estimates are made for a number of prespecified land subregions, shown in Fig. 1 as shaded squares. (These have been chosen as representative examples from a much larger set of regions for which calculations have been made.) The chosen regions include some tropical areas where predictability may be expected to be high, together with extratropical areas where internal atmospheric variability may be expected to obscure, at least partially, the influence of lower boundary forcing.

The method of analysis is as follows. For each ensemble, there is a unique nine-member ensemble-mean

value for the regionally averaged precipitation or 2-m temperature. On the other hand, there are 84 possible values for a three-member subensemble-mean of the same precipitation or 2-m temperature from the original nine-member ensemble (corresponding to the number of ways of choosing a three-member subset from a nine-member set). The number of possible values for a 4-, 5-, 6-, 7- or 8-member subensemble is 126, 126, 84, 36, and 9, respectively. For each n -member subensemble ($3 \leq n \leq 9$), we calculate the subensemble-mean difference between years in the third and fourth columns in Table 2, and the corresponding t statistic based on the null hypothesis H_0 that the subensembles are not significantly different. We then calculate a mean t statistic by averaging over all possible subensemble t values. In performing this averaging, the sign of the t variable is ignored, because the sense of the interannual variation has no relevance to the discussion here.

For reference, in Figs. 9–11 the t value corresponding to the rejection of H_0 at the 90% and 99% confidence levels is shown. We comment on the minimum ensemble

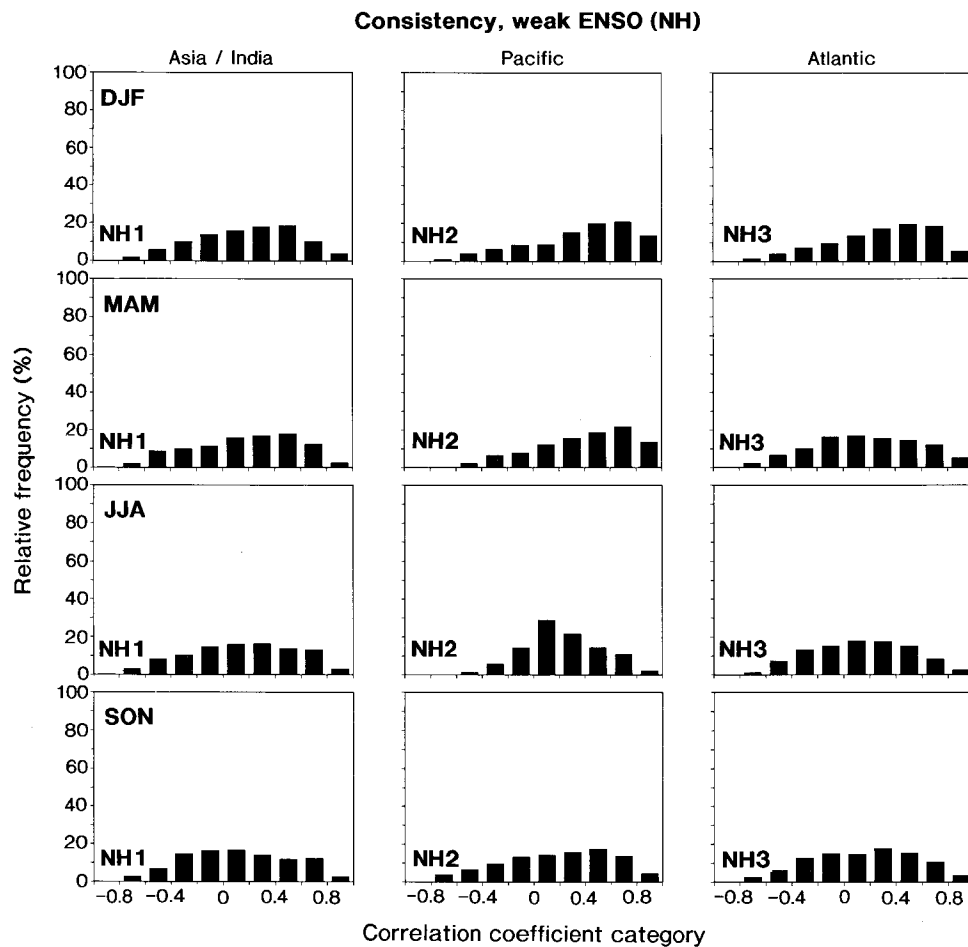


FIG. 8. Same as Fig. 7 but for the weak ENSO-index DJF, MAM, JJA, and SON seasons.

sizes required to exceed these t values, for different seasons, regions, and magnitude of the ENSO index. To aid the discussion below, we have also remarked on likely t values for hypothetical larger ensembles based on an extrapolation of the computed values. In all diagrams we note an almost linear relationship between ensemble size and the t value. This is due to a linear increase in the number of degrees of freedom in calculating t values as the size of the ensemble increases.

a. Extratropics

1) EUROPE

Figures 9a–d show the t values for 2-m temperature in the northern and southern European regions. For strong ENSO-index years (Figs. 9a,b), it can be seen that MAM and JJA are the most predictable seasons. These results are broadly in agreement with the NH3 500-mb height consistency distributions discussed in the previous section. For MAM, it appears that only four-member ensembles are required to reject H_0 with 99% confidence level. For SON and DJF, large ensemble sizes appear necessary to distinguish the two years.

In the weak ENSO-index years (Figs. 9c,d), only the MAM season in northern Europe seems to become predictable when the ensemble size increases significantly over nine members (again in agreement with NH3 height consistency diagrams). For this season and region, H_0 would appear to be rejected with 90% confidence, with approximately 20-member ensemble.

Figures 9e,f show the t values for rainfall for the northern and southern European regions for strong ENSO-index years. Consistent with the 2-m temperature results, MAM and JJA show evidence of potential predictability. However, for this variable, larger ensembles are required to reject H_0 with 99% confidence (in excess of 16 members for the northern Europe spring rainfall). The t values for weak ENSO-index rainfall are not shown but are generally smaller than those found in the strong ENSO-index years.

2) USA

For reason of space, predictability estimates for the USA are shown in Fig. 10 for the strong ENSO-index years only. The near-surface temperature appears to be

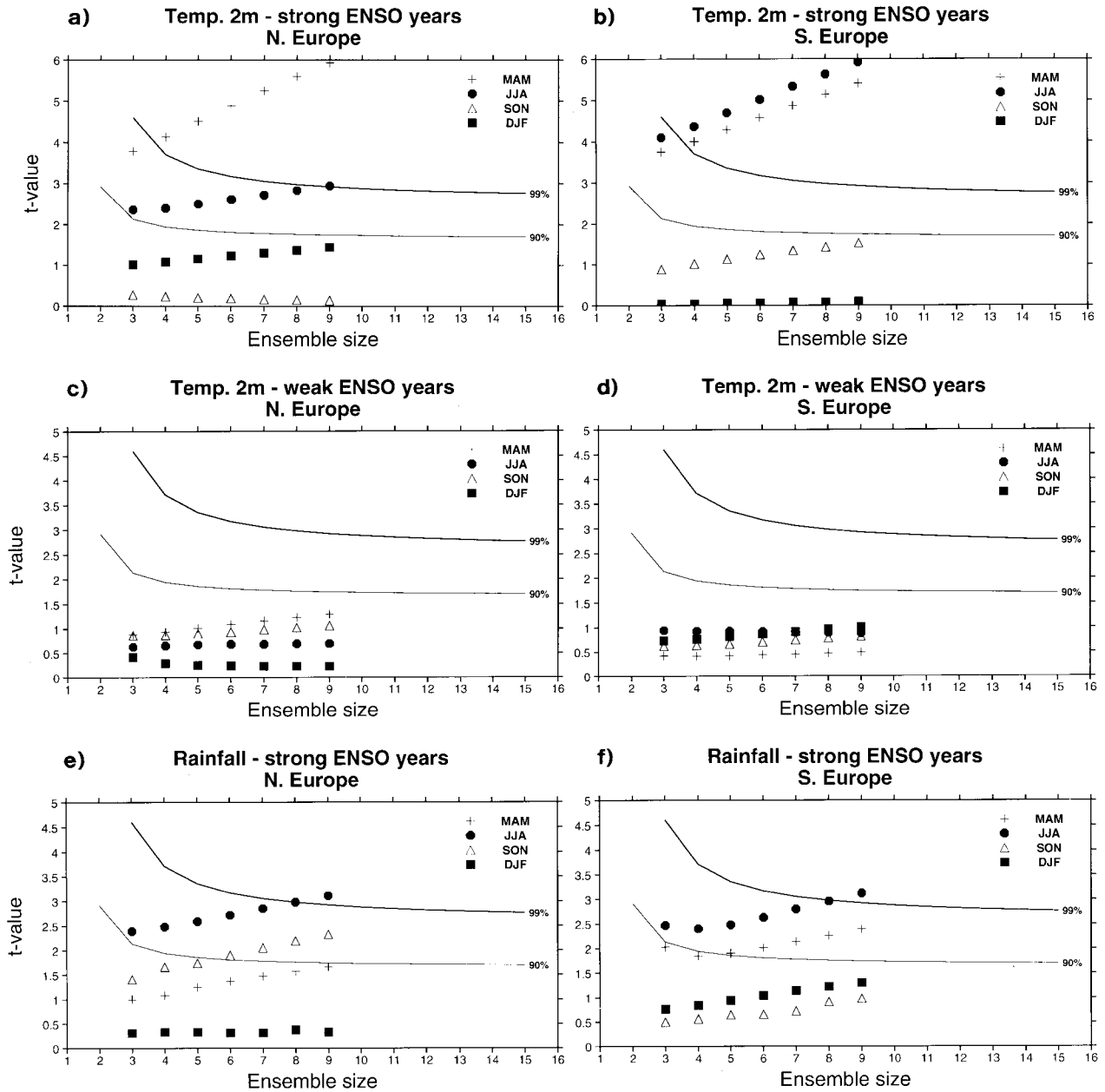


FIG. 9. The dependence of t values on ensemble size over northern Europe (left) and southern Europe (right). (a) and (b) Two-meter temperature, strong ENSO-index years; (c) and (d) 2-m temperature, weak ENSO-index years; (e) and (f) rainfall, strong ENSO-index years.

fairly predictable in general (Figs. 10a,b). Apart from JJA in the western United States, all seasons reach 90% confidence level of predictability with nine-member ensembles.

By contrast, during the weak ENSO-index years (not shown), significant predictability of 2-m temperature is found for the western United States only in summer, and in the eastern USA for winter.

For both the western and eastern United States rainfall (Figs. 10c,d), there is a dramatic seasonal cycle effect for the strong ENSO-index years, with spring (MAM)

showing much more significant values than for other seasons. For other seasons, there is relatively little difference between strong and weak ENSO-index years (not shown).

b. Tropics

To assess the relationship between predictability and ensemble size in the Tropics, we discuss the t value for the rainfall in four different regions (Sahel, India, Brazilian Nordeste, and northern Australia; see Fig. 1

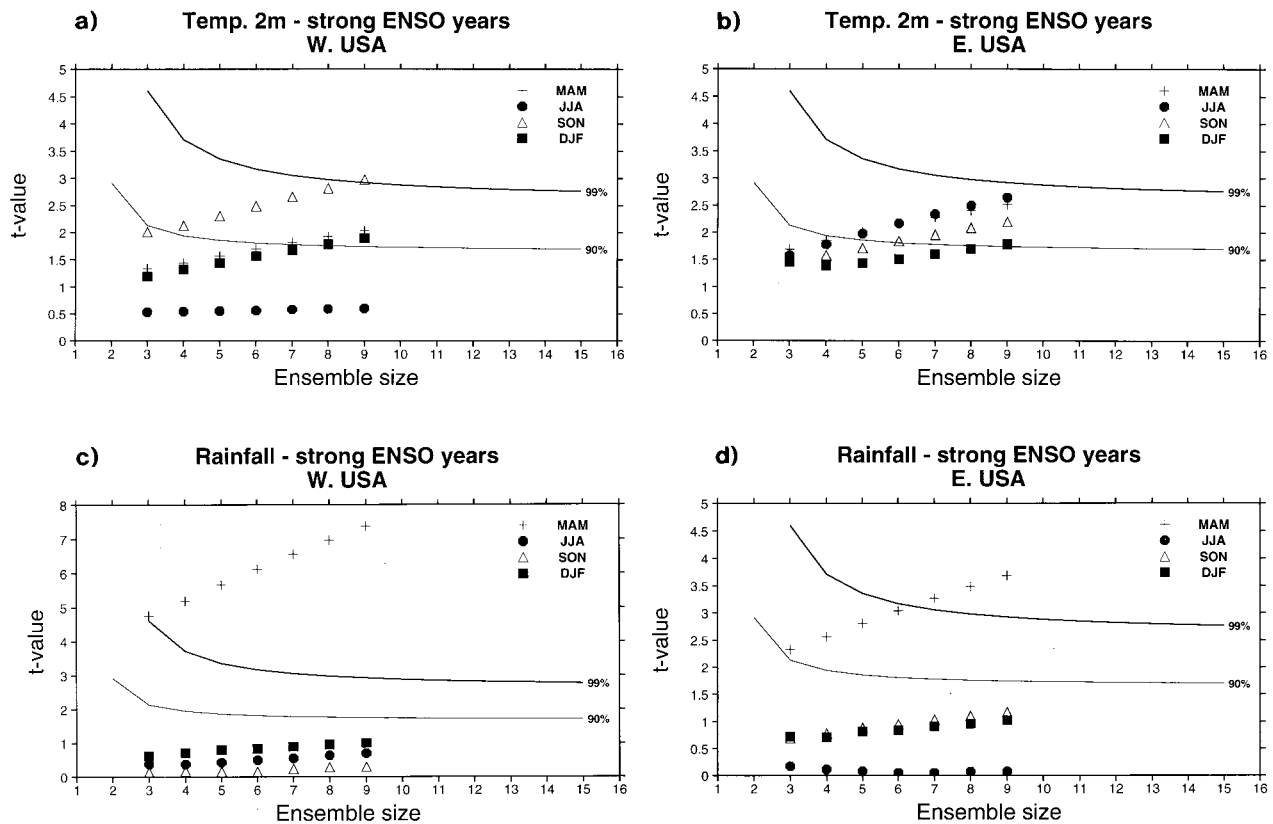


FIG. 10. Same as Fig. 9 but for the western United States (left) and eastern United States (right) and for the strong ENSO-index years only. (a) and (b) Two-meter temperature, (c) and (d) rainfall.

for the regions' boundaries) during the strong ENSO-index years (Fig. 11). As discussed by Charney and Shukla (1981), tropical rainfall should be intrinsically more predictable than extratropical rainfall. Clearly, as inferred from Fig. 11, in the Tropics relatively small-size ensembles are required to reach significant levels of predictability. Note the variable range of t values on the ordinate axes in Fig. 11. (In the Sahel during DJF there was no rainfall, hence all t values are zero.) The high level of skill over the Sahel for JJA and over the Nordeste for MAM is consistent with results from earlier studies (e.g., Rowell et al. 1995; Ward and Folland 1991). The Indian region has the highest proportion of t values below the 99% (and 90%) confidence level, indicating relatively lower predictability than in the other tropical regions shown. This could be associated with a model tendency to underestimate the overall level of the Indian rainfall, as shown in Branković and Palmer (1994), or with the influence of quasi-chaotic intraseasonal monsoon variations (Palmer 1994).

For the weak ENSO-index years (not shown), it can be noted that rather high levels of predictability were found for Nordeste and north Australian rainfall, while for the Sahel and India a noticeable deterioration in predictability for most seasons was found.

c. Summary for all regions

We have summarized the results for all seasons in Table 3. This shows the number of regions in the Tropics or the extratropics for which the rainfall or the near-surface temperature differences are statistically significant (i.e., H_0 is rejected) at the 90% confidence level, as a function of ensemble size. The regions in question are shown as shaded areas in Fig. 1. There are six regions in the extratropics (northern Europe, southern Europe, western United States, eastern United States, northeast China, and central China) and 6 in the Tropics (the Sahel, east Africa, northern Kalahari, India, northern Australia, and the Brazilian Nordeste). Table 3a is for the strong ENSO-index years; Table 3b is for the weak ENSO-index years.

In the strong ENSO-index years (Table 3a), only 9 out of 48 possible regionally averaged extratropical rainfall or temperature values are significant with 3-member ensembles.¹ This increases to 22 with a 9-member ensemble and a projected 32 (out of 48) with a 16-member

¹ The number 48 is obtained when the total number of extratropical regions considered (6) is multiplied by the number of seasons (4) and by the number of variables (2). The same is valid for the Tropics.

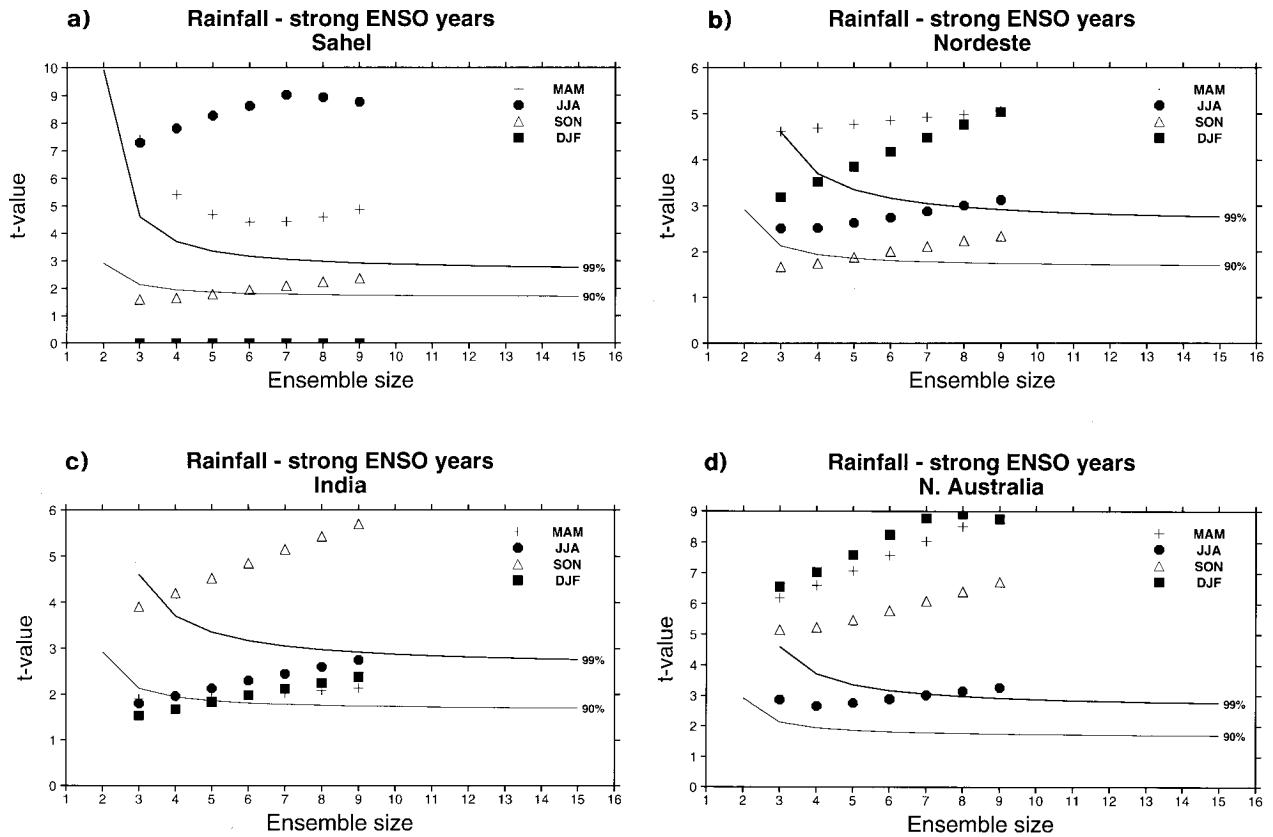


FIG. 11. Same as Fig. 9 but for selected tropical regions, rainfall, and strong ENSO-index years only.

ensemble. Hence, for extratropical prediction, the use of relatively large ensembles can be justified. For the tropical regions, the number of significant values increases from 31 to 40 to a projected 42 (out of 48) as the ensemble size increases from 3 to 9 to 16.

On the other hand, for the weak ENSO-index years (Table 3b), even tropical predictions appear to benefit from large ensembles, where the number of significant values increases from 15 to 24 to a projected 32 as the ensemble size increases from 3 to 9 to 16. For these weak ENSO-index years, a large ensemble benefits detection of extratropical predictability as well.

Although this information is not shown in Table 3, the greatest benefit from having a 16-member ensemble comes from detecting rainfall predictability in extratropical regions [for strong ENSO-index years from 5 to 12 to a projected 22 (out of 24) significant values with 3-, 9-, and 16-member ensembles, respectively].

5. Probability forecasts

As discussed in the introduction, the underlying SST fields can be thought of as influencing the geometry of the atmospheric attractor. We can represent this influence through changes in the PDF of atmospheric states. By focusing on changes to the PDF of specific weather

variables, such as rainfall, the discussion is relevant to the practice of operational seasonal weather prediction.

Figure 12 depicts the distribution of rainfall amounts over southern Europe and the Sahel for all individual integrations of the (strong ENSO-index) JJA 1987 (left-hand bars) and JJA 1988 (right-hand bars) ensembles. The experiment number on the x axis depicts individual members of ensembles with respect to their (ascending) initial dates shown in Table 1. For example, experiments denoted by the number 6 were initiated on 2 May of 1987 and 1988, respectively. The rainfall is averaged over land grid points only (the regions' boundaries are shown in Fig. 1). The variation of regional rainfall amounts within ensembles is seen in both diagrams, though the response of the model to the same SST forcing is much more stable in the Sahel. From such distributions we can estimate changes to rainfall PDFs.

It is interesting to note from Fig. 12a that a three-member subensemble {2, 7, 8} would actually give the opposite sign of the rainfall difference than other combinations of a three-member subensemble. However, for most of the three-member combinations in Fig. 12a, the average rainfall difference between JJA 1987 and JJA 1988 would be of the same sign and opposite to that for the {2, 7, 8} subensemble. This is why the value

TABLE 3a. Number of regions in the Tropics and extratropics (depicted as shaded areas in Fig. 1) for which interannual differences for either rainfall or 2-m temperature are statistically significant at 90% confidence level. Strong ENSO-index years.

Regions located at	Ensemble size		
	3	9	16
Extratropics	9	22	32
Tropics	31	40	42

for the JJA three-member ensemble shown in Fig. 9f reached a relatively high confidence level.

Figure 13a shows an example for the strong ENSO-index years over the African and south Asian region. From the incidence of both positive and negative rainfall differences, a gridpoint probability (or proportion of positive and negative differences) is assigned. For each grid point, the value obtained is an estimate of the probability that, for JJA, the El Niño year 1987 was wetter or drier than the La Niña year 1988. Probabilities for both positive and negative differences are then combined into one map with the discriminating contour of 60%.

Over the Sahel, the probability of negative rainfall difference (light stipple) is high, exceeding 90% over much of the region. This result implies that for most of the 81 differences, the ensemble simulations of JJA 1987 were drier than JJA 1988. This is in good agreement with verification differences for the two summers (Fig. 13b) obtained from the Global Precipitation Climate Centre (GPCC; Huffman et al. 1995).

Over much of India, the probability of negative differences is in excess of 60%, and more than 90% in the north and in the south. Generally smaller probability estimates over India imply lower predictability for that region than for the Sahel. This is associated with a relatively larger internal ensemble variability over the Indian subcontinent than over the Sahel and is consistent with estimates of the t variable shown in Fig. 11.

The GPCC rainfall verification data were available for years 1987 and 1988 only and at this stage it is difficult to validate precipitation forecasts in general; the production of global precipitation datasets is an ongoing effort. Hence, in order to validate such probability forecasts, we have created fields giving the probability that the local 500-mb geopotential height difference is either positive or negative. We would expect that for grid points enclosed by a given probability contour (here we choose the 60% probability contour), at least 60% of these grid points would validate correctly.

TABLE 3b. Same as Table 3a but for the weak ENSO-index years.

Regions located at	Ensemble size		
	3	9	16
Extratropics	2	10	15
Tropics	15	24	32

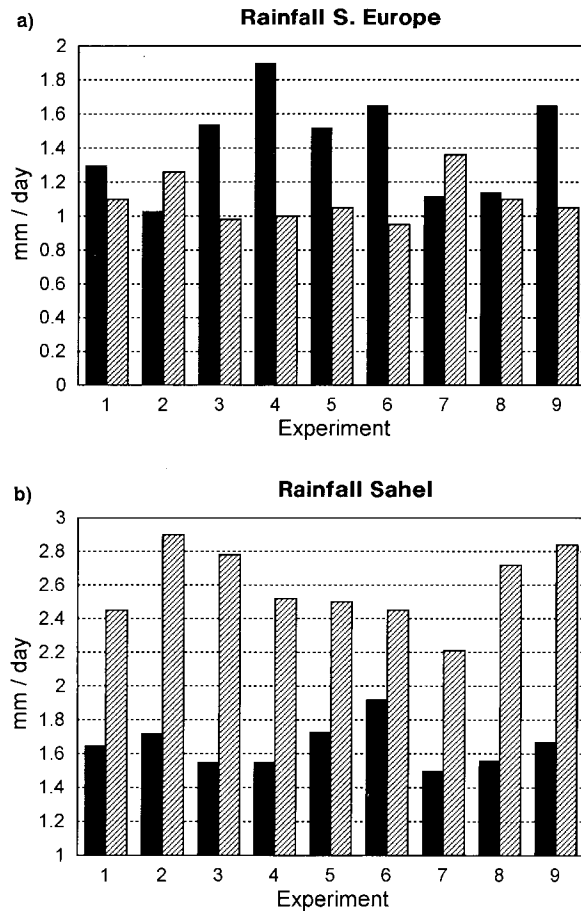


FIG. 12. Distribution of seasonally averaged rainfall (mm day^{-1}) in the JJA 1987 (left-hand, solid bars) and JJA 1988 (right-hand, hatched bars) ensembles for (a) southern Europe and (b) Sahel.

Table 4 shows a set of validations for the nine regions that span the globe (Fig. 1). For each region, and each season from the strong ENSO-index years, we create a 2×2 table. The two diagonal elements of each table show the percentage of points where (a) the probability of a positive difference exceeded 60%, and a positive difference occurred (top left element of the table), and (b) the probability of a negative difference exceeded 60%, and a negative difference occurred (bottom right element). These two diagonal elements essentially show the degree of agreement between the probability forecasts and verifying analyses. The off-diagonal elements correspond to forecasts that did not verify.

As an example, in Table 4 we show results for the strong ENSO-index MAM season. The best agreement between the model probabilistic fields and observed differences is found in the tropical Atlantic (TR3). In the Northern Hemisphere, this verification method gives satisfactory results for all three regions, whereas in the Southern Hemisphere, results are poor for SH1.

In JJA (and also in SON), the agreement in the three Northern Hemisphere regions is generally lower than in

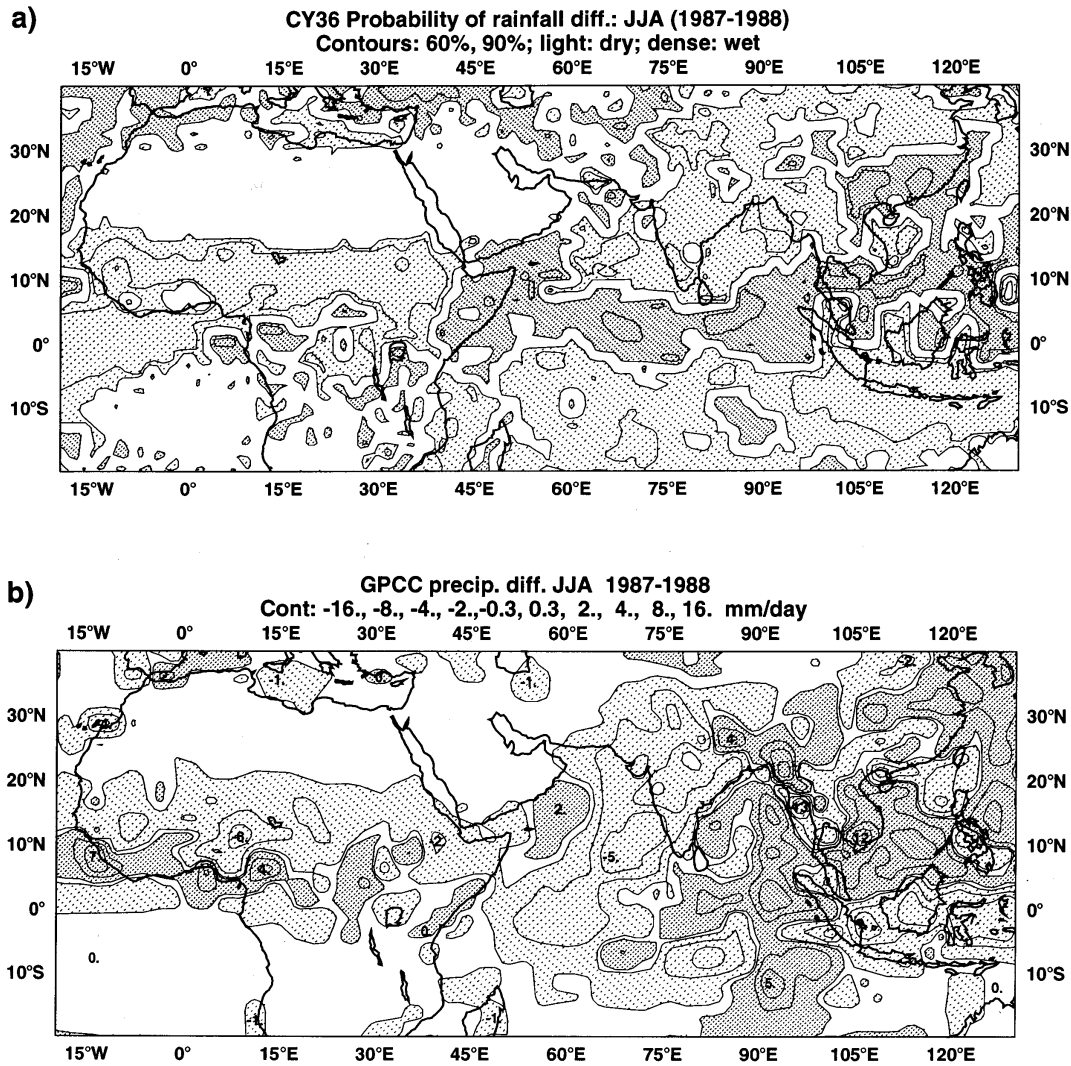


FIG. 13. (a) Probability of rainfall differences, JJA 1987 minus JJA 1988. Light stipple probability of negative differences, dense stipple probability of positive differences. Contours 60% and 90% for both probabilities. (b) The GPCP rainfall differences JJA 1987 minus JJA 1988. Contours are ± 0.3 , ± 2 , ± 4 , ± 8 , ± 16 mm day⁻¹.

the boreal cold seasons, DJF and MAM (not shown). In the Tropics, negative differences are entirely accurately predicted by the model, and in the Southern Hemisphere they are better predicted than positive differences.

6. Summary and conclusions

Results from a large set of 120-day, nine-member ensemble integrations of a T63L19 version of the ECMWF model have been presented. Individual ensemble members were initiated from consecutive operational ECMWF analyses, separated by 24 h. Integrations were made using specified observed SST, updated in the model every 5 days. The last three months of each individual integration, corresponding to conventional calendar seasons, were analyzed. This set of ensembles is

an extension of the three-member ensembles reported by Branković et al. (1994; BPF).

We focus on the ability of the model to simulate interannual atmospheric variations on seasonal timescales over the 5-yr period 1986–90. This period was characterized by significant variability in the El Niño–Southern Oscillation (ENSO). Based on ENSO, an index was defined that varied from large positive values in the first part of the 5-yr period to large negative values in the middle of the period with weak negative and weak positive values toward the end of the period. As in BPF, difference fields were computed from seasons when the ENSO index was large and opposite, and weak and opposite.

The skill of the model was first discussed in terms of anomaly correlation coefficients (ACCs). The ACCs have been derived with respect to the two different ref-

TABLE 4. Percentage of grid points for nine "global" regions (Fig. 1) for which probability of positive (top row) and negative (bottom row) 500-mb height difference exceeds 60%. Strong ENSO-index, MAM season.

Region	Anal diff > 0	Anal diff < 0
NH1		
Prob(+) > 60%	66%	34%
Prob(-) > 60%	21	79
NH2		
Prob(+) > 60%	59	41
Prob(-) > 60%	17	83
NH3		
Prob(+) > 60%	89	11
Prob(-) > 60%	25	75
TR1		
Prob(+) > 60%	97	3
Prob(-) > 60%	—	—
TR2		
Prob(+) > 60%	97	3
Prob(-) > 60%	33	67
TR3		
Prob(+) > 60%	95	5
Prob(-) > 60%	0	100
SH1		
Prob(+) > 60%	45	55
Prob(-) > 60%	88	12
SH2		
Prob(+) > 60%	84	16
Prob(-) > 60%	17	83
SH3		
Prob(+) > 60%	57	43
Prob(-) > 60%	37	63

erence fields (used in the computation of model anomalies). The first reference field was the observed climate, the second reference field was the model climate. It was shown that the model skill scores depend strongly on the choice of a reference climate. Therefore, in our analysis of skill scores we focus on difference fields between pairs of years with the opposite index of ENSO.

Distributions of ensemble skill scores for difference fields were calculated for nine regions over the globe. During strong ENSO-index years, the highest and most sharply peaked distribution of skill in the northern extratropics was found for the northern Pacific–North American region for the winter (DJF) season. Over the northern Atlantic–European region, the ensemble skill is highest and most sharply peaked in spring (MAM). This may have some important implications for the application of seasonal predictions in the growing season in Europe. In the tropical regions, skill is generally high with very sharply peaked distributions. In DJF and MAM, the highest skill is found in the tropical Pacific region, and in JJA and SON in the tropical Atlantic and tropical Indian Ocean regions.

For weak ENSO-index years, the distributions of skill are found to be generally broader than for the strong

ENSO-index years. However, for many regions there is a tendency of skill distributions to be skewed toward positive values. A shift toward negative correlation coefficients in the Southern Hemisphere may be associated with nonnegligible model systematic errors. For a given season, estimates of consistency for all possible differences between members of two ensembles were also made. These distributions can be thought of as giving some measure of intra-ensemble ("internal") variability. For strong ENSO-index years, consistency distributions in the tropical regions are peaked at high positive values. In the Northern Hemisphere, spring has the highest consistency between ensembles, similar to distributions of ensemble skill. For weak ENSO-index years the consistency distributions are generally broad, though clearly shifted toward positive correlations.

The above conclusions are generally consistent with those from BPF for three-member ensembles. However, some differences that may exist between the two papers could be attributed to a much poorer sampling in our earlier work. In BPF, for example, for the northern Atlantic–European area the model skill for both winter and spring seasons was found to be relatively high and almost identical. The increase in ensemble size as well as the analysis of skill distribution performed in this paper help to better discriminate between those two seasons.

Based on the above ensemble skill and consistency estimates, we can distinguish between the regions of the globe and seasons with relatively good prospects for seasonal prediction. Apart from the tropical regions in general, the northern Pacific–North American region in winter and the northern Atlantic–European region in spring appear to have such a potential during ENSO years. These results are not inconsistent with observational and empirical seasonal predictability studies (e.g., Halpert and Ropelewski 1992; Livezey 1990; Barnston 1994).

In this paper, an estimate, based on the t statistic, is given of the minimum size of an ensemble required to simulate with confidence the impact of the underlying SST anomalies on the probability distribution of atmospheric states. These t values were derived as a function of ensemble size, when the latter increases from three to nine members. In addition, an extrapolation of the t statistic for larger ensembles is discussed. This evaluation is performed for 2-m temperature and precipitation over a number of predefined regions in both extratropics and Tropics.

In general, relative large (≥ 20 member) ensembles may be needed for extratropical seasonal prediction of regional weather, even in the presence of a relatively strong tropical signal. On the other hand, in the Tropics during strong ENSO events, the same level of confidence can be attained with much smaller ensembles. However, these estimates may vary widely, depending on the region and season considered. For example, whereas the JJA rainfall prediction for the Sahel may

require only a two- to three-member ensemble, for India at least nine- or ten-member ensembles would be desirable. In weak ENSO years, a relatively large ensemble would be needed even for tropical regions.

Some examples of probability fields were shown, focusing on monsoon rainfall probability. An objective verification of such probabilities was given indicating good agreement between observed differences and probabilities inferred from model difference fields.

An extensive set of multimodel ensembles, based on a research project made jointly with the United Kingdom Meteorological Office, Météo France, and the French Electricity Board, currently in progress, will allow more conclusive investigation on the role of model formulation on seasonal predictability estimates.

Acknowledgments. We thank an anonymous reviewer for constructive and useful comments.

REFERENCES

- Barnett, T. P., 1995: Monte Carlo climate forecasting. *J. Climate*, **8**, 1005–1022.
- Barnston, A. G., 1994: Linear statistical short-term climate predictive skill in the Northern Hemisphere. *J. Climate*, **7**, 1513–1564.
- Branković, Č., and T. N. Palmer, 1994: Predictability of summer monsoons. *Proc. Int. Conf. on Monsoon Variability and Prediction*, Trieste, Italy, World Meteor. Org., 629–636.
- , —, F. Molteni, S. Tibaldi, and U. Cubasch, 1990: Extended-range predictions with ECMWF models: Time-lagged ensemble forecasting. *Quart. J. Roy. Meteor. Soc.*, **116**, 867–912.
- , —, and L. Ferranti, 1994: Predictability of seasonal atmospheric variations. *J. Climate*, **7**, 217–237.
- Charney, J. G., and J. Shukla, 1981: Predictability of monsoons. *Monsoon Dynamics*, J. Lighthill and R. Pearce, Eds., Cambridge University Press, 735 pp.
- Halpert, M. S., and C. F. Ropelewski, 1992: Surface temperature patterns associated with the Southern Oscillation. *J. Climate*, **5**, 577–593.
- Huffman, G. J., R. F. Adler, B. Rudolf, U. Schneider, and P. R. Keehn, 1995: Global precipitation estimates based on a technique for combining satellite-based estimates, rain gauge analysis, and NWP model precipitation information. *J. Climate*, **8**, 1284–1295.
- Kumar, A., and M. P. Hoerling, 1995: Prospects and limitations of seasonal atmospheric GCM predictions. *Bull. Amer. Meteor. Soc.*, **76**, 335–345.
- Livezey, R. E., 1990: Variability of skill of long-range forecasts and implications for their use and value. *Bull. Amer. Meteor. Soc.*, **71**, 300–309.
- Miller, M. J., A. C. M. Beljaars, and T. N. Palmer, 1992: The sensitivity of the ECMWF model to parametrization of evaporation from the tropical oceans. *J. Climate*, **5**, 418–434.
- Palmer, T. N., 1993: Extended-range atmospheric prediction and the Lorenz model. *Bull. Amer. Meteor. Soc.*, **74**, 49–65.
- , 1994: Chaos and predictability in forecasting the monsoons. *Proc. Indian Nat. Sci. Acad.*, **60**, 57–66.
- , 1996: Predictability of the atmosphere and oceans: From days to decades. *Decadal Climate Variability: Dynamics and Predictability*, NATO ASI Series, Vol. I 44, Springer, 83–155.
- , and D. L. T. Anderson, 1994: The prospects for seasonal forecasting—A review paper. *Quart. J. Roy. Meteor. Soc.*, **120**, 755–793.
- Rowell, D. P., C. K. Folland, K. Maskell, and M. N. Ward, 1995: Variability of summer rainfall over tropical north Africa (1906–92): Observations and modelling. *Quart. J. Roy. Meteor. Soc.*, **121**, 669–704.
- Simmons, A., D. M. Burridge, M. Jarraud, C. Girard, and W. Wergen, 1988: The ECMWF medium-range prediction models: Development of the numerical formulations and the impact of increased resolution. *Meteor. Atmos. Phys.*, **40**, 28–60.
- Stern, W., and K. Miyakoda, 1995: Feasibility of seasonal forecasts inferred from multiple GCM simulations. *J. Climate*, **8**, 1071–1085.
- Ward, M. N., and C. K. Folland, 1991: Prediction of seasonal rainfall in the north Nordeste of Brazil using eigenvectors of sea surface temperature. *Int. J. Climatol.*, **11**, 711–743.