

## An assessment of seasonal predictability using atmospheric general circulation models

By R. J. GRAHAM\*, A. D. L. EVANS, K. R. MYLNE, M. S. J. HARRISON and K. B. ROBERTSON  
*The Met. Office, UK*

(Received 26 February 1999; revised 25 February 2000)

### SUMMARY

Seasonal predictability is investigated using a 15-year set of 4-month range, 9-member ensemble integrations from atmospheric general circulation models (AGCMs) involved in the European PROVOST project (Prediction Of climate Variations On Seasonal to interannual Time-scales). The integrations were performed using prescribed ideal (observed) sea surface temperatures (SSTs), therefore skill attained (referred to as 'potential' skill) represents an estimated upper bound on skill achievable with current models using predicted SSTs. Most analysis is presented for The Met. Office Unified Model (UM), the European Centre for Medium-Range Weather Forecasts (ECMWF) T63 model (T63) and an 18-member multiple-model ensemble (JT2) constructed from these individual models. The benefits of higher-order multiple models (employing all four participating PROVOST AGCMs) are also investigated. Evaluation is focused on four assessment regions: the tropics; the northern extratropics; Europe; and North America. Probabilistic skill is assessed for the basic events: 3-month mean 850 hPa temperature above/below normal; 3-month mean precipitation accumulation above/below normal. Deterministic (ensemble mean) skill is also assessed. A summary of the main results is provided below.

- **Potential skill:** Skill scores for months 1–3 forecast 850 hPa temperature and precipitation calculated for the entire tropical and northern extratropical regions indicate that, while skill is highest in the tropics, it is also available over the northern extratropics in all seasons. Scores for the northern extratropics are highest in spring (March–April–May; MAM). Scores for precipitation are generally lower than for 850 hPa temperature, however, there is evidence of substantial potential for rainy season predictions in some tropical regions. Over Europe and North America skill scores for 850 hPa temperature are (for at least one of the UM, T63 and JT2 models) comparable to those of the northern extratropics in all seasons. Peak skill occurs over Europe in MAM (as found for the northern extratropics). In contrast, peak skill over North America occurs in December–January–February (DJF), apparently as a result of enhanced predictability during El Niño Southern Oscillation (ENSO) events. In non-ENSO years skill over Europe and North America is similar, suggesting that the greater predictability often attributed to the North American region relative to Europe may apply only during ENSO events. Skill for months 2–4 is generally lower than for months 1–3, though there is evidence that during ENSO events levels of skill in the first three months are maintained into the second three months. For precipitation, best skill over Europe and North America is found in MAM and DJF, with little evidence of any skill over Europe in summer and autumn.

- **Skill prediction:** Largest ENSO-related skill enhancements over North America are found in DJF and over Europe in the following (post-ENSO peak) MAM. Ensemble spread appears a useful indicator of ensemble-mean skill in some seasons over Europe and North America. Thus prospects for skill prediction appear promising, perhaps using strategies which combine information on both the state of ENSO and ensemble spread.

- **Benefits of multiple-model ensembles:** Multiple-model ensembles enhance prediction capabilities, allowing the strengths of the individual AGCMs to be exploited without extensive a priori calibration of each model. The multiple-model ensembles frequently provide a filter for the more skilful individual model (the identity of which varies with season and region). The key factor determining the skill of the multiple model appears to be the skill of the most skilful component ensemble, and does not appear to be strongly connected with the increased ensemble size.

- **Use of persisted SST anomalies:** Tests indicate that a substantial proportion of the skill achieved using observed SSTs is retained using persisted SST anomalies (SSTA) from the month preceding the initial date of the integration, indicating that use of persisted SSTA is a viable method for real-time seasonal prediction, at least for up to one season ahead.

- **User value:** A methodology for linking technical forecast quality with financial value for users has been outlined using the relative operating characteristic and the user cost/loss matrix. Results indicate promising potential for user value of probabilistic seasonal predictions not only over tropical areas but also in some extratropical areas, including Europe.

KEYWORDS: AGCM ensembles ENSO Multiple-models Seasonal predictability User value

### 1. INTRODUCTION

The scientific basis for seasonal prediction stems primarily from evidence that the atmosphere's lower boundary, particularly the sea surface temperature (SST), influences

\* Corresponding author: Ocean Applications (Seasonal Modelling and Prediction Group), CSc Division, The Met. Office, London Road, Bracknell, Berkshire RG12 2SZ, UK.

© Crown copyright 2000.

weather regime frequency statistics (Palmer and Anderson 1994), and consequently the seasonal-mean weather conditions. The SST field itself evolves slowly compared with individual synoptic-scale weather systems, and is often relatively predictable (at least in the tropics)—thus representation of the SST evolution in atmospheric general circulation models (AGCMs), so-called dynamical seasonal prediction, potentially provides a means of generating forecasts of seasonal-average weather. This paper describes investigations to evaluate that potential.

The work has been performed as part of the recent European PROVOST project (PRediction Of climate Variations On Seasonal to interannual Time-scales) which has provided an extensive simulation dataset that expands in a number of ways on those previously used. Briefly, the dataset (described in section 2) comprises 4-month range, 9-member ensemble AGCM integrations run over ideal (observed) SSTs, for all four seasons over the 15-year period 1979–93. The period contains five El Niño Southern Oscillation (ENSO) events, allowing assessment of the impact of ENSO on predictability. Moreover, integrations were repeated at four European centres using different AGCMs, allowing opportunities for comparing model performance and for assessing the potential performance benefits from combining predictions from different models. Here we describe research performed at The Met. Office as part of the PROVOST project. The key objectives and scope of the study are expanded below.

(i) *Assessment of AGCM ‘potential’ skill and benefits of multiple-model ensembles.* A main objective, addressed in section 3, is to obtain estimates of the upper-bound on skill (referred to as ‘potential’ skill) achievable using coupled AGCMs by measuring the skill obtained when uncoupled AGCMs are forced with observed (i.e. near-perfect predicted) SST. Emphasis is placed on assessment and comparison of The Met. Office Unified Model (UM; Cullen 1991) and the European Centre for Medium-Range Weather Forecasts (ECMWF) T63 ensembles (hereafter referred to as the UM and T63 ensembles). The purpose of the intercomparisons is to gain insight into the potential enhanced capability available through combining the strengths of two or more AGCMs.

Studies of medium-range prediction have shown that multiple-model ensembles, comprising a combination of members run with different AGCMs, provide significant performance benefits relative to the component individual ensembles (e.g. Evans *et al.* 2000). Here we use multiple-model configurations of up to four AGCMs to assess whether such benefits extend to the seasonal range. Benefits from combining ensembles derive potentially from both the inclusion of complementary predictive information (Brown and Murphy 1996; Evans *et al.* 2000) and the increased ensemble size.

(ii) *Assessment of the prospects for skill prediction.* Large amplitude warm/cold SST events associated with the El Niño/La Niña phases of the El Niño Southern Oscillation (ENSO) in the tropical Pacific may, through persistent thermal forcing on the atmosphere, give rise to enhanced predictability in some regions. Hereafter in this report we will refer to warm SST events as PW (Pacific Warm), cold events as PC (Pacific Cold) and both sets of events collectively as PC/W, to emphasise that it is specifically the presence (or otherwise) of anomalies in the tropical Pacific SST that is of interest. The degree to which established or developing PC/W events provide an advance indicator of relatively high seasonal prediction skill is investigated in section 4 for Europe and North America. The degree to which ensemble spread may be used as a predictor of ensemble-mean skill is also investigated.

(iii) *Assessment of skill available using persisted SST anomalies for boundary forcing.* Use of persisted SST anomalies (SSTA) for lower-boundary forcing is likely to be a competitive (and cheap) option relative to coupled ocean–atmosphere models, at least for one season ahead, because of the usual relatively slow evolution of the SST field. In section 5 we test the viability of using persisted SSTA for operational real-time prediction by comparing skill from hindcasts against the skill benchmarks obtained from the ‘near-perfect SST’ PROVOST simulations.

(iv) *Evaluation of the potential value of seasonal predictions to users.* For operational seasonal prediction to be a viable concern, it will be necessary to establish that potential users will be able to extract benefit from the forecasts, given the levels of technical skill that are available. To this end we outline in section 6 a methodology for estimating the potential user financial value. The method is based on the user cost/loss matrix associated with the outcome of probabilistic and deterministic predictions of specified weather events.

## 2. EXPERIMENTAL DETAILS

### (a) *The PROVOST simulations*

Ensemble integrations to a range of four months were performed for each season in the 15-year period December 1979 to March 1994. The integration periods were specified as: March to June (northern spring); June to September (northern summer); September to December (northern autumn); and December to March (northern winter). Each season was simulated with 3 different AGCMs: the UM run at climate resolution ( $3.75^\circ$  longitude by  $2.5^\circ$  latitude with 19 levels); the ECMWF T63 L31 model (referred to here as T63); and the Météo-France ARPÈGE T42 L31 model (referred to as AP1). A fourth set of integrations was run, for the winter season only, by Electricité de France using the ARPÈGE model at T63 L31 truncation (referred to as AP2). The Met. Office integrations were made using version 3.4 of the UM with HADAM2b physics; validation of this version of the climate model has been discussed in Hall *et al.* (1995).

Integrations with each model were made in 9-member ensembles initialized with the 1200 UTC analyses from the ECMWF re-analysis (ERA, Gibson *et al.* 1997) on the nine consecutive days before each season. Observed values of SST and ice cover, used for lower-boundary forcing, were updated at 5-day intervals during the integration. The observed SSTs and ice cover were obtained from The Met. Office Global sea-Ice and Sea Surface Temperature (GISST) analyses (Rayner *et al.* 1996) up to October 1981 and the Reynolds optimal interpolation (OI) analyses (Reynolds and Smith 1994) for the remaining period. In the UM experiments land surface initial conditions from climatology were specified, both for the PROVOST simulations and the persisted SSTA hindcasts.

The multiple-model configurations studied are: an 18-member combination of the UM and T63 ensembles (referred to as JT2); a 27-member combination of the above two models and the AP1 (referred to as JT3); and a 36-member combination of all four models (referred to as JT4, and available for winter simulations only).

Simulated monthly and seasonal averages are derived from daily model values, valid at 1200 UTC. In order to correct (*a posteriori*) for model bias, simulated anomalies are calculated by subtracting the model climate (defined over all 9 ensemble runs and all 15 years) from the individual ensemble fields, while the observed anomalies are derived from the ERA 15-year climatology.

(b) *AGCM hindcasts using persisted SST anomalies*

Hindcast experiments have been performed for twelve of the PROVOST DJF (December–January–February) and MAM (March–April–May) periods (1982–93) for the UM only, using persistence-based forecasts of SST (rather than observed SST) to force the model lower boundary. The persistence-based SST forecasts are produced by adding one-month SSTA from November (DJF hindcasts) and February (MAM hindcasts) to the GISST or Reynolds OI climatological SST fields. In the open ocean SSTA is kept constant until month four, and hence is effectively unchanged through the 4-month integrations considered in this paper. There is ice where the climatological SST plus SSTA falls below  $-1.8^{\circ}\text{C}$ ; thus positive (negative) initial SSTA can delay (hasten) ice formation.

(c) *Verification*

The variables assessed are 850 hPa temperature and precipitation, with emphasis on the first three months of the simulations. Potential probabilistic skill is evaluated for the events: 3-month mean 850 hPa temperature above/below normal; 3-month mean total precipitation above/below normal. Temperature and rainfall are selected for evaluation because of their interest to a wide range of users. Note that 850 hPa temperature anomalies may generally be considered to be a proxy for surface temperature anomalies. It is recognised that to be of benefit to many applications, skill at predicting higher-threshold events on these variables (rather than just the sign of the anomaly) will need to be proven. However, evaluation for ‘above/below normal’ events is considered a necessary first step, and skill at this level may be of practical use to some users. Deterministic (ensemble mean) skill is also evaluated for the same variables. Four main assessment areas are employed: the tropics,  $30^{\circ}\text{N}$  to  $30^{\circ}\text{S}$ ; the northern extratropics,  $20^{\circ}\text{N}$  to  $80^{\circ}\text{N}$ ; North America,  $130^{\circ}\text{W}$  to  $60^{\circ}\text{W}$ ,  $30^{\circ}\text{N}$  to  $70^{\circ}\text{N}$ ; and Europe,  $12.5^{\circ}\text{W}$  to  $42.5^{\circ}\text{E}$ ,  $35^{\circ}\text{N}$  to  $75^{\circ}\text{N}$ . Output from the T63, AP1 and AP2 simulations, archived at ECMWF, were interpolated onto the UM model grid prior to analysis.

The PROVOST simulations and the hindcasts are verified using the ERA dataset. ERA precipitation is based on the accumulation over a 24-hour forecast run from the ERA analyses. The use of model-based precipitation analyses is not ideal, for obvious reasons, but should be sufficient to provide large-scale estimates of potential skill for ‘above/below normal’ events.

### 3. ASSESSMENT OF POTENTIAL SKILL

(a) *Skill assessments for 850 hPa temperature*

In this section we present verifications of the simulated anomalies of average 850 hPa temperature over the first three months of the PROVOST integrations. Probabilistic skill is assessed using the relative operating characteristics (ROC; Stanski *et al.* 1989). Deterministic (ensemble-mean) skill is evaluated using temporal correlation and spatial anomaly correlation scores. For brevity discussion is focused on results from the UM simulations, with comparisons between the UM, T63, and multiple-models (JT2, JT3 and JT4) provided in summary form.

(i) *Probabilistic skill.* The ROC for a specific event is expressed in the form of a curve plotting hit rates against false-alarm rates for a specific event over a range of forecast probability thresholds. The probability thresholds considered here are nominally 20%, 40%, 60% and 80%. In practice the thresholds are defined according to the numbers of ensemble members which predict the event; definitions for the 9-member ensembles

TABLE 1(a). DEFINITION OF NOMINAL PROBABILITY THRESHOLDS, AS USED IN THE ROC ANALYSIS FOR 9-MEMBER ENSEMBLES

No. of members indicating the event	Corresponding nominal probability thresholds (%)
0 or more	0
2 or more	20
4 or more	40
6 or more	60
8 or more	80

TABLE 1(b). DEFINITION OF ROC HIT RATE AND FALSE-ALARM RATE

		Does ensemble probability for the event exceed threshold X?		
		YES	NO	
Is the event observed?	YES	Hit (H)	Miss (M)	H+M
	NO	False alarm (FA)	Correct rejection (CR)	FA+CR
		H+FA	M+CR	

Hit rate for probability threshold X is given by:  $HR = H/(H + M)$ .

False-alarm rate for probability threshold X is given by:  $FAR = FA/(FA + CR)$ .

This contingency table shows definitions of hit rate and false-alarm rate, for a given forecast probability threshold (X) for a binary event, used in the construction of the ROC curves. H, M, FA and CR are total numbers of hits, misses, false alarms and correct rejections at threshold X.

are given in Table 1(a), with an analogous procedure applied for the multiple-model ensembles. Note that the hit and false-alarm rates, for each probability threshold, are defined as proportions of the observed frequencies of the event and non-event, respectively (Table 1(b)). ROC evaluations of the PROVOST simulations for the four assessment regions have been constructed by calculating hit and false-alarm rates over the spatial/temporal domain represented by all grid points in the region and all 15 PROVOST years.

For UM MAM simulations of the event 850 hPa temperature below normal, hit rates exceed false-alarm rates for all threshold probabilities (20%, 40%, 60% and 80%) of the event in all four assessment regions (Fig. 1(a)–(d)), indicating that the ensemble has skill in detecting the event both in tropical and extratropical regions. Skill is greatest in the tropics (Fig. 1(a)) where the hit rate/false-alarm rate ratios are largest. Skill over Europe (Fig. 1(c)) is generally comparable to that over the northern extratropics as a whole (Fig. 1(b)) and is somewhat greater than found for North America (Fig. 1(d)), though the difference may not be significant.

Information on performance over a range of probability thresholds, as expressed in the ROC curve, can be extremely useful to the user of the forecast, particularly in the estimation of the (financial) value of forecasts—a concept explored in section 6. However, for the purposes of a general skill assessment the area under the ROC curve (referred to here as the ROC score) is often used as a summary statistic. The greater the skill of the ensemble, the more the ROC curve must bow up towards the top-left corner, and the greater the area under the curve. A ROC score of 0.5, the area below

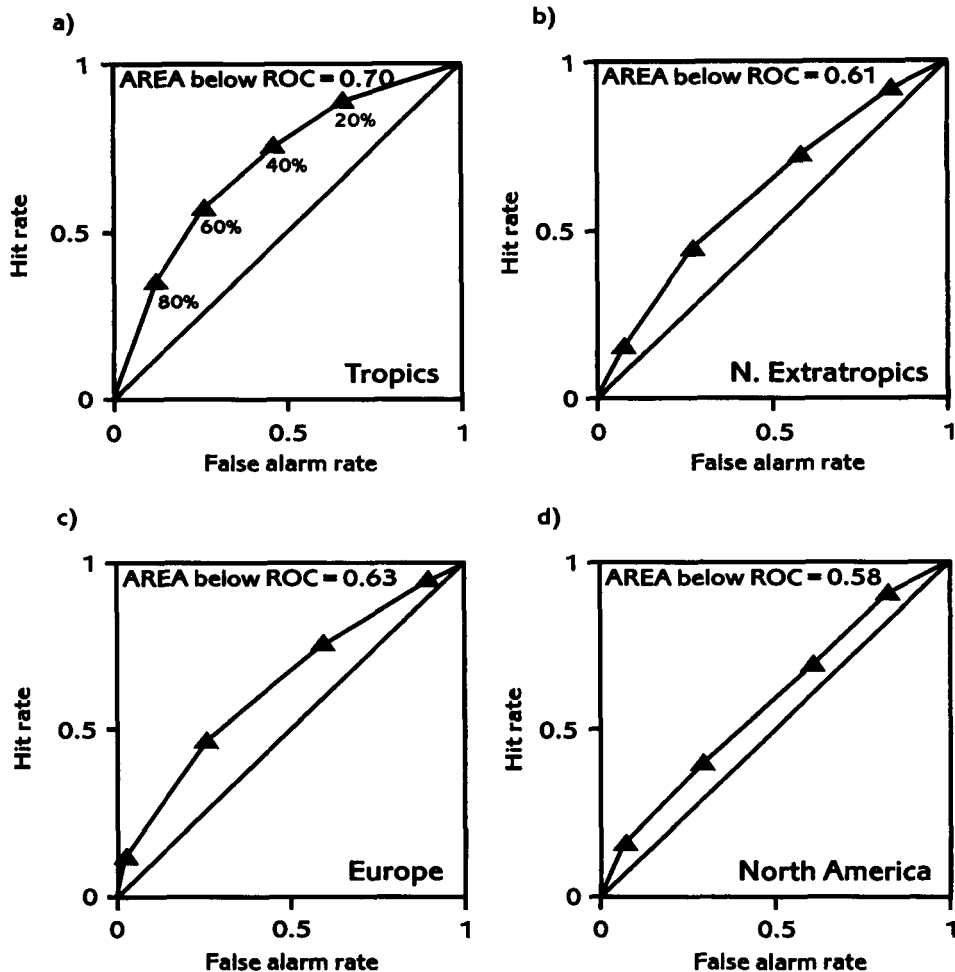


Figure 1. Relative operating characteristic (ROC) curves for The Met. Office Unified Model March–April–May (MAM) simulations of the event 850 hPa temperature below normal over: (a) the tropics, (b) the northern extratropics, (c) Europe, (d) North America. The curves are constructed from hit rates and false-alarm rates (see Table 1(b) for definitions) at four thresholds on the forecast probability of the event (20%, 40%, 60% and 80%, as indicated on panel (a)). The curve is bounded by the points (0,0) and (1,1) which correspond, respectively, to the false-alarm and hit rates achieved through never and always forecasting the event. The areas under the ROC curves (the ROC score) are included and give an overall measure of skill.

the diagonal when hit rates equal false-alarm rates, or less indicates no skill (i.e. skill is no better than that available from a climatology or a random forecast); while a score of unity indicates perfect deterministic skill (i.e. all members correctly predict the event in all years, represented by a single point at (0,1) on the ROC diagram). Note that ROC scores achieved for simulations of a ‘below normal’ event are identical to those obtained for ‘above normal’ events, because the events are complimentary.

ROC scores for 850 hPa temperature obtained with the UM, T63 and JT2 ensembles are compared for all seasons in Fig. 2. For all three ensembles the ROC score exceeds the 0.5 threshold for skill in all four regions and for all seasons, except for the T63 SON (September–October–November) simulations over Europe. Skill is highest in the tropics (Fig. 2(a)), where there is little seasonal variation and ROC scores with all

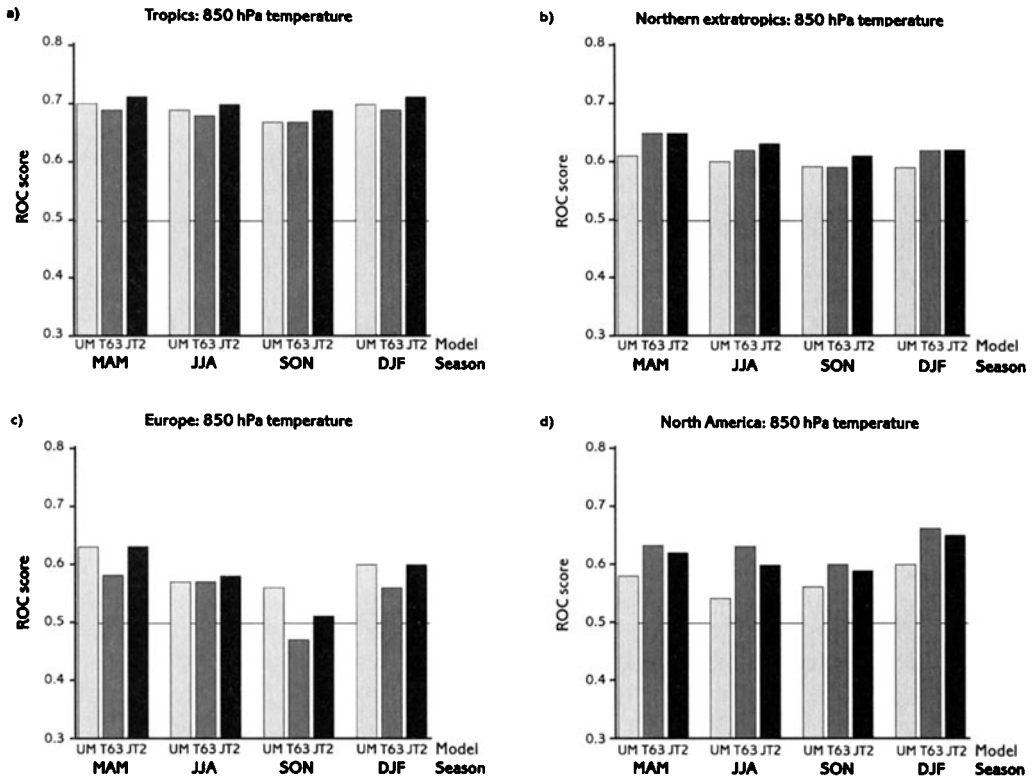


Figure 2. Relative operating characteristic (ROC) scores (areas under the ROC curves) for the event 850 hPa temperature below normal. ROC scores for the event 850 hPa temperature above normal are identical. Results are for: UM, The Met. Office Unified Model ensemble; T63, the ECMWF T63 ensemble; and JT2, a combination of the UM and T63 ensembles. Seasons are: March–April–May (MAM), June–July–August (JJA), September–October–November (SON), and December–January–February (DJF). Areas given are: (a) the tropics, (b) the northern extratropics, (c) Europe, (d) North America.

three ensembles are of order 0.7. In the northern extratropics (Fig. 2(b)) scores for all seasons/models are lower, of order 0.6. There is evidence, most notably in the T63 and JT2 simulations, of a MAM maximum in skill when scores with these models are of order 0.65; Branković *et al.* (1994) and Rowell (1998) have also found maximum skill in the northern extratropics in spring. The SON season appears to be a minimum in forecast skill. Note that of the three ensemble configurations JT2 achieves the best overall performance, achieving in each season and both regions a ROC score equal to or better than that of the best individual model (the UM appears better in the tropics, while the T63 has the best performance in the northern extratropics).

Seasonal and model differences in ROC score are more pronounced in the regional areas of Europe and North America (Fig. 2(c) and (d)). Over Europe the ROC score for all three ensembles is at a maximum in MAM and a minimum in SON—as found for the northern extratropics. In contrast, peak scores for North America in all three ensembles occur in DJF. Enhanced skill over North America in winter may reflect higher predictability in winters with PC/W events (see section 4).

To provide an estimate of the significance of the ROC scores, the calculations were repeated 500 times, each time scrambling the yearly order of the simulations and noting the frequency with which scores exceeded threshold values (in steps of 0.05) by chance.

Using this Monte Carlo method it was found that, for this (15-year) sample, the threshold for 90% significance lies between 0.60 and 0.65, and the threshold for 95% significance between 0.65 and 0.70. Thus scores for the tropics are significant at the 95% level in all seasons, while for the extratropics significance is evident in MAM but still is marginal in other seasons. Over Europe and North America scores in the seasons of peak skill (MAM and DJF respectively) are significant with 2 of the 3 ensembles. T63 scores are also significant over North America in MAM. In JJA (June–July–August) and SON scores are not significant over Europe and generally of marginal significance over North America. It is important to recall that the significance estimates discussed above are based on scores calculated over all 15 years. In section 4 we show that skill is notably enhanced in some seasons during PC/W years. Moreover, the assessment regions employed are very broad, and there is evidence (discussed later in this section) to suggest that skill levels in certain sectors of these regions will be substantially higher than indicated by the region-wide scores.

Comparisons of individual-model performance indicate that the UM performs better over Europe, while the T63 is more skilful over North America. The best model in each region achieves ROC scores comparable to those obtained for the northern extratropics, indicating potential gains in capability from the use of more than one AGCM in an operational environment. Note that the JT2 ensemble appears to act as a filter for the more skilful individual model, and thus provides a means of exploiting the strengths of each model. Even where the difference in the ROC score between models is relatively large (see e.g. MAM over Europe and JJA over North America (Fig. 2(c) and (d))), the value achieved by the JT2 system is similar to that obtained by the more skilful model. The main exception is the autumn season over Europe, where the T63 ROC score is below the 0.5 threshold for skill.

A corresponding analysis, including all four individual PROVOST models and the multiple-model combinations JT2, JT3 and JT4, is provided in Fig. 3 for simulations over Europe and North America. When individual-model skill is at similar levels the multiple-models provide improved skill (e.g. JJA simulations over Europe, Fig. 3(a)), which may derive both from the presence of additional models and from the increased ensemble size. However, the most striking benefit provided by the multiple-model ensembles is the skill-filtering property in regions/seasons when skill of the individual models varies widely; in this respect the benefits noted for JT2 are usefully extended by the JT3 and JT4 ensembles. Note, for example, the DJF simulations over North America (Fig. 3(b)) for which the skill of each multiple-model (JT2, JT3 and JT4) is consistently close to the skill of its most skilful individual component model: T63 for JT2, AP1 for JT3 and AP2 for JT4—resulting in the best overall skill (out of the multiple-models) for JT4. Similar results are found for the DJF simulations over Europe (Fig. 3(a)), where the higher skill of the UM and AP1 models is matched by the JT4 ensemble, despite lower skill from the T63 and AP2 models. The fact that JT4, a 36-member ensemble, provides similar skill to the best 9-member individual ensemble indicates that the increased ensemble size plays only a small role in producing the relatively high skill of JT4; the key element in determining the skill of the multiple-model appears to be the skill of the most skilful component ensemble.

The two examples discussed above show that the relative skill of individual models may differ markedly with season and region, and that a simple, unweighted, combination of ensembles from different models can usefully exploit the strengths of the individual models. This is an encouraging result, as it indicates that benefits are immediately available from a multiple-model approach. Note that benefits may be further improved by weighting the contribution of individual models according to performance over a



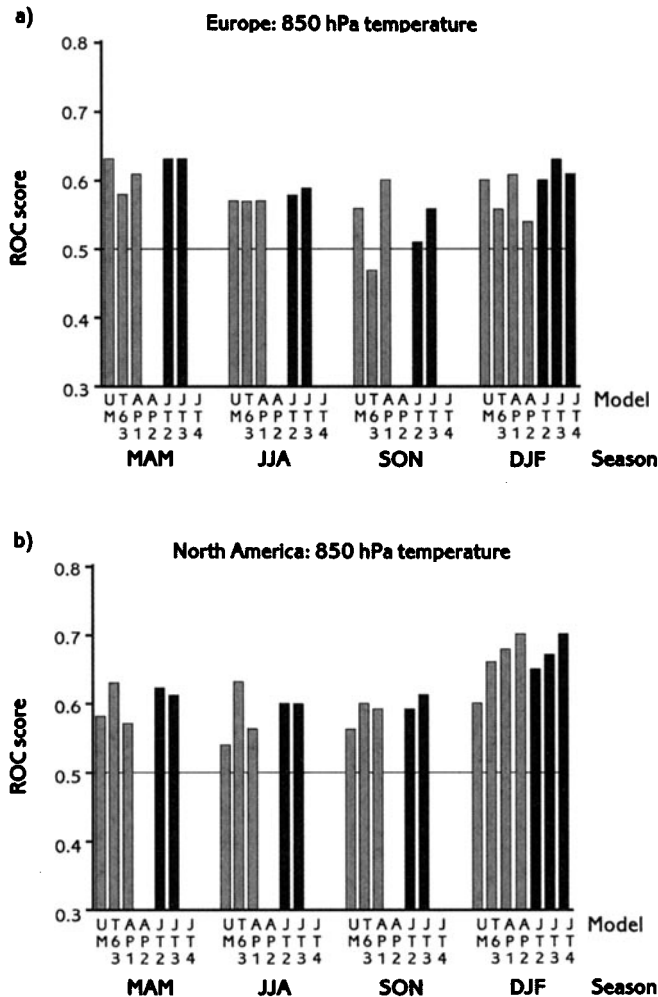


Figure 3. As Fig. 2 but for (a) Europe and (b) North America, and all four participating PROVOST models (hatched bars) and the multiple-model configurations JT2, JT3 and JT4 (solid bars). Scores are given for March–April–May (MAM), June–July–August (JJA), September–October–November (SON) and December–January–February (DJF). Model nomenclature is: UM, The Met. Office Unified Model; T63, ECMWF T63; AP1, Météo-France ARPÈGE T42 L31; AP2, the ARPÈGE T63 L31 (run at Electricité de France for DJF only); JT2, UM+T63 (18 members); JT3, UM+T63+AP1 (27 members); and JT4, UM+T63+AP1+AP2 (36 members, for DJF only).

hindcast period (see e.g. Krishnamurti *et al.* 1999). However, further assessment of relative weighting is required to find optimum strategies, particularly with regard to their use in ensemble-based probability forecasting.

The spatial distribution of ROC score has been calculated by obtaining ROC curves for each model grid point from the 15 available simulations. The distribution obtained for UM simulations of the MAM period is provided in Fig. 4 and gives further insight into the geographical variation of skill, though, because of the smaller sample size (15) individual values should be treated with caution. Geographical variations are broadly similar in the other seasons and are not shown. ROC scores exceed the 0.5 threshold for skill over many parts of the globe (shading starts at 0.55 in Fig. 4). Consistent with the regional analysis, high ROC scores are most widespread in the tropics where the

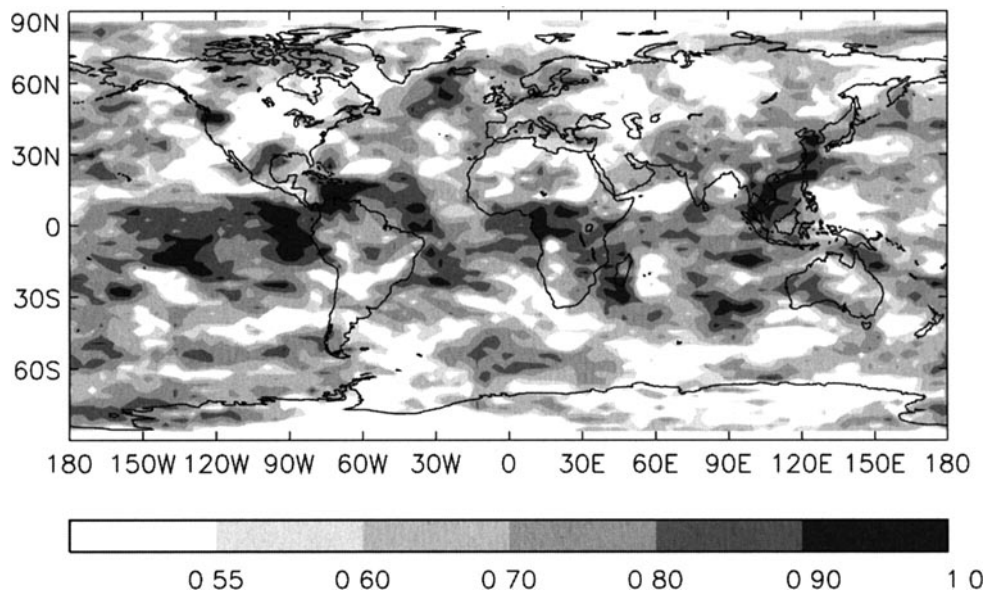


Figure 4. Geographical distribution of the relative operating characteristic (ROC) score (the area below ROC curves obtained at each model grid-point) for The Met. Office Unified Model simulations of the event March–April–May 850 hPa temperature below (or above) normal. Shading thresholds are 0.55, 0.6, and then at intervals of 0.1.

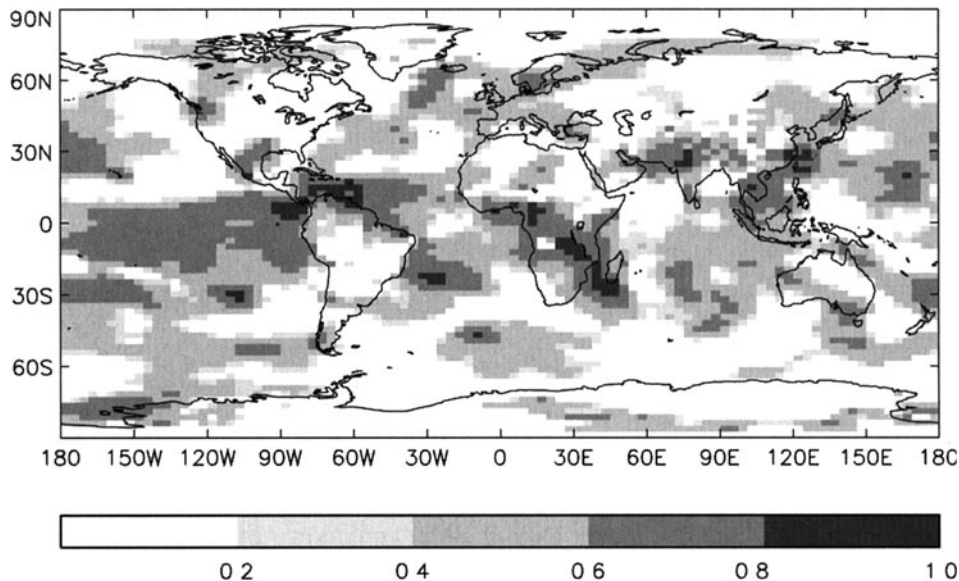


Figure 5. Spatial map of correlation coefficients between The Met. Office Unified Model ensemble-mean and observed 850 hPa temperature at each grid point over the 15 PROVOST March–April–May periods. The contour interval is 0.2; values are shown only where correlations are positive and significant at the 90% level or higher.

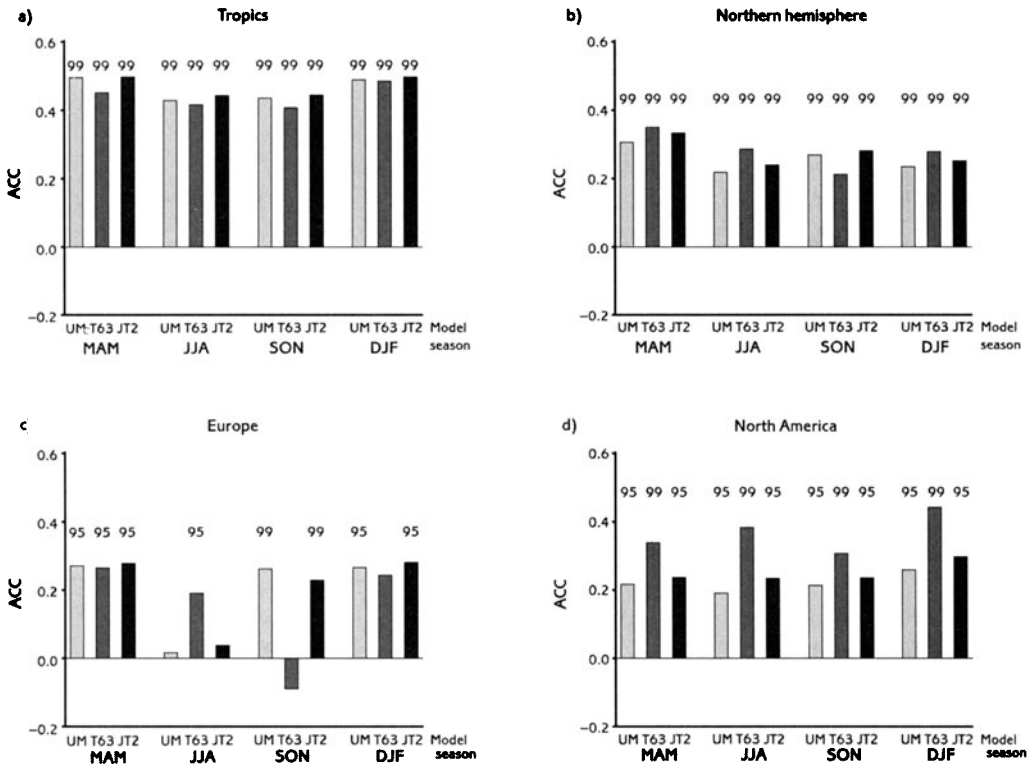


Figure 6. Seasonal average anomaly correlations of ensemble-mean and observed 850 hPa temperature for The Met. Office Unified Model (UM), ECMWF T63 L31 model (T63) and JT2 (UM+T63, giving 18 members) ensembles. The seasons are: spring (March–April–May, MAM), summer (June–July–August, JJA), autumn (September–October–November, SON) and winter (December–January–February, DJF). Averages are over the 15-year period December 1979 to March 1994 (14 years for T63 and JT2 DJF simulations), and are calculated using Fisher's z-transform (see text). When significance with which the average is different from zero exceeds a threshold of 95% or 99%, the threshold value is plotted above the bars. Results are for: (a) the tropics, (b) the northern extratropics, (c) Europe, (d) North America.

ensemble simulations frequently reduce to a single deterministic solution (e.g. consistent indication of the event in all nine ensemble members), and maximum deterministic skill over the 15-year period is approached or achieved in many areas. Marked regional variations in ROC score are evident both in the tropics and extratropics, with the regional assessment areas (Europe and North America) encompassing regions of widely different skill. Over Europe highest scores are generally found in northern and western regions. Over North America the higher scores are located over southern, western and northern continental fringes. In both the tropics and extratropics ROC scores are generally highest over the oceans and lowest over continental interiors.

(ii) *Temporal and spatial correlation of the ensemble mean.* Here we examine potential deterministic skill of the ensemble-mean 850 hPa temperature anomalies, and compare results with those described above for probabilistic forecasts of 2-category events. The distribution of point correlations of UM ensemble-mean and observed 850 hPa temperature anomalies is given for the MAM season in Fig. 5. The significance of the correlations has been estimated using a Monte Carlo technique to estimate the probability of achieving equivalent correlations by chance (500 correlations were calculated, each after randomly scrambling the yearly order of the ensemble-mean values). Correlations

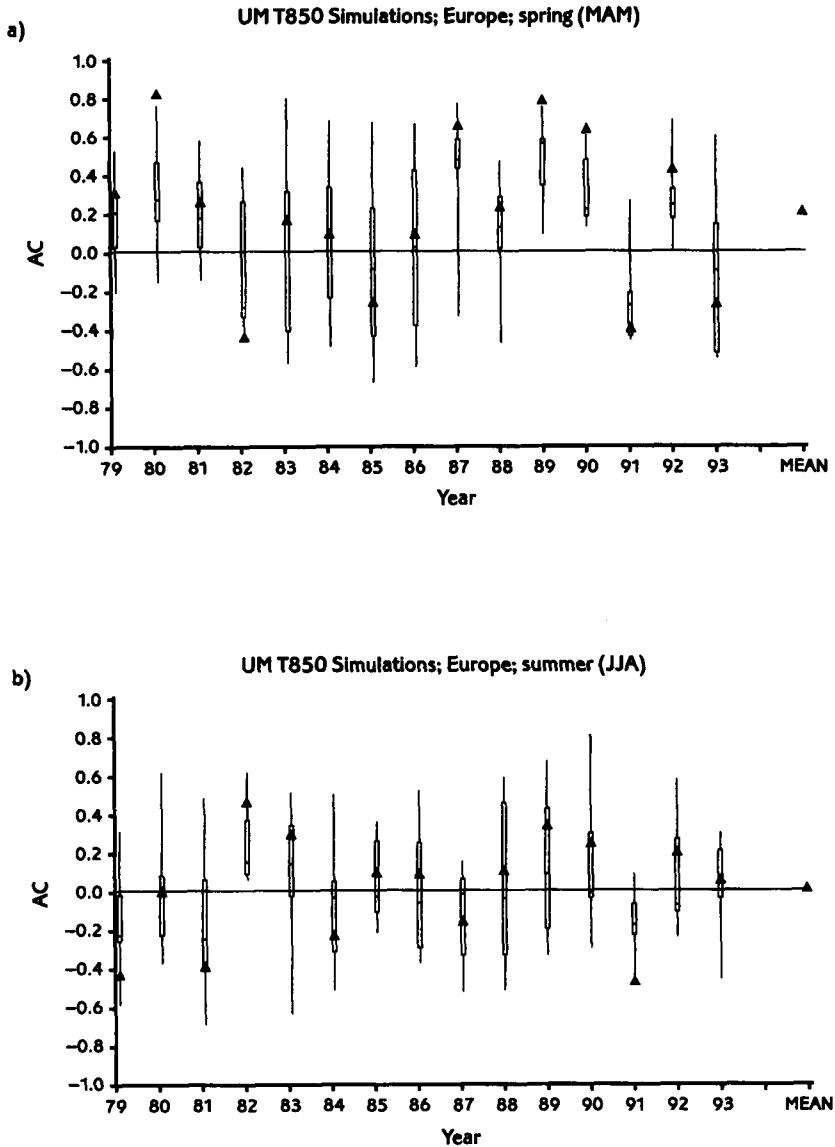


Figure 7. Annual anomaly correlation coefficients (ACs) for the nine Met. Office Unified Model ensemble members for simulations of 850 hPa temperature over Europe (1979–93). Each stem and whisker plot indicates median, quartile and extreme values. Ensemble-mean values are shown as solid triangles, with the 15-year average shown at the extreme right. (a) Spring, March–April–May; (b) summer, June–July–August; (c) autumn, September–October–November; (d) winter, December–January–February (DJF). Note that for DJF, the year refers to the December of the period, i.e. 79 is DJF79/80.

are plotted only when they are significant at the 90% level or higher. As found for the ROC score assessments, skill is best in the tropics where the correlation coefficient (CC) frequently exceeds a value of 0.6, with local peaks in excess of 0.8. Significant correlations are also present in the extratropics, though with lower CC values of order 0.4 with peaks to 0.6. In most regions significant correlations of the ensemble mean

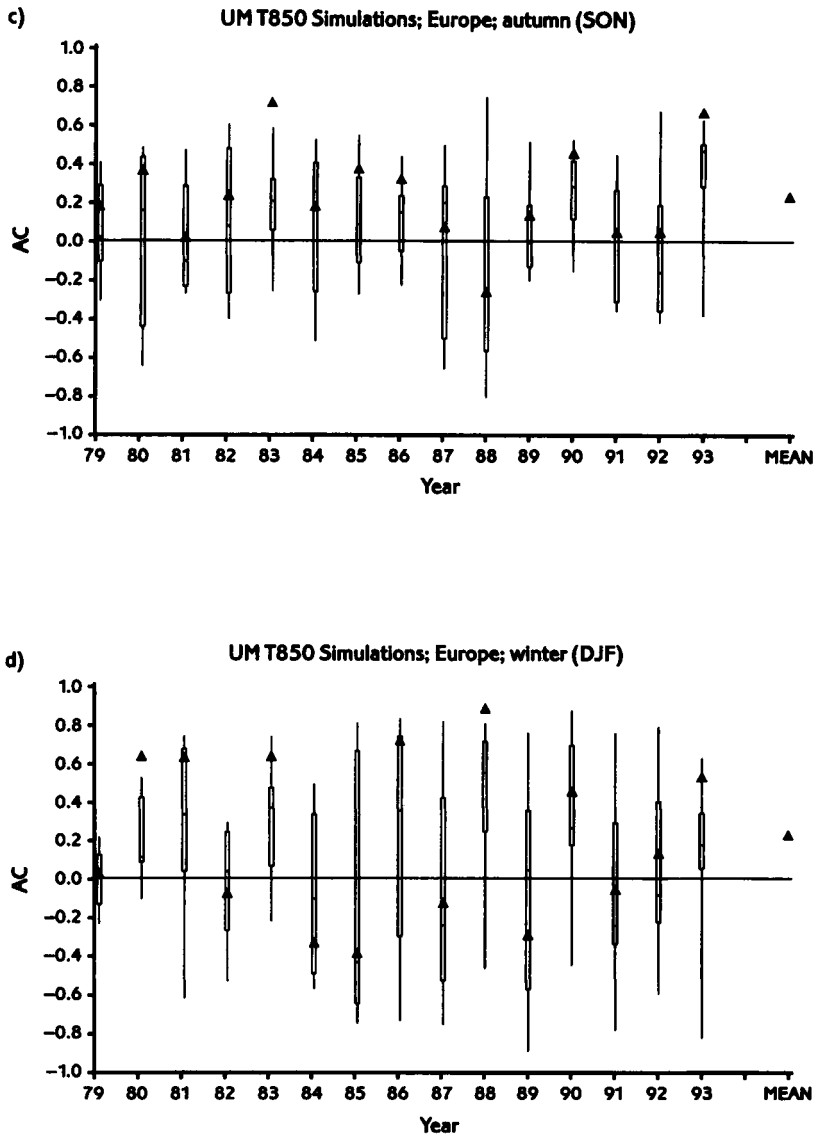


Figure 7. Continued.

correspond well with regions of high ROC score both in MAM (c.f. Fig. 4) and in other seasons (not shown).

Average spatial anomaly correlation coefficients (ACC scores) of ensemble mean values for the four assessment regions are provided in Fig. 6(a)–(d) for the UM, T63 and JT2 month 1–3 simulations; averages are over the 15 PROVOST years (14 for T63 and JT2 in DJF). (Discussion of ACC scores for months 2–4 is provided in section 4(a).) As the ACC is not a normally distributed quantity, averages are calculated (following e.g. Buizza 1997) by applying Fisher's  $z$ -transform to the individual yearly values, averaging, then applying the reverse transform. ACC scores are positive in all seasons and all regions, except for T63 SON simulations over Europe. ACC scores differ from zero with high levels of significance (using a  $t$ -test) in the tropics and

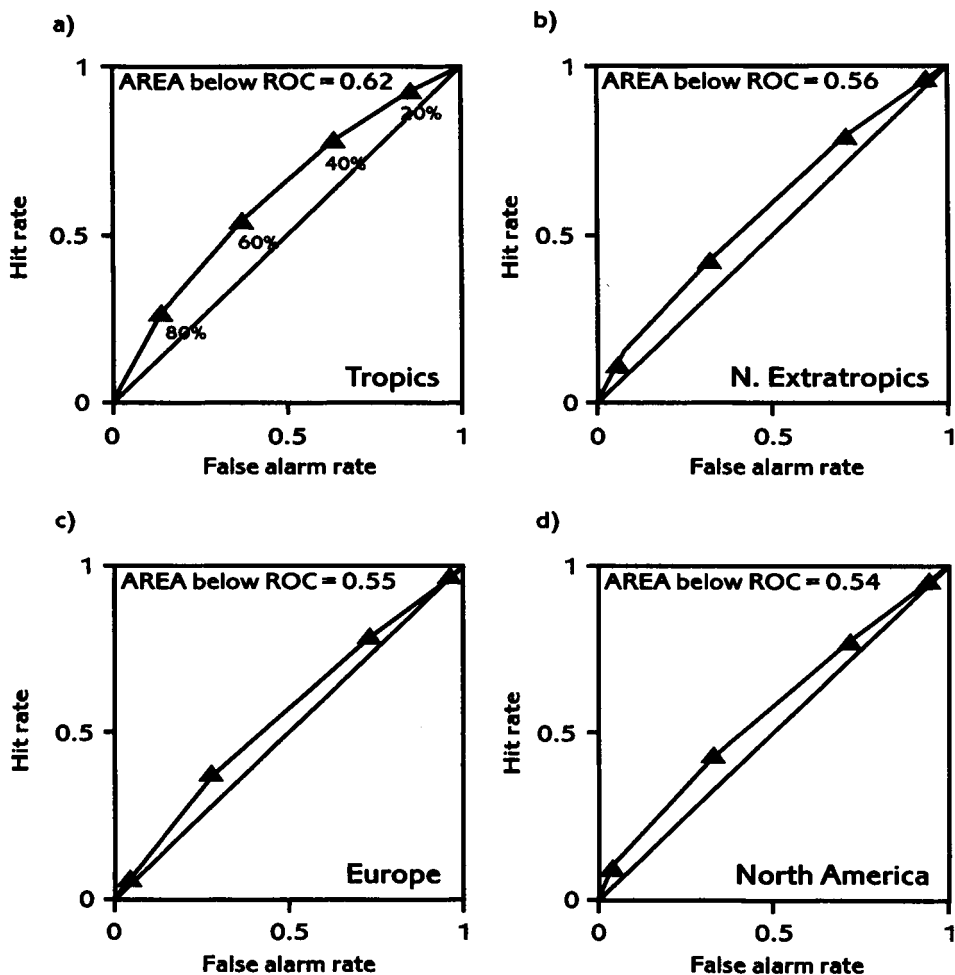


Figure 8. As Fig. 1 but for precipitation.

northern extratropics in all seasons and with all three ensembles. Over North America significant ACC scores are obtained in all seasons with all three ensembles, with the best performance from T63. Over Europe significant ACC scores are obtained with at least one model in all seasons indicating, as found for the ROC scores, the potential benefits of using two or more prediction models in an operational environment.

Seasonal and model differences in mean ACC scores show broad similarities with the results of the ROC analysis. In particular skill scores are highest in the tropics (Fig. 6(a)—average ACC score of order 0.4–0.5), where seasonal and model differences are small; in the northern extratropics (Fig. 6(b)) all three ensembles indicate best skill in MAM (ACC scores are of order 0.2–0.3); over Europe skill is most consistent over the three ensembles in MAM and DJF, with scores similar in both periods and comparable to the northern hemispheric values; over North America all three ensembles indicate best skill in DJF. Significance levels for the ACC scores are generally higher, but broadly consistent with those found for ROC scores using the Monte Carlo method (section 3(a)(i)). Benefits from the JT2 ensemble are evident in most seasons over the

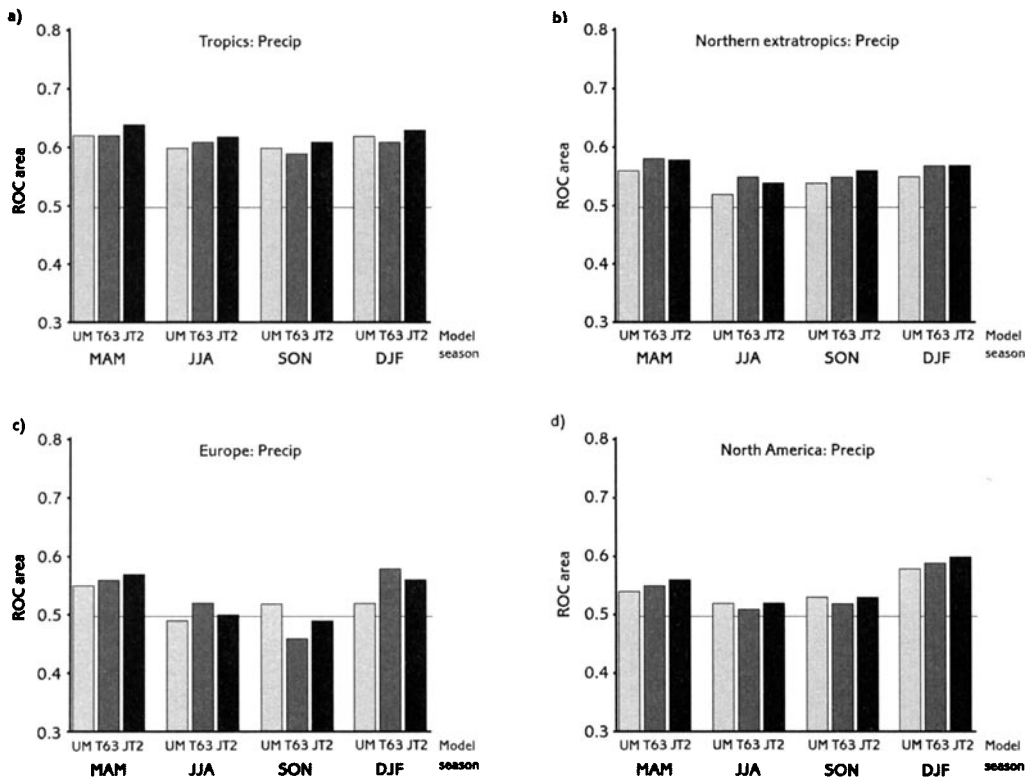


Figure 9. As Fig. 2 but for precipitation.

tropics and Europe, but (unlike from the ROC score assessments) are not evident over the North American region where the better individual model performance obtained with the T63 ensemble is not matched by JT2.

Over the northern extratropics UM ACC scores for individual seasons are positive in most years (not shown); of the 15 years for each season, negative scores are present in only one MAM, two JJA, and three DJF seasons (correlations in SON are positive in all years). For the regional areas the interannual variability of skill is larger, particularly over Europe (Fig. 7) where negative scores are present in four MAM, five JJA, one SON and six DJF seasons. Note, however, that even in years when the ensemble mean is negatively correlated, members with relatively large positive correlations are often present (e.g. DJF 1985/86 and 1989/90, Fig. 7(d)), indicating the greater potential of a probabilistic approach to seasonal forecasting.

#### (b) Skill assessments for precipitation

For reasons of brevity, we restrict diagnosis of skill for precipitation to results obtained with the ROC analysis. ROC curves for UM simulations of the event 3-month mean MAM total precipitation below normal, and ROC scores (for both 'above' and 'below' events) for all three ensembles over all seasons are provided in Figs. 8 and 9 respectively. In all regions and seasons skill is lower than for 850 hPa temperature (c.f. Figs. 1 and 2). A lower level of predictability for precipitation is to be expected, since its production is sensitive to a greater range of 'chaotic' processes than acts on the 850 hPa temperature field. As for 850 hPa temperature, skill is highest in the tropics (Fig. 9(a)),

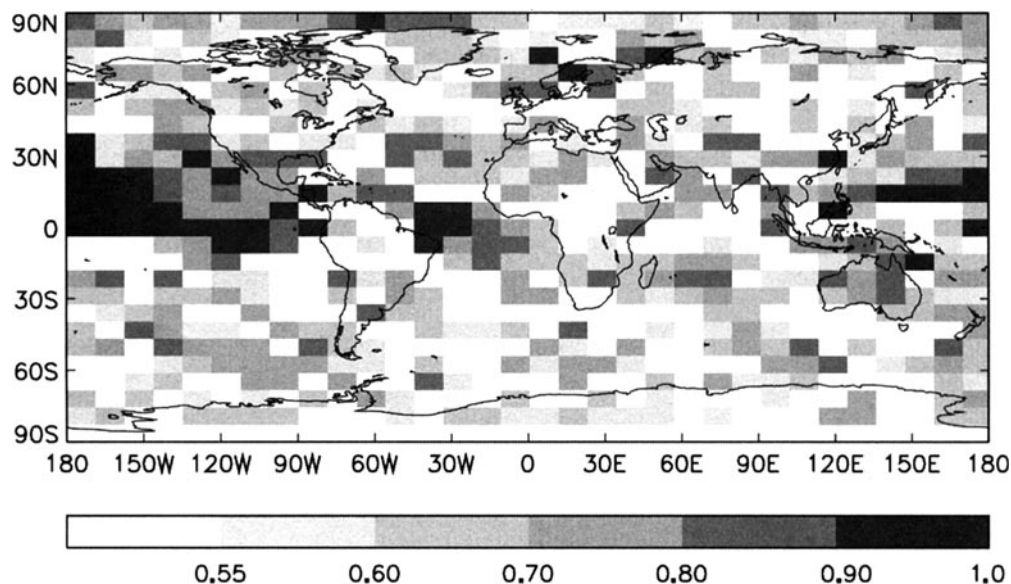


Figure 10. As Fig. 4, but for precipitation. Simulated and observed (ECMWF re-analysis) precipitation is smoothed over  $3 \times 3$  grid-boxes (corresponding to  $7.5^\circ$  latitude by  $11.25^\circ$  longitude) prior to calculation of the ROC scores.

where there is little variation with season and scores with all three ensembles are of order 0.6. In the northern extratropics (Fig. 9(b)) the ROC scores for all seasons/models are lower at order 0.55, with an indication of enhanced skill in DJF and MAM relative to JJA and SON. As noted for the 850 hPa temperature simulations, the JT2 ensemble achieves the best overall performance, obtaining in most seasons in both regions a score similar to or better than that of the best individual model.

Differences in performance between seasons and models are more pronounced over Europe and North America (Fig. 9(c) and (d)). In both regions ROC scores are highest in spring and winter, with values comparable to the northern hemispheric values. Over Europe (Fig. 9(c)) scores in winter and spring are generally similar, while over North America (Fig. 9(d)) DJF appears to be the season of highest skill (as found for 850 hPa temperature). In both regions scores for JJA and SON are below the hemispheric values, notably over Europe where for both seasons scores fail to exceed the 0.5 threshold for skill in two of the three ensembles. On average there is little difference in the performance of the UM and T63 ensembles. Note, however, that with the exception of the relatively low-skill seasons (JJA and SON) over Europe, the ROC score achieved by the JT2 ensemble exceeds or is similar to that obtained by the most skilful individual model. ROC scores are sufficiently high to be deemed significant (at the 90% level, scores between 0.6 and 0.65) only in the tropics. However, as for 850 hPa temperature, skill can vary considerably within the assessment regions and is likely to be significant in some extratropical regions with peak skill (see below).

The geographical variations in ROC score for UM simulations of the MAM season are provided in Fig. 10, and are broadly representative, in most regions, of the variations in other seasons. To reduce 'noise' in the observed and simulated precipitation fields, both have been smoothed over  $3 \times 3$  grid boxes ( $7.5^\circ$  latitude by  $11.25^\circ$  longitude) prior to calculation of the ROC curves. The distribution of skill is broadly similar to that found for 850 hPa temperature (c.f. Fig. 4); however, skill coverage is generally lower,



particularly at higher ROC scores (e.g. 0.8 and above). As for 850 hPa temperature, highest scores (exceeding 0.9) occur in tropical regions. However, scores greater than 0.5 (shading starts at 0.55 in Fig. 10) are present over many parts of the extratropics with notable peaks in some regions.

In some tropical regions high ROC scores (of order 0.9) correspond with the local wet season, indicating substantial potential for seasonal rainfall prediction. Note, for example, the relatively high ROC score over north-east Brazil in MAM. Relatively high ROC score values are also present (not shown) over parts of equatorial Africa (notably the Guinea coast) and the Indian sub-continent in JJA—implying probabilistic skill in simulating average monsoon rainfall in these regions. Evidence of potential model skill in these areas has led to production of experimental real-time predictions (see e.g. Harrison *et al.* 1997a,b and Evans *et al.* 1998).

Over North America highest scores are generally found over southern and western regions, with the notable exception of DJF (not shown) when relatively high scores are also present over central and north-western regions. Over Europe highest scores are found in north-western regions and over parts of the Mediterranean region. The caveat that these distribution patterns are based on a small sample should be emphasised; however, they serve to indicate where regional enhancements in skill may be expected.

#### 4. PROSPECTS FOR SKILL PREDICTION

Although 15-year average ensemble-mean ACC scores (Fig. 6, 850 hPa temperature) are below values usually considered useful for medium-range numerical weather prediction (a threshold of 0.6 is frequently quoted for instantaneous fields), particularly skilful years are present in the time series (see e.g. Fig. 7 for Europe). However, the presence of such skilful years is only of value if the occurrence of relatively high skill can be predicted. In this section we examine the prospects for skill prediction by investigating the degree to which the state of ENSO and internal ensemble spread may be used as predictors of deterministic skill over North America and Europe. Analysis is restricted to 850 hPa temperature.

##### (a) *The impact of PC/W events on predictability for SON, DJF and MAM*

ACC scores for UM simulations over Europe and North America have been calculated separately for the five DJF periods with PC/W events: 1982/3, 1986/7, 1991/2 (PW, El Niño) and 1984/5, 1988/9 (PC, La Niña); and for the five SON/MAM periods preceding/following these DJF periods when substantial SSTAs are also present (see Fig. 1 of Branković and Palmer 2000). Note, however, that SSTAs are also substantial in two SON periods (1985 and 1987) not considered here as pre El Niño or La Niña. For both regions scores for SON 1985 were above average, while scores for 1987 were below average (see Fig. 7 for Europe). Thus classification of these seasons as non-PC/W makes little difference to the average score.

Average ACC scores for UM simulations of DJF periods with and without PC/W events, and for the preceding/following SON/MAM periods are provided for Europe and North America in Fig. 11(a) and (b) for month 1–3 simulations. Corresponding scores for month 2–4 simulations are provided in Fig. 12(a) and (b). Enhancement of UM month 1–3 ACC scores during PC/W years peaks in MAM over Europe and in DJF over North America (Fig. 11(a) and (b)). Over Europe (Fig. 11(a)) the ACC scores in PC/W years increase throughout the SON–DJF–MAM period (reaching 0.4 in MAM), while the score for non-PC/W years remains relatively constant (at order 0.2). Note that in MAM the average score in PC/W years is significantly different from zero (at

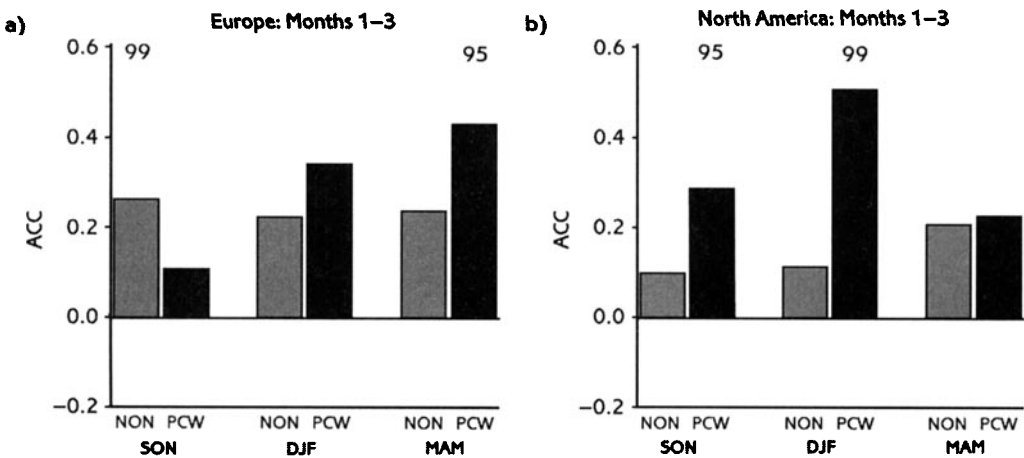


Figure 11. Average anomaly correlation coefficients (ACC) of The Met. Office Unified Model (UM) months 1–3 simulations of 850 hPa temperature for PC/W (bars labelled ‘PCW’) and non-PC/W (bars labelled ‘NON’) September–October–November (SON), December–January–February (DJF) and March–April–May (MAM) periods (see text for definition of PC/W and non-PC/W seasons). When significance (using a *t*-test) with which the average is different from zero exceeds a threshold of 95% or 99%, the threshold value is plotted above the bars. Results are for: (a) Europe, (b) North America.

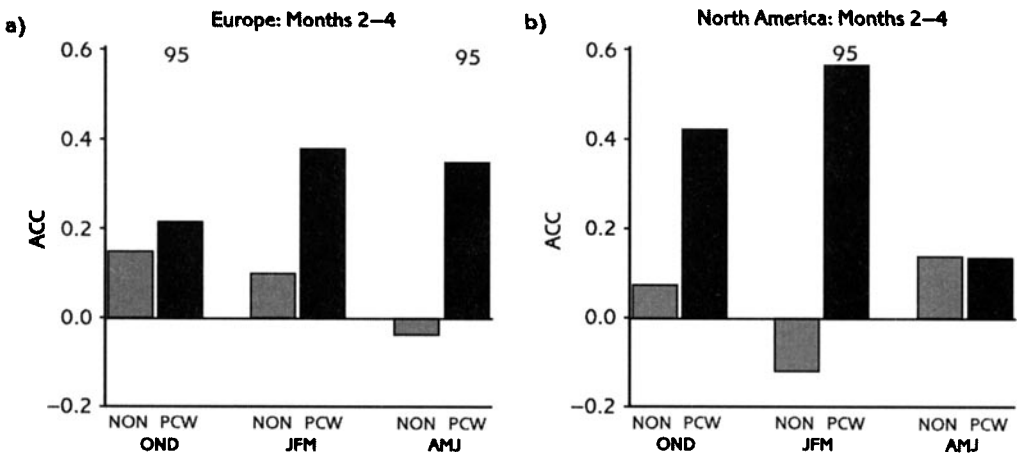


Figure 12. As Fig. 11, but for UM months 2–4 simulations. Periods are now: October–November–December (OND), January–February–March (JFM) and April–May–June (AMJ).

the 95% level), while the average for non-PC/W years is not significant. The only other ACC score (over both PC/W and non-PC/W years) significantly different from zero over Europe occurs for non-PC/W SON seasons.

Over North America (Fig. 11(b)) the ACC score in PC/W seasons increases from order 0.3 in SON to a maximum (of order 0.5) in DJF; in both periods the score is significantly greater than zero (at the 95% and 99% levels, respectively). In contrast the ACC score for non-PC/W years is of order 0.1–0.2, and not significant in all three periods. The average anomaly correlation for PC/W seasons is lowest in MAM, when scores are similar to those of non-PC/W years (order 0.2). This is in direct contrast to the European area, where MAM appears to be the period of maximum enhancement from

PC/W forcing (Fig. 11(a)). Note that skill over North America and Europe is similar in non-PC/W DJF periods (ACC scores are in fact slightly better over Europe), but greater over North America in PC/W DJF periods (compare Fig. 11(a) and (b)). The greater predictability often attributed to the North American region relative to Europe would thus appear to apply (at least with the UM) during PC/W DJF events only.

An investigation of the geographical variation of ENSO impact on skill was conducted by assessing the frequency with which the sign of the anomaly (above or below) indicated as most probable by the ensemble members was correct. It was found that skill enhancements in the DJF simulations over North America during PC/W events are focused over north-western regions, suggesting increased predictability of the Pacific North American (PNA) mode. Barnett *et al.* (1997) also find higher predictability of the winter PNA during strong tropical Pacific SST events. The enhanced skill over North America in PC/W periods explains (at least in part) why best overall skill for this region is found in DJF (Figs. 2(d) and 6(d)), in contrast to the northern hemisphere as a whole (and Europe) for which best skill is found in MAM.

ACC scores for months 2–4 of the UM simulations are provided for Europe and North America in Fig. 12(a) and (b) respectively. In non-PC/W years average skill decreases in months 2–4 relative to the months 1–3 period, substantially in some cases (e.g. compare MAM and AMJ over Europe, Figs. 11(a) and 12(a)), suggesting that when SST forcing is weak a substantial part of the month 1–3 skill may derive from the first month of the integrations, when memory of initial conditions will contribute significantly to skill. In contrast the average skill in PC/W years for seasons with enhanced skill in months 1–3 is generally maintained, or even increases into the later period (e.g. compare DJF and JFM over North America, Figs. 11(b) and 12(b)), indicating the dominant role of SST forcing throughout both periods in these seasons.

For the seasons with best skill enhancement in PC/W years, ACC scores for the months 1–3 and months 2–4 periods compare favourably with the 15-year average score for the first month of the integration. For example, over North America the 15-year average ACC score for December is 0.34 compared to a score of order 0.5 for PC/W DJF and JFM periods, while for Europe the average ACC score for March is 0.48 compared to one of order 0.4 for PC/W MAM and AMJ periods. Thus in PC/W years the skill for forecasts one-season ahead, at zero- and one-month lead, appears comparable to the overall average expected skill for forecasts one-month ahead. In contrast, skill for forecasts one-month ahead exceeds that for those one-season ahead in non-PC/W years.

Corresponding results from the T63 and JT2 simulations are similar in most respects to those presented above for the UM, and are therefore not shown. In particular, all three ensembles indicate best skill enhancement in PC/W years over North America in DJF and over Europe in MAM.

#### (b) *Relationship between ensemble spread and ensemble-mean skill*

Prediction of forecast skill may also be approached through the relationship between ensemble spread and the skill of the ensemble mean. For a given case, the degree of spread of the ensemble members about the ensemble mean is a measure of the sensitivity to initial conditions, and thus should allow assessment of the intrinsic predictability; low spread is ideally associated with high ensemble-mean skill (see e.g. Molteni *et al.* 1996). In this section we evaluate ensemble skill/spread relationships for simulations of 850 hPa temperature over Europe and North America. In order to investigate possible links between ensemble spread and PC/W events, emphasis is given to analysis of the seasons for which PC/W events appear to have most impact, i.e. MAM over Europe and DJF over North America.

(i) *Europe*. Scatterplots of ensemble-mean skill (anomaly correlation of the ensemble mean) and ensemble spread are provided in Fig. 13 for UM, T63 and JT2 simulations of 850 hPa temperature over Europe. Ensemble spread is defined here as the average anomaly correlation of the ensemble members with the ensemble mean (the average is again calculated using Fisher's  $z$  transform, see section 3(a)(ii)); note that with this definition high correlation corresponds to low ensemble spread. A measure of ability to distinguish relatively skilful from unskilful predictions is provided by comparing the total entries in the diagonal quadrants with the total in the off-diagonal quadrants (here the quadrants are constructed using the ensemble median values of ACC skill and spread). A positive skill/spread correlation, and thus potential for skilful prediction, is indicated if entries are maximized in the diagonal quadrants. The linear correlation is also provided; however, because of the sensitivity of linear correlation to outlying points, the nonlinear measure described above is preferred. Years with PC/W events are indicated by symbols in Fig. 13 and non-PC/W years by letters.

The UM ensemble spread shows clear potential for distinguishing relatively skilful and unskilful MAM predictions (Fig. 13(a)). A total of 11 simulations fall in the diagonal quadrants representing correct assessments of ensemble-mean relative skill, and two fall in off-diagonal quadrants, representing incorrect assessments. We shall express this measure of the skill-spread correlation as a positive correlation of 11/2 (when diagonal entries are fewer than the off-diagonal entries we refer to negative correlation). Note that the total of the correct and incorrect assessments may be different from the total number of years because simulations with skill or spread equal to the median value are not counted.

For Europe, MAM is the season for which ACC scores are most enhanced, on average, during PC/W events. It is interesting, therefore, to compare ACC scores for individual PC/W and non PC/W years. The four UM simulations with ACC scores of order 0.6 or greater are split equally between PC/W years (1987 and 1989) and non-PC/W years (1980 and 1990). Thus a PC/W event is not a requirement for relatively high skill. Moreover, the presence of a PC/W event does not guarantee above median skill, with below median scores obtained in the PC/W years 1983 and 1985, the score in 1985 being amongst the lowest (of order  $-0.2$ ) achieved in all years.

In contrast to the UM, the T63 MAM simulations (Fig. 13(b)) show a negative skill/spread correlation (4/8). However, the enhancement of ACC scores in PC/W years appears somewhat more consistent than for the UM; four of the five simulations in these years achieve above median skill, compared with three out of five for the UM. Note that the JT2 ensemble (Fig. 13(c)) retains much of the ability for skill prediction exhibited by the UM simulations, achieving a skill/spread correlation of 9/3. There is no strong signal in either the T63 or UM simulations to suggest that spread in the ensemble is lower during PC/W events.

The UM simulations (Fig. 13(a)) provide tentative evidence that, for some seasons, successful skill prediction strategies may be formulated through reference both to the state of ENSO and to ensemble spread. For example, a 'cautious' approach requiring both a PC/W event and low ensemble spread identifies the skilful MAM seasons 1987, 1989 and 1992 and filters out all seven MAM seasons with below median skill—though this strategy would also 'miss' two particularly skilful years (1980 and 1990). However, the above strategy does not appear to be useful with the T63 simulations (Fig. 13(b)), the 1985 MAM season being associated both with a La Niña event and low spread but having a relatively low ensemble-mean skill. For any AGCM it would be necessary to test proposed strategies for skill prediction on a larger sample of cases.

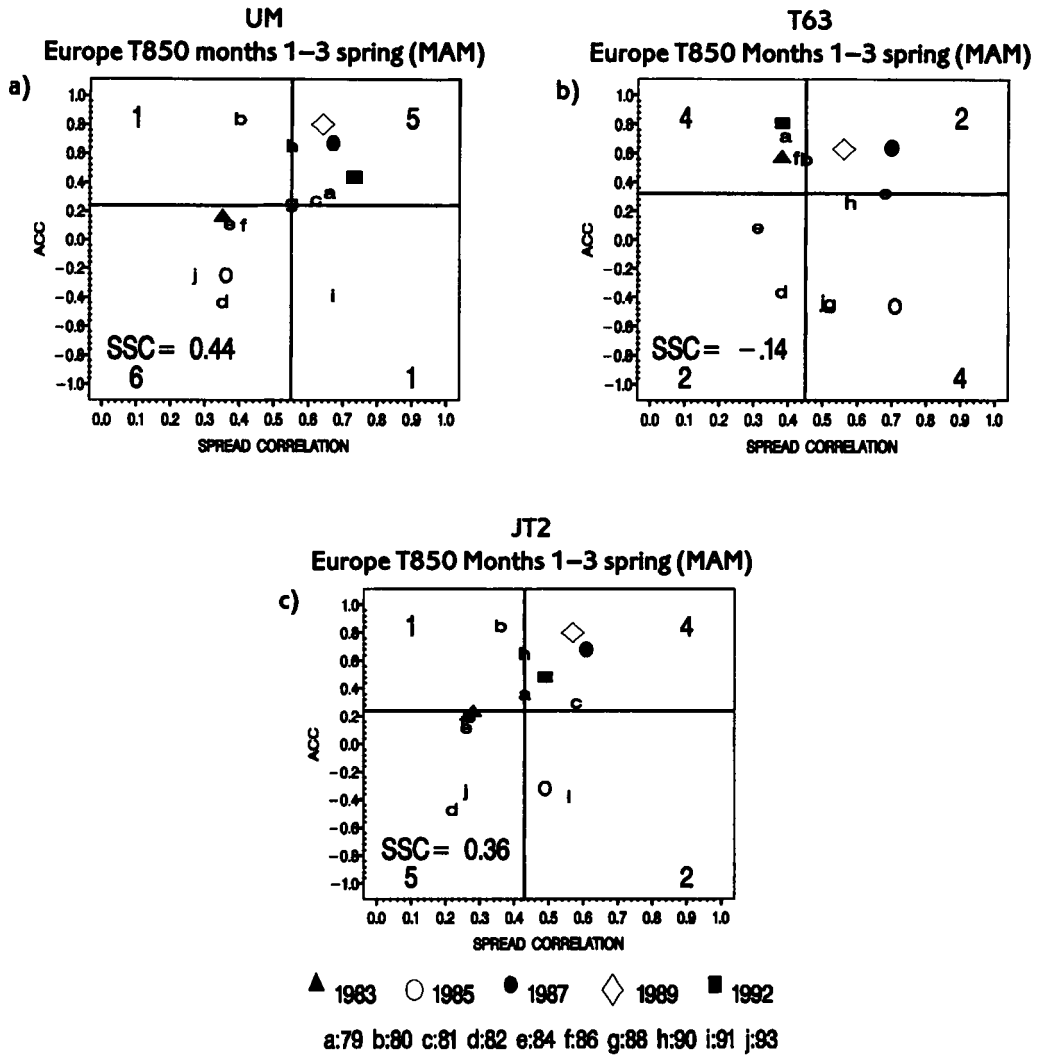


Figure 13. Scatterplots of ensemble-mean skill (defined using anomaly correlation, ACC) and ensemble spread for months 1–3 March–April–May (MAM) simulations of 850 hPa temperature over Europe. Ensemble spread is defined as the average anomaly correlation, calculated using Fisher's  $z$ -transform, of the ensemble members with the ensemble mean; note that high average correlation of ensemble members with the ensemble mean corresponds to low ensemble spread. (a) The Met. Office Unified Model (UM); (b) ECMWF T63 L31 model (T63); (c) JT2 ensemble, a combination of the UM and T63 ensembles. The total entries in the diagonal quadrants (constructed using the ensemble median values of ACC skill and spread) are given. Years corresponding with PC/W events (see text for definition) are indicated with symbols, and the identities of other individual years are denoted with letters. The linear skill/spread correlation (SSC) is also given.

The skill/spread correlations found for all four seasons over Europe are summarized in Table 2. For the UM, correct skill assessments exceed incorrect assessments in all seasons except SON, while for the T63 a positive skill/spread correlation is found only for DJF. There is evidence that the JT2 ensemble provides skill/spread correlations that are similar to or better than the best individual model (note improvements in the DJF season from skill/spread correlations of 8/6 to 10/4 with JT2).

TABLE 2. SKILL-SPREAD CORRELATIONS FOR EUROPE

	UM	T63	JT2
MAM	<b>11/2</b>	4/8	<b>9/3</b>
JJA	<b>8/6</b>	6/8	<b>8/5</b>
SON	5/7	4/9	6/6
DJF	<b>8/6</b>	<b>8/6</b>	<b>10/4</b>

TABLE 3. SKILL-SPREAD CORRELATIONS FOR NORTH AMERICA

	UM	T63	JT2
MAM	6/7	6/7	4/9
JJA	<b>8/5</b>	<b>7/6</b>	<b>9/4</b>
SON	3/10	<b>8/6</b>	5/9
DJF	<b>8/6</b>	<b>8/6</b>	<b>10/4</b>

Total number of ensemble-mean simulations correctly/incorrectly identified as either above median or below median anomaly correlation skill, by reference to the median spread of the ensemble. Bold type is used where the number of correct skill assessments exceeds incorrect assessments. Results are for 15 MAM, JJA and SON seasons and 14 DJF seasons. Note that for MAM, JJA and SON the sum of the correct and incorrect assessments varies, as simulations with median skill and/or spread, of which there may be more than one, are not counted.

(ii) *North America.* The skill/spread correlations found for all four seasons over North America are summarized in Table 3. Positive skill/spread correlations are obtained with all three ensembles in DJF and JJA with the best correlation achieved in both periods by the JT2 ensemble. Skill enhancements during PC/W years in this region are a maximum in DJF (section 4(a)), and indeed ACC scores in all PC/W years are close to or above the median value in all three ensembles (Fig. 14(a)–(c)). As for the MAM simulations over Europe, there is tentative evidence that reference to both ENSO and ensemble spread might provide the most effective skill prediction strategies. For the T63 model, for example, requiring a PC/W event *or* low spread identifies eight of the 11 simulations with ACC score better than 0.2 while still successfully rejecting the remaining three years with negative ACC scores (however, the same strategy is less successful with the UM).

Comparison of Tables 2 and 3 indicates that the UM achieves better skill/spread correlations over Europe, while the T63 performs better over North America. In both regions the JT2 ensemble provides improvements to the skill/spread relationship in cases when both individual models have positive skill/spread correlations.

## 5. USE OF PERSISTED SSTA FOR LOWER-BOUNDARY FORCING

A relatively cheap option for a ‘two-tier’ representation of the ocean/atmosphere system is to use a persistence-based forecast of SST, in which SSTA over a period preceding the initial time of the forecast are persisted throughout the integration. In this section we compare UM skill obtained using persisted SSTA from the month preceding the initial time, with the skill obtained using observed SST. Details of the method used to produce the persistence-based SST forecasts are given in section 2(b). The comparison has been performed for 12 of the 15 PROVOST MAM and DJF periods for 850 hPa temperature and precipitation using the ROC score discussed in section 3.

In the tropics ROC scores obtained using persisted SSTA are consistently lower than, but nevertheless comparable with, the estimated upper limit on scores (available with the current UM) provided by the observed-SST runs (Table 4(a)). Scores with persisted SSTA are depressed by no more than 5%, and are above the 0.5 threshold for skill for both 850 hPa temperature and precipitation. For the northern hemisphere (Table 4(a)) and the European and North American regions (Table 4(b)), results are mixed; scores obtained using persisted SSTA are equivalent to, or even higher than, the observed-SST runs in about half the observed/persisted pairs (note that DJF scores with persisted SSTA

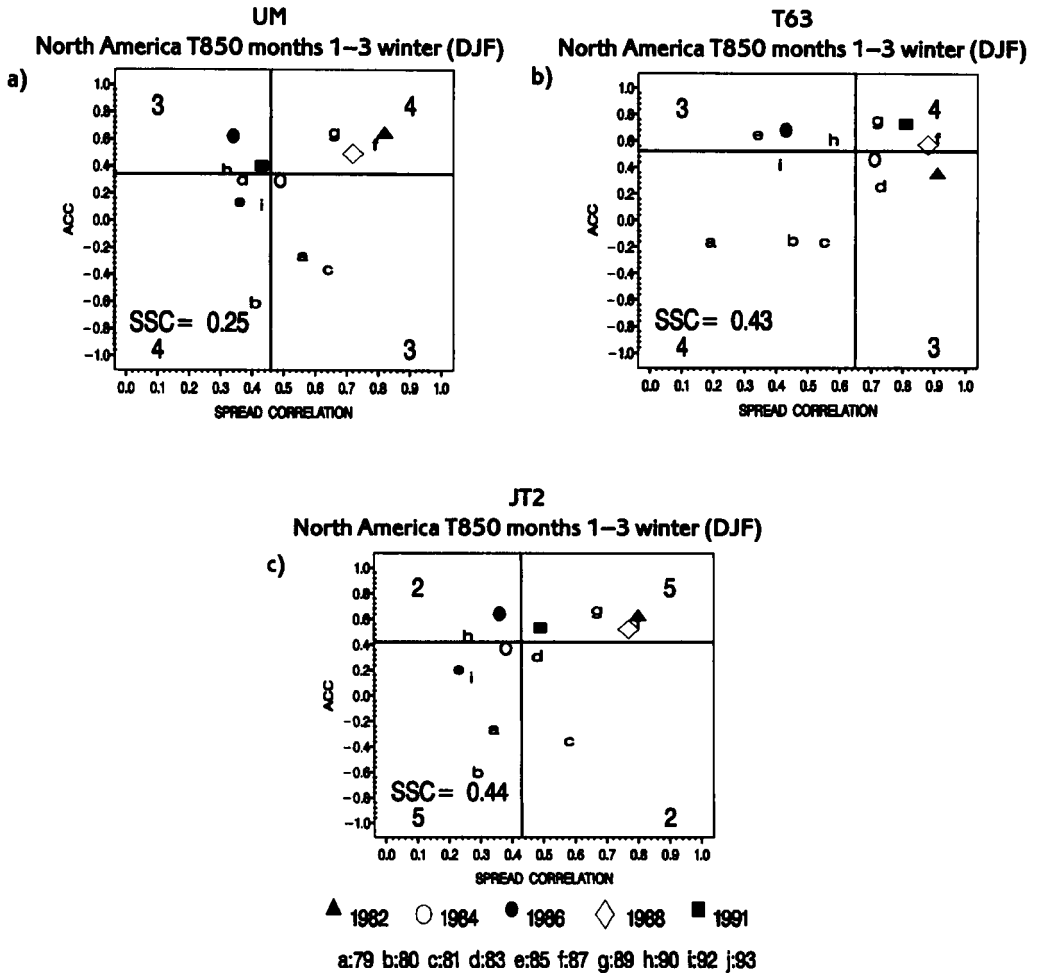


Figure 14. As Fig. 13 but for December–January–February (DJF) over North America.

are equal or higher in 5 out of 6 pairs). Differences (positive or negative) are mainly small, except for Europe in MAM for which the ROC score for 850 hPa temperature drops from 0.59 (observed SST) to 0.46 (persisted SSTA). Case-study investigations of individual JJA and SON seasons indicate that skills with persisted SSTA and observed SST are comparable regardless of the season.

The apparent near equivalence of skill with persisted SSTA and observed SST in extratropical regions is striking, since the persistence strategy is likely to be less effective at these latitudes. The results suggest that the extratropical SST field has little impact on predictability, and that the model representation of teleconnection processes between tropical SST and the extratropical atmosphere, most important in PC/W years (see section 4(a)), is not noticeably degraded, on average, by typical errors in tropical SST associated with the persisted SSTA method.

Although the accuracy of persisted SSTA is higher in the tropics than in the extratropics the dependence of predictability on SST is stronger, and errors in the SST field will impact on skill through both local and teleconnection processes. Thus the

TABLE 4. COMPARISON OF ROC SCORES OBTAINED WITH OBSERVED AND PERSISTED SSTA

(a)		TROPICS		N. HEMISPHERE	
		T850	PRECIP	T850	PRECIP
MAM	OBS	0.72	0.64	0.62	0.56
	PERS	0.69	0.61	0.59	0.55
DJF	OBS	0.72	0.63	0.61	0.56
	PERS	0.69	0.60	0.62	0.56
(b)		EUROPE		N. AMERICA	
		T850	PRECIP	T850	PRECIP
MAM	OBS	0.59	0.55	0.60	0.56
	PERS	0.46	0.52	0.60	0.52
DJF	OBS	0.55	0.52	0.68	0.57
	PERS	0.57	0.54	0.71	0.56

Comparisons are made using observed SST (OBS) and persisted SSTA (PERS) over 12 MAM and DJF periods from 1982 to 1993. Scores are for the events: T850, months 1–3 850 hPa temperature, above (or below) normal; and PRECIP, months 1–3 total precipitation, above (or below) normal.

relatively small average reductions of skill with persisted SSTA for the tropics are very encouraging. In individual cases loss of skill is likely to be greatest when marked anomalies in the SST field (e.g. in PC/W years) develop rapidly after the forecast is initialized.

Noting the above caveat, the results of this section suggest that forecasts of SST based on persisted SSTA are of sufficiently quality, at least up to one season ahead, to achieve average skill approaching the estimated potential upper bound (with current AGCMs). We can conclude that use of persisted SSTA as forcing to AGCM ensemble integrations appears to be a competitive (and cost effective) method for real-time seasonal forecasting, at least for a range of up to one season ahead.

## 6. USER VALUE OF SEASONAL PREDICTIONS

In this section we develop a methodology for assessing the user value of forecasts, following the work of Murphy (1977, 1985, 1994), and apply it to investigate the potential value of dynamical seasonal predictions over Europe. A similar approach is also discussed in Palmer *et al.* (2000) and for medium-range forecasts by Richardson (2000). Table 5 gives the cost/loss matrix for a user wishing to act on predictions of an adverse weather event. Each of the four contingencies in Table 5, hit, miss, false alarm or correct rejection, is associated with a financial impact, or loss; they are denoted here by  $L_h$ ,  $L_m$ ,  $L_f$  and  $L_c$ , respectively, and explained in Table 5. For convenience the losses are measured relative to the 'normal' loss associated with a correct rejection, i.e. the event is not forecast and is not observed, so that  $L_c = 0$ ; however, results obtained without this assumption are identical. The frequency of the four contingencies must be established (through 'track record' validation of the prediction system) for a range of forecast probability thresholds of the event. For any one probability threshold, these frequencies are denoted here by  $h$  (hit),  $m$  (miss),  $f$  (false alarm) and  $c$  (correct rejection); and the expected mean expense of taking action on the forecast ( $ME_{fx}$ ) whenever the forecast probability exceeds the given threshold is thus:

$$ME_{fx} = hL_h + mL_m + fL_f. \quad (1)$$



TABLE 5. THE USER COST/LOSS MATRIX

		FORECAST	
		YES	NO
OBSERVED	YES	<b>HIT, <math>h</math></b> Losses occur but are mitigated through taking appropriate protective action Loss = $L_h$	<b>MISS, <math>m</math></b> Full losses occur, as no protective action is taken. Loss = $L_m$
	NO	<b>FALSE ALARM, <math>f</math></b> Costs of protection are borne, but protection is not required. Loss = $L_f$	<b>CORRECT REJECTION, <math>c</math></b> No losses; no costs; outcome is that expected from normal activities without event. Loss = $L_c = 0$

The cost/loss table is for assessing the financial value of predictions for an adverse weather event. Lower-case letters represent the frequency of occurrence of each contingency. Losses are denoted by  $L$  with corresponding subscripts. Note that  $L_c = 0$  by definition, and other losses are scaled relative to this baseline.

Note that  $h$ ,  $m$ ,  $f$  and  $c$  are related to the hit rates ( $HR$ ) and false-alarm rates ( $FAR$ ) of the ROC verification method (Table 1(b)) by:

$$h = oHR; m = o(1 - HR); f = (1 - o)FAR; c = (1 - o)(1 - FAR),$$

where  $o$  is the overall frequency of the event. Thus the ROC contingency tables (Table 1(b)), and associated curves (e.g. Fig. 1) provide an expedient technical validation of the prediction system, as the results may also be used in estimating forecast value. Note also that the definition of the mean expense (1) is an extension of the cost/loss ratio formulation (see e.g. Murphy 1985) in which protective action is assumed to eliminate all loss associated with an adverse weather event (i.e. imposing  $L_h = L_f$ ). Relaxing this restriction allows representation of circumstances in which protection is not fully effective ( $L_h > L_f$ ). Moreover, it also allows representation of financial benefits that may result from actions taken on advance warnings of an event which later occurs (i.e. a hit); such benefits may offset expenditure on protection so that  $L_h < L_f$ . For example, protective action for a farmer might include switching part or all of a planned seed crop to one more suited to the probable expected seasonal conditions. Resulting high yields from the new crop if the event occurs may offset the additional expenditure incurred.

It is readily shown that the value,  $V$ , of the forecast system relative to use of climatology, and expressed as a fraction of the maximum possible value obtained by following a perfect forecast system in which protective action is taken only when the event occurs (no misses or false alarms), is given by:

$$V = \frac{ME_{cl} - ME_{fx}}{ME_{cl} - ME_p},$$

where  $ME_{cl} = \min\{oL_m, oL_h + (1 - o)L_f\}$ , in which  $o$  is the climatological frequency of the event, gives the mean expense of the best climate option; the first term in the argument gives the mean expense of never protecting, the second term the mean expense of always protecting, and  $ME_p = oL_h$  is the mean expense incurred on following a

perfect forecast system.  $V$  has a maximum value of 1 for a perfect system, while for a forecast system that is no better than climate  $V = 0$ ; note, however, that there is no lower bound, and for forecasts with low hit rates or high false-alarm rates (large  $m$  or  $f$ ), or when large losses are associated with misses and false alarms ( $L_m$  or  $L_f$ ),  $V$  may be negative.

Forecast value may be evaluated in this way for a range of forecast probability thresholds of the event, and plotted as a function of the threshold. In this way the user may select the probability threshold that provides the greatest value. An example based on the ROC verifications of the PROVOST simulations presented in section 3 is provided in Fig. 15. The weather event considered is MAM 850 hPa temperature above normal over Europe. Value obtained both from probability forecasts and from deterministic forecasts based on the ensemble mean is considered for the three ensembles UM, T63 and JT2. The cost/loss matrix is defined as  $L_h = 1$ ,  $L_m = 4$  and  $L_f = 2$ . The derivation of the loss values will in general be complex, depending on details of the user sensitivity and planned response to each of the four contingencies. For a simple interpretation, however, we may assume in this example that on forecasts of an 'adverse' weather event the user spends 2 units on protection, a sum equivalent to half the expected loss in the event of a miss, and that financial benefits in the event of a hit offset expenditure on protection by 1 unit. Figure 15 shows that with these loss estimates potential skill of the probability forecasts is indeed sufficient to obtain value. Highest value in this case, equal to approximately 7% of the value of a perfect forecast system, is achieved by the UM ensemble employing a probability threshold for protection of 20%. Similar results are found with the JT2 ensemble. Note that in this case value is not available from the deterministic forecasts, for which all three models show negative value; in general it is found that value obtained from use of the optimum probability threshold exceeds that from use of deterministic forecasts (see e.g. Richardson 2000). This result demonstrates that issuing seasonal predictions probabilistically is likely to extend the range of applications that can extract value from the forecasts.

Note that, in this example, the user is not served by a cautious approach of waiting for higher certainty of the event, since value is negative for probability thresholds in excess of about 40%. Moreover, examination of a wide range of examples demonstrates (as might be expected) that value is very sensitive to the user losses, and will therefore vary with forecast application.

In practice, users of seasonal forecasts are likely to require more precision than the simple two-category (above/below) classification used in the above analysis, and may require more resolution in geographical region. In addition the existence of favourable cost/loss matrices vis-a-vis the forecast events is yet to be fully explored. Nevertheless, the examples discussed demonstrate that there is potential for user value from seasonal predictions over Europe. For optimum value, close co-operation between forecast supplier and user appears necessary, to reach agreement on the meteorological events which are both relevant to the user and sufficiently predictable on seasonal time-scales. In addition the user may need to revise and adapt the cost/loss estimates (e.g. through changed responses to the four contingencies) to make the most of the forecast system performance.

The above examples were chosen to demonstrate potential value of seasonal predictions in the extratropics. Clearly, potential value will be at greater levels in tropical areas where model skill is at a higher level. Making use of a similar method, with cost/loss values supplied by a range of users, M. Harrison and N. Graham (personal communication) have found that given current marginal costs within regional National Meteorological Services in developing and supplying information on seasonal predictions, and

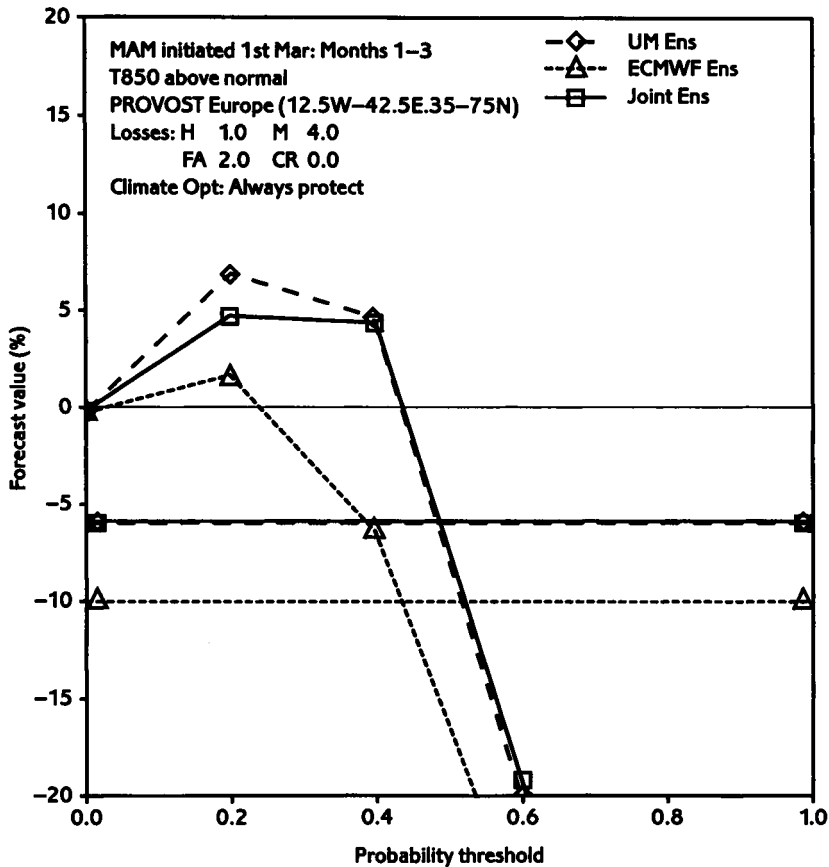


Figure 15. Potential value ( $V$ ) of months 1–3 MAM simulations of the event 850 hPa temperature above normal over Europe, obtained with the UM (The Met. Office Unified Model), T63 (ECMWF T63 model) and JT2 (UM and T63 combined) ensembles assuming user losses of  $L_h = 1$ ,  $L_m = 4$ ,  $L_f = 2$ .  $V$  is expressed as a percentage of the value obtained by following a perfect forecast system. Plotted curves show the potential values of probabilistic predictions at thresholds of 0%, 20%, 40% ... 100%, while horizontal lines show corresponding potential values of deterministic predictions based on the ensemble mean. UM, diamonds and dashed line; T63, triangles and dotted line; JT2, squares and solid line. See text for full explanation and details.

assuming current levels of skill, the cost/benefit ratio of seasonal precipitation forecasts over the southern African region can be estimated to lie between 1:20 and 1:200.

## 7. SUMMARY AND CONCLUSIONS

Seasonal simulations from the 15-year PROVOST database have been analysed to assess the potential skill of 9-member ensemble integrations of the UM and the T63 model for seasonal prediction. A joint 18-member ensemble (JT2) produced by combining all members of the UM and T63 ensembles, and higher-order multiple-model ensembles employing all four participating PROVOST models, have also been assessed. The ensembles were integrated using observed ('near-perfect' predicted) SSTs to force the lower boundary, and the skill attained therefore represents an upper bound with the AGCMs employed. UM integrations were also performed replacing observed SST forcing with persisted SSTA.

All skill measures calculated for the entire tropical and northern extratropical regions indicate that, while skill is highest in the tropics, it is also available over the northern extratropics with all models in all seasons. Skill is present for both 850 hPa temperature and precipitation, though for the latter it is generally at a lower level than for the former. Nevertheless, there is evidence of substantial potential for rainy season predictions in some tropical regions. In both the tropics and extratropics, highest skill tends to be centred over oceanic regions and lowest skill over continental interiors. Best skill in the northern extratropics is found in the MAM season for both 850 hPa temperature and precipitation.

Over North America skill is found for both 850 hPa temperature and precipitation in all four seasons with both models. Scores are highest, both for 850 hPa temperature and precipitation, in the DJF season, apparently as a result of enhanced winter predictability of the PNA mode during PC/W events. Scores are lowest for precipitation in JJA and SON. Over Europe skill for 850 hPa temperature is present in all seasons with at least one model. Skill is highest and most consistent between models in MAM (as for the northern extratropics). As for North America, ROC scores for precipitation over Europe indicate best skill in MAM and DJF. Little skill for precipitation is evident in the JJA and SON periods. Spatial variations of skill within the assessment regions, both for 850 hPa temperature and precipitation, indicate that over Europe highest skill tends to be concentrated in northern and western regions, while over North America highest levels of skill are located over southern, western and northern regions.

Performance differences between the UM and T63 models are most pronounced over the regional areas of North America and Europe, the UM generally achieving better scores over Europe and the T63 achieving better scores over North America. Thus useful complementary skill is present between the individual models which could be exploited in an operational environment. In these regions the skill filtering property of the JT2 ensemble provides substantial benefits, achieving ROC scores similar to or better than the more skilful individual model, even when skill differentials between the individual models are relatively large. Benefits from the JT2 ensemble are less apparent for the ACC skill scores, but notable improvements are found in some seasons. Higher-order multiple-models (JT3 and JT4) show further improvements over JT2. The skill of the multiple-models appears to be mainly a function of the skill of the most skilful component ensemble, rather than principally being related to the increased ensemble size. An important benefit of the multiple-model method is that it improves potential capability without the need for *a priori* identification of the strengths of the individual component models, as would be needed, for example, if a strategy of choosing the best model for each region were adopted.

PC/W events in the tropical east Pacific enhance model predictability, on average, over North America and Europe in some seasons, with similar results found for the UM and T63 ensembles. The largest average skill enhancement for North America is found for the DJF season, apparently as a result of increased predictability of the PNA pattern. Over Europe skill enhancement in PC/W years is largest, on average, in MAM. In non-PC/W years skill over Europe and North America is similar, suggesting that the greater predictability frequently attributed to the North American region applies mainly to years with PC/W events. In non-PC/W years skill over North America and Europe is generally lower for months 2–4 than for months 1–3; however, in PC/W years skill for months 1–3 appears to be maintained into the months 2–4 period.

Ensemble consistency appears to be promising as a predictor of ensemble-mean skill in some seasons over Europe and North America. Best results are achieved with different individual models in different regions/seasons, however in most cases the

JT2 multiple-model ensemble improves the skill/spread correlations attained with the individual models.

The sample of 15 simulations is too small to define clear skill prediction strategies. However, there is evidence that skill prediction strategies, perhaps based on the state of ENSO and on ensemble consistency, can be developed. As anticipated, optimum strategies are found to be likely to depend on the prediction model, the region and the season.

Comparisons of integrations using persisted SSTA and observed SST as boundary forcing indicate that, on average, a substantial proportion of the skill achieved using observed SST is retained using persisted SSTA, both in the tropics and in the extratropics (though there will be local variations in the proportion of skill retained). Thus the use of persistence-based forecasts of SST appears to be a viable method for real-time seasonal forecasting, at least for a range of one season ahead.

For future development of operational seasonal prediction it will be crucial to establish the levels of technical skill (i.e. as measured using skill scores) required in order for seasonal predictions to be of value to users. A methodology for linking technical forecast quality with financial value for users has been outlined using the ROC and the user cost/loss matrix. Results employing an assumed user cost/loss matrix indicate promising potential for user value of seasonal predictions not only over tropical areas but also in some extratropical areas, including Europe.

#### ACKNOWLEDGEMENTS

The help of other participants in the PROVOST project is acknowledged; specifically Michael Davey (The Met. Office), Tim Palmer (PROVOST co-ordinator and ECMWF), Michel Déqué (Météo-France) and Jean-Yves Canneil (Electricité de France). Discussions with Tim Palmer and Čedo Branković are also gratefully acknowledged. The persisted SSTA fields referred to in section 5 were provided by Sarah Ineson and Michael Davey of the Ocean Applications branch of The Met. Office. This work was supported by the European Union Environment and Climate Programme under contract ENV4-CT95-0109.

#### REFERENCES

- |   |      |   |
|---|------|---|
| Barnett, T. P., Arpe, K., Bengtsson, L., Ji, M. and Kumar, A.   | 1997 | Potential predictability and AMIP implications of midlatitude climate variability in two general circulation models. <i>J. Climate</i> , <b>10</b> , 2321–2329                                      |
| Brown, B. H. and Murphy, A. H.                                  | 1996 | Improving forecasting performance by combining forecasts: the example of road-surface temperature forecasts. <i>Met. Apps.</i> , <b>3</b> , 257–265   |
| Branković, Č. and Palmer, T. N.                                 | 2000 | Seasonal skill and predictability of ECMWF PROVOST ensembles. <i>Q. J. R. Meteorol. Soc.</i> , <b>126</b> , 2035–2067   |
| Branković, Č., Palmer, T. N. and Ferranti, L.                   | 1994 | Predictability of seasonal atmospheric variations. <i>J. Climate</i> , <b>6</b> , 217–237   |
| Buizza, R.  | 1997 | Potential forecast skill of ensemble prediction and spread and skill distributions of the ECMWF prediction system. <i>Mon. Weather Rev.</i> , <b>125</b> , 99–119                                   |
| Cullen, M. J. P.  | 1991 | 'Introduction to the unified forecast/climate model'. The Met. Office FR Division Scientific Paper No. 1. Available from the National Meteorological Library, London Road, Bracknell, Berkshire, UK |
| Evans, T., Graham, R., Harrison, M., Davey, M. and Colman, A.   | 1998 | A dynamic one-month lead seasonal rainfall prediction for July to September 1998 for North Africa from 20°N to the equator. <i>Experimental Long-lead Forecast Bull.</i> , <b>7</b> , No. 2, 38–42  |
| Evans, R. E., Harrison, M. S. J., Graham R. J. and Mylne, K. R. | 2000 | Joint medium-range ensembles from The Met. Office and ECMWF systems. <i>Mon. Weather Rev.</i> , <b>128</b> , in press   |

- Gibson, J. K., Kållberg, P., Uppala, S., Hernandez, A., Nomura, A. and Serano, E. 1997 'ERA description'. ECMWF re-analysis project report series. Available from ECMWF, Shinfield Park, Reading, UK
- Hall, C. D., Stratton, R. A. and Gallani, M. L. 1995 'Climate simulations with the Unified Model: AMIP runs'. The Met. Office Climate Research Technical Note 61. Available from the National Meteorological Library, London Road, Bracknell, Berkshire, UK
- Harrison, M., Evans, T., Evans, R., Davey, M. and Colman, A. 1997a A dynamical one-month lead seasonal rainfall prediction for March to May 1997 for the north-eastern area of South America. *Experimental Long-lead Forecast Bull.*, **6**, No. 1, 25–28
- Harrison, M., Soman, M. K., Davey, M., Evans, T., Robertson, K. and Ineson, S. 1997b Dynamical seasonal prediction of the Indian summer monsoon. *Experimental Long-lead Forecast Bull.*, **6**, No. 2, 29–32
- Krishnamurti, T. N., Kishtawal, C. M., LaRow, T. E., Bachiochi, D. R., Zhang, Z., Williford, C. E., Gadgil, S. and Surrendan, S. 1999 Improved weather and seasonal climate forecasts from a multi-model superensemble. *Science*, **285**, 1548–1550
- Molteni, F., Buizza, R., Palmer, T. N. and Petroliagis, T. 1996 The ECMWF Ensemble Prediction System: Methodology and validation. *Q. J. R. Meteorol. Soc.*, **122**, 73–119
- Murphy, A. H. 1977 The value of climatological, categorical and probabilistic forecasts in the cost–loss ratio situation. *Mon. Weather Rev.*, **105**, 803–816
- 1985 Decision making and the value of forecasts in a generalised model of the cost–loss ratio situation. *Mon. Weather Rev.*, **113**, 362–369
- 1994 Assessing the economic value of weather forecasts: an overview of methods, results and issues. *Met. Apps.*, **1**, 69–73
- Palmer, T. N. and Anderson, D. L. T. 1994 The prospects for seasonal forecasting—A review paper. *Q. J. R. Meteorol. Soc.*, **120**, 755–793
- Palmer, T. N., Branković, Č. and Richardson, D. S. 2000 A probability and decision-model analysis of PROVOST seasonal multi-model ensemble integrations. *Q. J. R. Meteorol. Soc.*, **126**, 2013–2033
- Rayner, N. A., Horton, E. B., Parker, D. E., Folland, C. K. and Hackett, R. B. 1996 'Version 2.2 of the global sea-Ice and sea surface temperature dataset, 1903–1994'. CTRN 74, Hadley Centre for Climate Prediction and Research, The Met. Office, Bracknell, Berkshire, UK
- Reynolds, R. W. and Smith, T. M. 1994 Improved global sea surface temperature analyses using optimum interpolation. *J. Climate*, **7**(6), 929–948
- Richardson, D. S. 2000 Skill and relative economic value of the ECMWF ensemble prediction system. *Q. J. R. Meteorol. Soc.*, **126**, 649–667
- Rowell, D. P. 1998 Assessing potential seasonal predictability with an ensemble of multidecadal GCM simulations, *J. Climate*, **11**, 109–120
- Stanski, H. R., Wilson, L. J. and Burrows, W. R. 1989 'Survey of common verification methods in Meteorology'. World Weather Watch Technical Report. No. 8, WMO/TD 358, World Meteorological Organization, Geneva, Switzerland