

1 **Multi-model forecast skill for mid-summer rainfall over**
2 **southern Africa**

3
4 **Willem A. Landman^{a,b}**

5 **Asmerom Beraki^c**
6

7 ^a Council for Scientific and Industrial Research, Natural Resources and the
8 Environment, Pretoria, South Africa

9 ^b Department of Geography, Geoinformatics and Meteorology, University of
10 Pretoria, South Africa

11 ^c South African Weather Service, Pretoria, South Africa
12

13 Re-submitted to: *International Journal of Climatology*

14 29 October 2010
15

16 Correspondence to: Council for Scientific and Industrial Research, P.O. Box 395,
17 Pretoria, 0001, South Africa; E-mail: WALandman@csir.co.za

18 Tel: +27-12-841-3395

19 Fax: +27-12-841-4863
20

21 *Key words:* multi-model, downscaling, seasonal forecasting, ENSO, southern
22 Africa.
23

1 **ABSTRACT**

2 Southern African December-January-February (DJF) probabilistic rainfall forecast
3 skill is assessed over a 22-year retro-active test period (1980/81 to 2001/02) by
4 considering multi-model ensembles consisting of downscaled forecasts from
5 three of the DEMETER models, the ECMWF, Météo-France and UKMO coupled
6 ocean-atmosphere general circulation models. These models are initialized in
7 such a way that DJF forecasts are produced at an approximate 1-month lead-
8 time, i.e., forecasts made in early November. Multi-model forecasts are obtained
9 by 1) downscaling each model's 850 hPa geopotential height field forecast using
10 canonical correlation analysis (CCA) and then simply averaging the rainfall
11 forecasts, and 2) by combining the three models' 850 hPa forecasts and then
12 downscaling them using CCA. Downscaling is performed onto the 0.5°x0.5°
13 resolution of the CRU rainfall data set south of 10° south over Africa. Forecast
14 verification is performed using the relative operating characteristic (ROC) and the
15 reliability diagram. The performance of the two multi-model combinations
16 approaches are compared with the single model downscaled forecasts and also
17 with each other. It is shown that the multi-model forecasts outperform the single
18 model forecasts, that the two multi-model schemes produce about equally skilful
19 forecasts, and that the forecasts perform better during El Niño and La Niña
20 seasons than during neutral years.

21

22

1 **1. Introduction**

2

3 The scientific basis for the existence of seasonal climate predictability originates
4 from the observation that slowly evolving sea-surface temperature (SST)
5 anomalies influence seasonal-mean weather conditions (Palmer and Anderson
6 1994). Therefore, estimation of the evolution of SST anomalies, which are often
7 relatively predictable, and subsequently employing them in atmospheric general
8 circulation models (GCMs), potentially provides means of generating forecasts of
9 seasonal-average weather (Graham *et al.* 2000). With the advent of fully coupled
10 ocean-atmosphere models (Stockdale *et al.*, 1998; Saha *et al.*, 2006; Weisheimer
11 *et al.* 2009), evidence that the ocean models participating in fully coupled GCMs
12 can predict the evolution of SSTs to elevated levels of skill has been presented.
13 This notion has been demonstrated conclusively through the DEMETER
14 (Development of a European Multimodel Ensemble system for seasonal to
15 interannual prediction) project (Palmer *et al.* 2004), and recently the usefulness
16 of these forecasts over the mid-latitudes has been further demonstrated (Coelho
17 *et al.* 2006; Frias *et al.* 2010). In theory coupled models should eventually
18 outperform using GCMs as a second step in a 2-tiered system in which SSTs are
19 first predicted since the former is able to describe the feedback between ocean
20 and atmosphere while the latter assumes that the atmosphere responds to SST
21 but does not in turn affect the oceans (Copsey *et al.*, 2006; Troccoli *et al.*, 2008).

22

1 Although GCMs, commonly configured with an effective resolution of 100-300
2 km, have demonstrated skill at global or even continental scale, they are unable
3 to represent local sub-grid features, subsequently overestimating rainfall over
4 southern Africa (Joubert and Hewitson 1997; Mason and Joubert 1997). Also, the
5 representation of rainfall at mid-to-high latitudes is complex and often not well
6 estimated (Graham *et al.* 2000; Goddard and Mason 2002). Such systematic
7 biases have created the need to downscale GCM simulations over southern
8 Africa. Semi-empirical relationships exist between observed large-scale
9 circulation and rainfall, and assuming that these relationships are valid under
10 future climate conditions and also that the large-scale structure and variability is
11 well characterized by GCMs, mathematical equations can be constructed to
12 predict local precipitation from the forecast large-scale circulation (Landman and
13 Goddard, 2002; Wilby and Wigley 1997). Empirical remapping of GCM fields to
14 regional rainfall has been demonstrated successfully over southern Africa
15 (Bartman *et al.* 2003; Landman and Goddard 2002, 2005; Landman *et al.* 2001;
16 Shongwe *et al.*, 2006).

17

18 The chaotic inherent variability of the atmosphere requires seasonal climate
19 simulations to be expressed probabilistically. Probabilistic forecasts are made
20 possible through the proper use of GCM ensembles since ensemble forecasting
21 is a feasible method to estimate the probability distribution of atmospheric states
22 (Branković and Palmer 2000). In addition, errors in the initial conditions as well as
23 deficiencies in the parameterizations and systematic or regime-dependent model

1 errors can be to a large part accounted for through ensemble forecasting (Evans
2 *et al.* 2000). Moreover, there is inevitable growth in differences between forecasts
3 started from very slightly different initial conditions suggesting that there is no
4 single valid solution but rather a range of possible solutions (Tracton and Kalnay
5 1993). Information contained in the distribution of the ensemble members can
6 subsequently be used to represent forecast probabilities by calculating the
7 percentage of ensemble members that fall within a particular category (e.g.
8 below-normal, near-normal or above-normal). Similarly, forecast probabilities can
9 be produced indicating the percentage of ensemble members in the upper or
10 lower extremes, e.g., 15th percentiles (Mason *et al.* 1999).

11

12 There are advantages in combining ensemble members of a number of GCMs
13 into a multi-model ensemble since GCMs differ in their parameterizations and
14 therefore differ in their performance under different conditions (Hagedorn *et al.*,
15 2005). Using a suite of several GCMs not only increases the effective ensemble
16 size; it also leads to probabilistic simulations that are skilful over a greater portion
17 of the region and a greater portion of the time series. Multi-model ensembles are
18 nearly always better than any of the individual models (Dirmeyer *et al.* 2003,
19 Doblas-Reyes *et al.* 2000, 2005, Hagedorn *et al.*, 2005; Krishnamurti *et al.* 2000).
20 The benefits from combining ensembles are a result of the inclusion of
21 complementary predictive information since the forecast scheme is able to
22 extract useful information from the results of individual models from local regions
23 where their skill is higher (Krishnamurti *et al.* 2000). In fact, the most striking

1 benefit obtained from multi-model ensembles is the skill-filtering property in
2 regions or seasons when the performance of the individual models varies widely
3 (Graham *et al.* 2000). Moreover, increased ensemble size leads to further
4 benefits (Brown and Murphy 1996), but the multi-model approach is only
5 beneficial if the individual models produce independent skilful information
6 (Graham *et al.* 2000). A number of ensemble combining algorithms exists. The
7 most simple of these is the unweighted combination of ensembles from different
8 models (Hagedorn *et al.* 2005; Graham *et al.* 2000, Mason and Mimmack 2002;
9 Peng *et al.* 2002; Tippet and Barnston 2008). The improvements of a multi-model
10 over the individual ensemble systems are attributed to the collective information
11 of all the models used in the mean of probabilities algorithm. However, the
12 forecast quality of a simple multi-model ensemble is often difficult to improve on
13 when the available sample size is relatively small (Doblas-Reyes *et al.* 2005).

14

15 An association exists between South Africa's summer seasonal rainfall and the
16 equatorial Pacific Ocean. However, the association in the middle to late austral
17 summer season is higher than earlier in the summer rainy season (e.g., Tyson
18 and Preston-Whyte, 2000), and it is also non-linear (Fauchereau *et al.* 2008).
19 Notwithstanding, in the mid-summer months South Africa tends to be
20 anomalously dry during El Niño years and anomalously wet during La Niña years,
21 although wet El Niño seasons and dry La Niña seasons are not uncommon.
22 Indian and Atlantic Ocean SST also have a statistically detectable influence on
23 South African rainfall variability (e.g., Mason, 1995; Reason *et al.*, 2006).

1 Moreover, while the El Niño-Southern Oscillation (ENSO) has a control on rainfall
2 variability over the southern African region, Indian Ocean SST anomalies,
3 sometimes varying independently of ENSO, are important for the skilful
4 simulation of southern African seasonal rainfall variability using atmospheric
5 GCMs (e.g., Washington and Preston, 2006). Since ENSO is the dominant mode
6 of seasonal and interannual climate variability globally, and since ENSO has a
7 strong influence on southern African rainfall, it needs to be investigated to what
8 extent ENSO influences coupled model performance over southern Africa.

9

10 The paper consists of three parts: 1) single coupled model downscaled forecast
11 performance during mid austral summer over southern Africa compared with that
12 of multi-models, 2) the comparison between unweighted and weighted
13 combination of forecasts, and 3) multi-model performance during ENSO and
14 during neutral years. For the second part, the unweighted combination involves
15 downscaling and correcting GCM output first before combining, while for the
16 weighted combination weighting is done and then combined before downscaling
17 and correcting.

18

19

1 **2. Data, models and methods**

2

3 2.1. Rainfall data

4

5 The season of interest is December-January-February (DJF) when southern
6 Africa is being dominated by influences mainly from the tropics and so is a
7 season of relatively high predictability and ideal for seasonal predictability studies
8 over the region. The University of East Anglia Climatic Research Unit (CRU)
9 global 0.5° x 0.5° monthly data, Version 2.1 (Mitchell and Jones, 2005) are used
10 to construct DJF seasonal averaged rainfall totals for southern Africa south of 10°
11 south for the period 1959/60 to 2001/02. This data set is used for both empirical
12 downscaling and for forecast verification.

13

14 2.2. Coupled general circulation models

15

16 The atmosphere-ocean models used in this study are from the DEMETER project
17 (Palmer *et al.*, 2004) and in particular are the ECMWF, Météo-France and UKMO
18 coupled models. These models were selected since they each have 43 years of
19 available hindcast data, and the longer the record of archived model data the
20 better the chance is to develop robust empirical downscaling equations.
21 Hindcasts had been started from 1 November and nine ensemble members
22 created. Seasonal means are used in the study.

23

1 2.3. Model output statistics

2

3 Given the low spatial resolution of the coupled models (Palmer *et al.*, 2004) there
4 is a need to downscale the global model output to a higher resolution to satisfy
5 end-user needs and to further improve on the forecasts (Landman and Goddard,
6 2002) through the correction of systematic deficiencies in the global models
7 (Tippet *et al.*, 2005). Model output statistics (MOS; Wilks, 2006) equations are
8 developed here because they can compensate for these errors in the model
9 fields directly in the regression equations. The reason why these errors can be
10 overcome is because MOS uses predictor values from the global models in both
11 the development and forecast stages. Notwithstanding, the selection of the
12 appropriate model field require careful consideration: Raw model forecast of
13 rainfall that is a result of, for example, the interaction between atmospheric
14 circulation and topography is poorly resolved, and may therefore not be a good
15 predictor of rainfall observed at ground level. Rainfall fields, even when totalled
16 over a season, are noisy, and normally contain structures on spatial scales well
17 below those resolved by the models. However, variables such as large-scale
18 circulation are more accurately simulated by models than rainfall and should
19 therefore be used instead in a MOS system to predict seasonal rainfall totals
20 (Landman and Goddard, 2002).

21

22 The MOS equations are developed by using the canonical correlation analysis
23 (CCA; Barnett and Preisendorfer, 1987) option of the Climate Predictability Tool

1 (CPT). This tool was developed at the International Research Institute for Climate
2 and Society (IRI; <http://iri.columbia.edu>). The forecast fields from each GCM used
3 in the MOS are restricted over a domain that covers an area between the
4 Equator and 40°S, and Greenwich to 60°E. Empirical orthogonal function (EOF)
5 analysis is performed on both the predictor (model forecast fields) and predictand
6 sets (CRU data over southern Africa) prior to CCA, and the number of EOF and
7 CCA modes to be retained in the CPT's CCA procedure is determined using
8 cross-validation skill sensitivity tests. Both the models' ensemble mean rainfall
9 and 850 hPa geopotential height fields were separately considered over the
10 available 43-year period (1959/60 – 2001/02) to find out which of the two fields
11 provide the best first estimate for the downscaled forecasts. A 5-year-out cross-
12 validation design was selected and it was found that for both the ECMWF and
13 UKMO models, the height field is the better option, but for the Météo-France
14 model, rainfall was a slightly better performer. Notwithstanding, 850 hPa
15 geopotential heights were selected for all three models for consistency and
16 because of the potential problems mentioned above when rainfall as a
17 downscaling predictor field is used. Considering other model fields such as
18 moisture and geopotential heights at levels other than 850 hPa showed no further
19 benefits over only using the 850 hPa geopotential fields as a single predictor field
20 either.

21

22

1 2.4. Model combination

2

3 A number of forecast combining algorithms exists, but only two are considered
4 here. The first is the most simple of all combination schemes and involves
5 unweighted averaging of the forecast probabilities (e.g., Hagedorn *et al.*, 2005).

6 For this simple combination approach, the 850 hPa height forecasts from the
7 three coupled models are first separately downscaled to DJF rainfall at the 0.5° x
8 0.5° CRU resolution and then averaged, and is referred to here as a combination
9 using equal weights (MMeqw). The second approach allows the models to be
10 weighted by combining the 850 hPa geopotential height forecasts fields from the
11 models prior to EOF pre-filtering in the CCA process. Downscaling is then
12 performed as before, but with combined forecast fields (MMcca) as opposed to
13 individual model fields.

14

15 2.5. Retro-active forecasts

16

17 In order to minimize artificial inflation of forecast skill, the performance of the
18 individual models and the two multi-model systems (MMeqw and MMcca) should
19 be verified over a test period that is independent of the training period and should
20 involve evaluation of predictions compared to their matching observations
21 excluding any information following the forecast year. Such a system mimics a
22 true operational forecasting environment where no prior knowledge of the coming
23 season is available. The individual models and two multi-model systems are first

1 trained with information from 1959/60 and leading up to and including 1979/80.
2 The seasonal rainfall of the next year (1980/81) is subsequently predicted using
3 the trained models. The various MOS sets of equations are subsequently
4 retrained using information leading up to and including 1980/81 to predict for
5 1981/82 conditions. This procedure is continued until the 2001/02 DJF rainfall is
6 predicted using MOS systems trained with data from 1959/60 to 2000/01,
7 resulting in 22 years (1980/81 – 2001/02) of independent forecast data. In
8 estimating the skill in predicting DJF rainfall over southern Africa, the observed
9 and predicted fields are separated into three equi-probable categories based on
10 the preceding years' climatology defining above-normal, near-normal and below-
11 normal seasonal rainfall totals.

12

13 The distribution of individual ensemble members is intended to be able to
14 indicate forecast uncertainty. However, only a finite ensemble is available (9
15 members from each coupled model) suggesting that the forecast distribution may
16 be poorly sampled – and so the uncertainty associated with the forecasts has to
17 be estimated. Probabilistic MOS forecasts for each of the 22 retro-active years
18 are obtained here from the error variance of the cross-validated predictions using
19 the ensemble mean (Troccoli *et al.*, 2008) for each of the various training periods.
20 The errors in the predictions are assumed to be Gaussian. Cross-validation is
21 performed using a (large) 5-year-out window, which means that 2 years on either
22 side of the predicted year are omitted, in order to minimize the chance of
23 obtaining biased results.

1

2 This modelling study also focuses on one of the major sources of predictability
3 over southern Africa, namely the El Niño – Southern Oscillation (ENSO)
4 phenomenon, and how forcing from the equatorial Pacific Ocean influences
5 predictability over the region. The El Niño, La Niña and neutral years considered
6 are those listed by Coelho et al. (2006). Rainfall retro-active forecast skill over the
7 subcontinent is then assessed during El Niño (1982/83, 1986/87, 1987/88,
8 1990/91, 1991/92, 1992/93, 1994/95 and 1997/98 = 8 seasons), La Niña
9 (1983/84, 1984/85, 1988/89, 1995/96, 1998/99, 1999/00 and 2000/01 = 7
10 seasons) and neutral (1980/81, 1981/82, 1985/86, 1989/90, 1993/94, 1996/97
11 and 2001/02 = 7 seasons) events.

12

13 2.6. Estimating true forecast performance

14

15 For the generation of verification data we adopt an approach that minimizes the
16 inflation of forecast skill by testing the models in an environment that mimics that
17 of an operational centre, i.e. a retro-active forecast setting (Wilks, 2006).
18 However, owing to the limited archived model data set available the MOS
19 equations used for the prediction of the first part of the verification set may not
20 display a robust relationship between the predictor (850 hPa heights) and
21 predictand (rainfall at the surface) throughout the retro-active process, but this
22 problem should become less of an issue as the forecast process progresses

1 beyond about 30 years of training data. Notwithstanding, here we assume that
2 the relationships remain robust, a notion that will be tested later on in the paper.

3

4 Since seasonal climate is inherently probabilistic, seasonal forecasts should be
5 judged probabilistically. The main attributes of interest for probabilistic forecasts
6 are: 1) reliability (is the confidence communicated in the forecast appropriate and
7 are there systematic biases in the forecast probabilities?), 2) resolution (is there
8 any useable information in the forecast?), 3) discrimination (are the forecasts
9 discernibly different given different outcomes?), and 4) sharpness (what is the
10 confidence level that is communicated in the forecast?) (Troccoli *et al.*, 2008;
11 Wilks, 2006). The forecast verification measures are the reliability diagram
12 (Hamill 1997; Wilks, 2006) and the relative operating characteristic (ROC; Mason
13 and Graham, 1999; Wilks, 2006). A forecast system is deemed reliable if there is
14 consistency between predicted probabilities of an event such as drought/floods
15 (or below/above-normal rainfall in this paper) and the observed relative
16 frequencies of drought/floods. Reliability diagrams will be used here to assess
17 the reliability and confidence of the forecasts. ROC applied to probabilistic
18 forecasts indicates whether the forecast probability was higher when an event
19 such as drought occurred compared to when it did not occur, and therefore
20 identifies whether a set of forecasts has the attribute of discrimination. Here the
21 area underneath the ROC curve is used as a measure of discrimination in the
22 prediction of below-normal and above-normal DJF rainfall totals.

23

1 **3. Results**

2

3 3.1. Deterministic assessment of forecasts

4

5 Although the seasonal climate is inherently probabilistic and therefore seasonal
6 forecasts globally are for the most part issued probabilistically, it is often
7 informative to investigate deterministic forecast performance. Figure 1 shows
8 area-averaged (Africa south of 10°S) deterministic cross-validated (5-year-out
9 approach) multi-model DJF rainfall (mm) forecasts over the available 43-year
10 period (1959/60 – 2001/02) compared with the observed. The cross-validation
11 procedure is designed in such a way that the data is “wrapped” around in order to
12 make a 5-year-out approach possible while at the same time producing cross-
13 validated forecasts for the whole period. Forecasts for both MMcca and MMeqw
14 are shown, and El Niño and La Niña seasons are respectively marked with “E”
15 and “L”. The vertical line on the figure divides the time series into two parts: The
16 initial training period for the creation of retro-active forecasts (1959/60 – 1979/80;
17 21 years) and the retro-active test period (1980/81 – 2001/02; 22 years) for which
18 probabilistic forecasts are generated. The Spearman’s correlation between the
19 area-averaged 22-year forecasts and observations for MMcca and MMeqw are
20 respectively 0.4783 and 0.4873, suggesting about equally skilful area-averaged
21 deterministic forecasts from the two multi-model methods. The Spearman’s
22 correlation is used here since the 1997/98 rainfall predictions are considered
23 outliers (Figure 1). The four driest years during the 22-year test period (1982/83,

1 1986/87, 1991/92, 1994/95) are associated with El Niño seasons and the four
2 wettest with La Niña seasons (1988/89, 1995/96, 1998/99, 1999/00). For the
3 most part, the forecasts do not capture the size of the observed anomalies for
4 these extreme seasons, but this is often found with linear regression-based
5 downscaling techniques such as the one used here. Notwithstanding, no attempt
6 was made here to inflate the forecasts since variance adjustment of forecasts are
7 generally discouraged (Troccoli *et al.*, 2008).

8

9 The length of the training period may have an effect on the robustness or stability
10 of the MOS equations (Doblas-Reyes *et al.*, 2005; Wilks, 2006). For stability it is
11 understood that the fitted equations are also applicable to independent data.
12 Since the initial training period (for making the 1980/81 rainfall forecasts) is only
13 21 years long, investigation into the variation of forecast performance over the
14 various training periods is warranted. Figure 2 shows area-averaged Spearman's
15 correlations (adjusted with the Fisher Z transformation (Wilks, 2006)) for various
16 cross-validation training periods ranging from 12 years to 43 years, using MMcca,
17 and using August-September-October averaged SSTs as predictor in a statistical
18 model (CCA). The SST predictor field is between 170°E to 80°W and 20°N to
19 20°S in order to capture central and eastern equatorial Pacific SST variability. A
20 4th order polynomial is fitted to the averaged Spearman's correlations and a
21 gradual improvement in forecast skill can be seen towards a training set
22 consisting of 32 years when MMcca is used, and throughout the whole period
23 when using SSTs as predictor in the statistical model. A skill plateau could have

1 been attained with the MMcca were it not for the large errors associated with the
2 rainfall prediction of the 1997/98 El Niño season and of the two preceding years.
3 Thereafter a gradual decrease is seen until 43 years are included in the MOS
4 training period. Using the DJF 850 hPa geopotential field predicted at the end of
5 October by the coupled ECHAM4.5-MOM3-DC2
6 ([http://iridl.ldeo.columbia.edu/SOURCES/.IRI/.MP/.RESEARCH/.COUPLED/.GL
8 OBAL/.ECHAM4p5-MOM3-DC2/](http://iridl.ldeo.columbia.edu/SOURCES/.IRI/.MP/.RESEARCH/.COUPLED/.GL
7 OBAL/.ECHAM4p5-MOM3-DC2/)) as predictor in the same MOS downscaling
9 approach for southern Africa, a similar shape is found in the variation of skill
10 (Figure 2). Here the initial training period is from 1982/83 to 1991/92. It is
11 suggested that the decrease in skill towards the 2001/02 season is therefore not
12 a function of the DEMETER data used here, since a differently configured
13 coupled model produces similar results. Forecast skill using physical models may
14 thus not be constant in time. However, the dominant modes of CCA (Barnett and
15 Preisendorfer, 1987) for the multi-model considered here remain the same (not
16 shown) regardless of the training period used (e.g. Landman and Goddard,
17 2002), which suggests stability in the selected dominant modes of variability
18 included in the MOS equations, and therefore implies stability in the MOS
19 prediction equations even though forecast skill may not be constant in time.

19

20 3.2. Multi-model vs. single model results

21

22 By knowing the probability of a predicted category occurring, additional forecast
23 value is obtained (Mason and Graham, 1999), since probabilistic forecasts exhibit

1 reliability considerably in excess of that achieved by corresponding deterministic
2 forecasts (Murphy, 1998). Probabilistic rainfall forecasts are produced here for
3 three equi-probable categories of above-normal, near-normal and below-normal.
4 Only the verification results for the above- and below-normal categories are
5 presented here since there is little skill to be derived from predicting the near-
6 normal category (Van den Dool and Toth 1991).

7

8 A ROC graph is made by plotting the forecast hit rates against the false alarm
9 rates (Wilks, 2006). The area beneath the ROC curve is used as a measure of
10 discrimination here and is referred to as a ROC score. If the area would be ≤ 0.5
11 the forecasts have no skill, and for a maximum ROC score of 1.0, perfect
12 discrimination has been obtained. The ROC score can be interpreted here as a
13 probability of the forecast system successfully discriminating respectively above-
14 or below-normal seasons from other seasons.

15

16 The ROC graph and its score can be meaningfully applied in seasonal
17 forecasting given the small sample size normally associated with these forecasts
18 (Troccoli *et al.*, 2008). Figure 3 shows the area-averaged ROC scores for above-
19 and below-normal DJF rainfall for each of the individual downscaled models
20 (Météo-France – MF; ECMWF and UKMO) and for the two multi-models (MMeqw
21 and MMcca) as calculated over the 22-year test period in a retro-active design.
22 All area-averaged scores are above 0.5, which means that on average there is
23 more than a 50% chance that all the forecast systems have the ability to

1 successfully discriminate respectively wet and dry seasons from other seasons.
2 Two of the three single models have a greater ability to discriminate the below-
3 normal category as opposed to the above-normal one, but both the multi-models
4 are better able to discriminate the below-normal category. Moreover, the multi-
5 models have higher averaged ROC scores than any of the individual models. In
6 fact, based on the area-averaged scores the multi-models each have at least a
7 61% chance of discriminating the above-normal category and at least a 63%
8 chance of discriminating the below-normal DJF rainfall. The outperformance by
9 the multi-models over southern Africa confirms what has been found with many
10 other studies that multi-model forecasts usually outscore single model forecasts
11 (e.g. Barnston *et al.*, 2003; Doblas-Reyes *et al.*, 2005; Hagedorn *et al.*, 2005;
12 Coelho *et al.*, 2006; Weigel *et al.*, 2008; Wang and Fan, 2009).

13

14 The improvement in forecast performance of the multi-models over the single
15 models is further demonstrated in Figure 4 that shows the geographical
16 distribution of ROC score differences between the multi-models and the
17 individual models. Figure 4(a) shows where the multi-model that uses equal
18 weights (MMeqw) outscore each of the individual models, and Figure 4(b) where
19 the weighted forecast combination multi-model (MMcca) outscores them. Shaded
20 areas are where the multi-models outperform the single models. Both sets of
21 maps show that most of southern Africa is associated with positive ROC score
22 differences, thus providing further evidence that the multi-models are outscoring
23 the single models.

1

2 The ROC score is sometimes criticized as a measure of forecast performance
3 because of its insensitivity to reliability (Troccoli *et al.*, 2008). Figure 5 shows the
4 reliability diagrams for the individual models. In addition to the respective
5 reliability curves for the two categories, their least-squares regression lines are
6 presented on the diagrams. The regression lines are calculated with weighting
7 relative to how frequently forecasts are issued at a given confidence. When these
8 regression lines lie along the diagonal, the forecasts are perfectly reliable. When
9 the regression line lies above the diagonal observed above- or below-normal DJF
10 rainfall tends to occur more frequently than forecast, but when it lies below the
11 diagonal the observed categories respectively tend to occur less frequently than
12 forecast, indicating under- and over-forecasting respectively. The most common
13 slope of the regression line found for seasonal forecasting is one that is shallower
14 than the diagonal line (Troccoli, *et al.*, 2008) – the forecasts are said to be over-
15 confident. Histograms are also included in the figures, and they show the
16 frequencies with which forecasts occur in probability intervals of 10%, starting at
17 5%.

18

19 All the forecasts made by the single models for both above- and below normal
20 DJF are over-confident (Figure 5). However, forecasts for below-normal rainfall
21 totals are less over-confident than forecasts for above-normal rainfall for all three
22 single models. Since the single models are over-confident, multi-model
23 ensembles can enhance prediction skill regardless of which combination

1 approach is used since multi-model combination reduces over-confidence
2 (Weigel *et al.*, 2008). Figure 6 shows the reliability diagrams of the two multi-
3 models, and here improved reliability over the single models is in fact seen (the
4 regression lines for both categories tend to be closer to the diagonal). However,
5 for both multi-models the high-probability above-normal forecasts are not reliable,
6 as well as the high-probability below-normal forecasts of the MMeqw model. This
7 result suggests that a simple equal weighting scheme to combine forecasts may
8 not sufficiently reduce over-confidence (Barnston *et al.*, 2003) for high-probability
9 forecasts. Difference maps (not shown) of ROC scores (MMcca minus MMeqw)
10 for the two categories show more or less an even split in terms of the areas of
11 positive and negative score differences. This result indicates that both multi-
12 model approaches are not much different in their ability to discriminate events
13 from non-events, and that the MMcca is only slightly better able to produce
14 reliable high-probability below-normal rainfall forecasts. However, such forecasts
15 are often made during El Niño seasons

16

17 It has been shown that both the single and multi-models have the ability to
18 discriminate between different observed situations. However, the multi-models
19 outscore the single models, both in terms of discrimination and reliability. Since
20 southern African mid-summer rainfall is influenced by the state of the equatorial
21 Pacific Ocean, there is a need to investigate how skilful a multi-model predicts
22 the two rainfall categories during ENSO and during neutral events separately.

23

1 3.3. Multi-model forecast performance during ENSO years

2

3 CCA pattern and time series analysis (Barnett and Preisendorfer, 1987) of the
4 multi-model (MMcca) forecast system suggests that the dominant modes of
5 predictor variability (three or four canonical modes that produce the best forecast
6 results over the retro-active forecast period) are partly related to different
7 influences of ENSO on southern African mid-summer rainfall (Fauchereau *et al.*,
8 2008) since the correlations between the Oceanic Niño Index (ONI;
9 www.cpc.noaa.gov/products/analysis_monitoring/ensostuff/ensoyears.shtml) and
10 the three leading canonical temporal scores of the predictor (combined 850 hPa
11 geopotential height fields) are respectively 0.5017 ($p < 0.01$), -0.5337 ($p < 0.01$) and
12 -0.3023 ($p < 0.05$) over the 43-year period. The question may arise then what
13 added benefit there may be in running multi-model systems that consist of
14 physical models that are primarily ENSO driven, over a simple statistical model
15 that uses Pacific Ocean SSTs as predictors and is much cheaper to run. This
16 question is answered by referring back to Figure 2. The gray dashed line is the
17 4th order polynomial that is fitted to the area-averaged Spearman's correlation
18 obtained by using a simple statistical model (CCA) with central and eastern
19 equatorial Pacific Ocean SST (170° E to 80° W; 20° N to 20° S) as predictor.
20 Although there is convergence in the performance of the forecasting systems
21 towards the end of the cross-validation period, the multi-model outscores the
22 simple model throughout. This result suggests that the coupled models'
23 downscaled forecasts include additional forecast information that cannot be

1 derived from equatorial Pacific SST alone, which justifies the use of physical
2 forecast models to predict seasonal rainfall variability over southern Africa. Take
3 note that the introduction here of the statistical model was not to set an easy to
4 beat baseline skill level, but to demonstrate that the skill of the GCMs comes
5 from climatological forcings beyond the central and eastern equatorial Pacific
6 Ocean.

7

8 The multi-model DJF rainfall forecast performance during the El Niño (8
9 seasons), La Niña (7 seasons) and neutral (7 seasons) years over the 22-year
10 retro-active period are shown in Figure 7 to 9. The forecasts for the ENSO and
11 non-ENSO years are separately taken from the retro-active forecasts prior to
12 calculating the verification statistics for these years. Since the skill calculations
13 are based on only a few cases (7 or 8) they may be sensitive to sampling errors.
14 ROC calculations are however less sensitive to sampling errors than reliability
15 diagrams (Troccoli *et al.*, 2008). Figure 7 presents area-averaged ROC scores
16 and it is shown that on average the multi-model is able to discriminate the above-
17 normal and below-normal rainfall categories during ENSO years, but fails to do
18 so during neutral years (averaged ROC scores are below 0.5 for both
19 categories). Moreover, the multi-model performs best predicting drought during El
20 Niño years and floods during La Niña years, but there is skill in predicting wet El
21 Niño and dry La Niña seasons over southern Africa too. This result is further
22 manifested in the geographical distribution of ROC scores for the above- and
23 below-normal rainfall categories and for ENSO and neutral years as shown in

1 Figure 8. Large patterns of ROC scores in excess of 0.5 are seen for the El Niño
2 and La Niña cases, but much smaller areas associated with neutral years are
3 found. The multi-model therefore performs poorly during neutral years. The
4 reliability diagrams for rainfall prediction during El Niño and La Niña years are
5 shown in Figure 9. Forecasts are again over-confident, but as is found with the
6 ROC scores there is skill in predicting both drought and wet seasons during El
7 Niño years and predicting wet and drought seasons during La Niña years. The
8 forecasts at least correctly indicate increases and decreases in the probabilities
9 of the wet and dry events.

10

11 **4. Discussion and conclusions**

12

13 Southern African mid-summer probabilistic rainfall prediction skill has been
14 assessed by using forecasts from state-of-the-art fully coupled models that are
15 empirically downscaled and combined in order to produce multi-model forecasts.
16 Forecast performance was tested over a retro-active period of 22 years that
17 mimics an operational forecast configuration. Multi-model forecasts outscore
18 single model forecasts and can be used with confidence during El Niño and La
19 Niña seasons. In addition, the two multi-model forecast approaches produce
20 about equally skilful forecasts.

21

22 The robustness of the MOS equations was tested and found that although
23 forecast skill may not be constant in time, especially with short training periods,

1 the dominant modes of variability included in the equations remain similar for a
2 variety of training periods. Regardless of this variation in skill, multi-model
3 performance consistently outscored a simple statistical model that only includes
4 equatorial Pacific Ocean SST variability as predictor. The improved multi-model
5 forecasts are therefore a result of the system's ability to include forecast
6 information in addition to the signal originating from the central and eastern
7 equatorial Pacific Ocean. Both single model downscaled forecasts and multi-
8 model forecasts seems to be able to discriminate between different observed
9 situations such as below-normal and above-normal DJF rainfall seasons,
10 notwithstanding the result that forecasts are overconfident. Prediction of wet or
11 dry conditions during ENSO years is also skilful, but little skill has been found
12 predicting DJF rainfall when the equatorial Pacific Ocean is in a neutral state.
13 Predictions during El Niño seasons are strongly overconfident, but are less so for
14 rainfall predictions during La Niña seasons.

15

16 The paper has demonstrated that multi-model systems are able to provide useful
17 operational mid-summer rainfall forecasts over southern Africa, but only during
18 ENSO years. Rainfall forecasts for southern Africa produced by the EUROSIP
19 multi-model, that consists of later versions of the three coupled GCMs discussed
20 here, made near the end of 2009 for the 2009/10 DJF El Niño season show
21 mostly enhanced probabilities for dry conditions to occur. A similar forecast was
22 also issued by other international centres such as the IRI, and also by the South
23 African Weather Service. Moreover, summer rainfall forecasts for 2009/10 issued

1 to the South African public was made with high confidence, partly based on the
2 result that multi-models can produce reliable drought forecasts and because of
3 the confidence in summer rainfall forecasts during El Niño seasons. However,
4 DJF rainfall over South Africa was anomalously high, especially over the central
5 and western parts of that country (<http://www.weathersa.co.za>) and so the
6 observed wet 2009/10 austral summer season over the region was largely
7 missed by most forecasting systems. Further model development (e.g.
8 Engelbrecht *et al.*, 2007) and modelling studies on how models represent the
9 coupled system over southern Africa are therefore warranted.

10

11 **Acknowledgements**

12

13 This work was partly sponsored by the Water Research Commission of South
14 Africa (project K5/1492).

15

16 **References**

17

18 Barnett TP, Preisendorfer RW. 1987. Origins and levels of monthly and seasonal
19 forecast skill for United States air temperature determined by canonical
20 correlation analysis. *Monthly Weather Review* **115** : 1825-1850.

21

1 Barnston AG, Mason SJ, Goddard L, DeWitt DG, Zebiak SE. 2003. Multimodel
2 ensembling in seasonal climate forecasting at IRI. *Bulletin of the American*
3 *Meteorological Society*, 1783-1796. DOI: 10.1175/BAMS-84-12-1783.
4
5 Bartman AG, Landman WA, Rautenbach CJ deW. 2003. Recalibration of general
6 circulation model output to austral summer rainfall over southern Africa.
7 *International Journal of Climatology* **23**: 1407-1419.
8
9 Branković Č, Palmer TN 2000. Seasonal skill and predictability of ECMWF
10 PROVOST ensembles. *Quarterly Journal of the Royal Meteorological Society*
11 **126**: 2035-2067.
12
13 Brown BH, Murphy AH. 1996. Improving forecasting performance by combining
14 forecasts: the example of road-surface temperature forecasts. *Meteorological*
15 *Applications* **3**: 257-265.
16
17 Coelho CAS, Stephenson BD, Balmaseda M, Doblas-Reyes FJ, van Oldenborgh
18 GJ. 2006. Toward an integrated seasonal forecasting system for South
19 America. *Journal of Climate* **19**: 3704-3721.
20
21 Copsey D, Sutton R, Knight JR. 2006. Recent trends in sea level pressure in the
22 Indian Ocean region. *Geophysical Research Letters* **33**: L19712,
23 doi:10.1029/2006GL027175.

1

2 Doblas-Reyes FJ, Déqué M, Piedelieve J-P. 2000. Multi-model spread and
3 probabilistic seasonal forecasts in PROVOST. *Quarterly Journal of the Royal*
4 *Meteorological Society* **126**: 2035-2067.

5

6 Doblas-Reyes FJ, Hagedorn R, Palmer TN. 2005. The rationale behind the
7 success of multi-model ensembles in seasonal forecasting – II. Calibration and
8 combination. *Tellus* **57A**: 234-252.

9

10 Dirmeyer PA, Fennessy MJ, Marx L. 2003. Low skill in dynamical prediction of
11 boreal summer climate: Grounds for looking beyond sea surface temperature.
12 *Journal of Climate* **16**: 995-1002.

13

14 Engelbrecht FA, McGregor JL, Rautenbach CJdeW. 2007. On the development
15 of a new nonhydrostatic atmospheric model in South Africa. *South African*
16 *Journal of Science* **103**: 127-134.

17

18 Evans RE, Harrison MSJ, Graham RJ, Mylne KR. 2000. Joint medium-range
19 ensembles from the Met Office and ECMWF systems. *Monthly Weather Review*
20 **128**: 3104-3127.

21

22 Fauchereau N, Pohl B, Reason CJC, Rouault M, Richard Y. 2008. Recurrent
23 daily OLR patterns in the southern African/southwest Indian Ocean region,

1 implications for South African rainfall and teleconnection. *Climate Dynamics*.
2 DOI:10.1007/s00382-008-0426-2.

3

4 Frías MD, Herrera S, Cofiño AS, Gutiérrez JM. 2010. Assessing the skill of
5 precipitation and temperature seasonal forecasts in Spain: Windows of
6 opportunity related to ENSO events. *Journal of Climate* **23**: 209-220.

7

8 Goddard L, Mason SJ. 2002. Sensitivity of seasonal climate forecasts to
9 persisted SST anomalies. *Climate Dynamics* **19**: 619-631.

10

11 Graham RJ, Evans ADL, Mylne KR, Harrison MSJ, Robertson KB. 2000. An
12 assessment of seasonal predictability using atmospheric general circulation
13 models. *Quarterly Journal of the Royal Meteorological Society* **126** : 2211-2240.

14

15 Hagedorn R, Doblas-Reyes FJ, Palmer TN. 2005. The rationale behind the
16 success of multi-model ensembles in seasonal forecasting – I. Basic concept.
17 *Tellus* **57A**: 219-232.

18

19 Hamill TM. 1997. Reliability diagrams for multicategory probabilistic forecasts.
20 *Weather and Forecasting* **12** : 736-741.

21

1 Joubert AM, Hewitson BC. 1997. Simulating present and future climates of
2 southern Africa using general circulation models. *Progress in Physical*
3 *Geography* **21**: 51-78.

4

5 Krishnamurti TN, Kishtawal CM, Zang Z, LaRow T, Bachiochi D, Williford E,
6 Gadgil S, Surendran S. 2000. Multimodel ensemble forecasts for weather and
7 seasonal climate. *Journal of Climate* **13**: 4196-4216.

8

9 Landman WA, Goddard L. 2002. Statistical recalibration of GCM forecasts over
10 southern Africa using model output statistics. *Journal of Climate* **15**: 2038-2055.

11

12 Landman WA, Goddard L. 2005. Predicting southern African summer rainfall
13 using a combination of MOS and perfect prognosis. *Geophysical Research*
14 *Letters* **32**: L15809. DOI: 10.1029/2005GL022910.

15

16 Landman WA, Mason SJ, Tyson PD, Tennant WJ. 2001. Retro-active skill of
17 multi-tiered forecasts of summer rainfall over southern Africa. *International*
18 *Journal of Climatology* **21**: 1-19.

19

20 Mason SJ. 1995. Sea-surface temperature – South African rainfall associations,
21 1910-1989. *International Journal of Climatology* **15**: 119-135.

22

1 Mason SJ, Graham NE. 2002. Areas beneath the relative operating
2 characteristics (ROC) and levels (ROL) curves: Statistical significance and
3 interpretation. *Quarterly Journal of the Royal Meteorological Society* **128**: 2145-
4 2166.

5

6 Mason SJ, Joubert AM. 1997. Simulated changes in extreme rainfall over
7 southern Africa. *International Journal of Climatology* **17**: 291-301.

8

9 Mason SJ, Mimmack GM. 2002. Comparison of some statistical methods of
10 probabilistic forecasting of ENSO. *Journal of Climate* **15**: 8-29.

11

12 Mason SJ, Goddard L, Graham NE, Yulaeva E, Sun L, Arkin PA. 1999. The IRI
13 seasonal climate prediction system and the 1997/98 El Niño event. *Bulletin of the*
14 *American Meteorological Society* **80**: 1853-1873.

15

16 Mitchell TD, Jones PD. 2005. An improved method of constructing a database of
17 monthly climate observations and associated high-resolution grids. *International*
18 *Journal of Climatology* **25**: 693-712. DOI: 10.1002/joc.1181

19

20 Murphy AH. 1998. The early history of probability forecasts: Some extensions
21 and clarification. *Weather and Forecasting* **13**: 5-15.

22

1 Palmer TN, Anderson DLT. 1994. The prospects of seasonal forecasting – a
2 review paper. *Quarterly Journal of the Royal Meteorological Society* **120**: 755-
3 793.

4

5 Palmer TN, Coauthors. 2004. Development of a European multimodel ensemble
6 system for seasonal-to-interannual prediction (DEMETER). *Bulletin of the*
7 *American Meteorological Society*, DOI: 10.1175/BAMS-85-6-853.

8

9 Peng PT, Kumar A, van den Dool H and Barnston AG. 2002. An analysis of multi-
10 model ensemble predictions for seasonal climate anomalies. *Journal of*
11 *Geophysical Research*, **107**.

12

13 Reason CJC, Landman W, Tennant W. 2006. Seasonal to decadal prediction of
14 southern African climate and its links with variability of the Atlantic Ocean,
15 *Bulletin of the American Meteorological Society* : DOI:10.1175/BAMS-87-7-941.

16

17 Saha S, and Coauthors. 2006. The NCEP climate forecast system. *Journal of*
18 *Climate* **19**: 3483-3517.

19

20 Shongwe ME, Landman WA, Mason SJ. 2006. Performance of recalibration
21 systems for GCM forecasts for southern Africa. *International Journal of*
22 *Climatology* **26**: 1567-1585.

23

1 Stockdale TN, Anderson DLT, Alves JOS, Balmaseda MA. 1998. Global
2 seasonal rainfall forecasts using a coupled ocean-atmosphere model. *Nature*
3 **392**: 370-373.

4

5 Tippet MK and Barnston AG. 2008. Skill of multimodel ENSO probability
6 forecasts. *Monthly Weather Review*, **136**, 3933-3946.

7

8 Tippet MK, Goddard L, Barnston AG. 2005. Statistical-dynamical seasonal
9 forecasts of central-southwest Asian winter precipitation. *Journal of Climate* **18**:
10 1831-1843.

11

12 Tracton MS, Kalnay E. 1993. Operational ensemble prediction at the National
13 Meteorological Center: Practical aspects. *Weather and Forecasting* **8**: 379-398.

14

15 Troccoli A, Harrison M, Anderson DLT, Mason SJ. 2008. *Seasonal Climate:*
16 *Forecasting and managing risk*. NATO Science Series. Earth and Environmental
17 Sciences Vol 82. Springer.

18

19 Tyson PD, Preston-Whyte RA. 2000. *The Weather and Climate of Southern*
20 *Africa*. Oxford University Press.

21

22 Van den Dool HM and Toth Z. 1991. Why do forecasts for near normal often fail?
23 *Weather and Forecasting*, **6**, 76-85.

1

2 Wang H, Fan K. 2009. A new scheme for improving the seasonal prediction of
3 summer precipitation anomalies. *Weather and Forecasting* **34**: 548-554. DOI:
4 10.1175/2008WAF2222171.1.

5

6 Washington R, Preston A. 2006. Extreme wet years over southern Africa: Role of
7 the Indian Ocean sea surface temperatures. *Journal of Geophysical Research*
8 **111** : D15104, DOI: 10.1029/2005JD006724.

9

10 Weigel AP, Liniger MA, Appenzeller C. 2008. Can multi-model combination really
11 enhance the prediction skill of probabilistic ensemble forecasts? *Quarterly*
12 *Journal of the Royal Meteorological Society* **134**: 241-260.

13

14 Weisheimer A, Dobals-Reyes FJ, Palmer TN, Alessandri A, Arribas A, Déqué M,
15 Keenlyside N, MacVean M, Navarra A, Rogel P. 2009. ENSEMBLES: A new
16 multi-model ensemble for seasonal-to-annual predictions – Skill and progress
17 beyond DEMETER in forecasting tropical Pacific SSTs. *Geophysical Research*
18 *Letters* **36**: L21711, doi:10.1029/2009GL040896.

19

20 Wilby RL, Wigley TML. 1997. Downsclaing general circulation model output: A
21 review of methods and limitations. *Progress in Physical Geography* **21**: 530-548.

22

1 Wilks DS. 2006. *Statistical Methods in the Atmospheric Sciences, 2nd Edition*.
2 Academic Press.

3

4 **Figure captions**

5

6 Figure 1. Area-averaged observed (thick line) DJF rainfall (mm) over Africa south
7 of 10° S, versus cross-validation forecasts (thin lines) from the two multi-models
8 described in the text. El Niño (E) and La Niña (L) seasons are also shown. The
9 arrow indicates where the retro-active test period starts. The years on the x-axis
10 refer to the December months of the DJF seasons.

11

12 Figure 2. Variation in cross-validation forecast skill predicting DJF rainfall over
13 southern Africa as reflected by area-averaged Spearman's correlation values.
14 The thick black solid line (4th order polynomial) and associated thin black solid
15 line show the MMcca multi-model's performance as a function of cross-validation
16 training period, while the thick black dotted and thin black dotted lines represent
17 the ECHAM4.5-MOM3-DC2 coupled model. The remaining gray lines represent
18 the statistical model that uses equatorial Pacific Ocean SST as predictor. The
19 arrow indicates where the retro-active test period starts.

20

21 Figure 3. ROC scores, averaged over the southern African domain, for the
22 above-normal and below-normal rainfall categories. Scores for the single models
23 and for the two multi-models are shown.

1

2 Figure 4. ROC score differences between the a) MMeqw multi-model and the
3 single models, and b) MMcca multi-model and single models. Positive ROC
4 score differences are where the multi-models are superior.

5

6 Figure 5. Reliability diagrams and frequency histograms for above- and below-
7 normal DJF rainfall forecasts produced by the single models. The thick black
8 curves and black bars of the histogram represent the below-normal rainfall
9 category, while the thick black dotted curves and white bars of the histogram
10 represent the above-normal rainfall category. For perfect reliability the curves
11 should fall on top of the thick black diagonal line. The thin solid and dotted lines
12 are respectively the weighted least-squares regression lines of the above-normal
13 and below-normal reliability curves.

14

15 Figure 6. As in Figure 5, but for the two multi-models.

16

17 Figure 7. ROC scores, averaged over the southern African domain, for the
18 above-normal and below-normal rainfall categories during El Niño, La Niña and
19 neutral seasons. Scores for the MMcca multi-model are shown.

20

21 Figure 8. ROC scores of the MMcca multi-model, for El Niño, La Niña and neutral
22 seasons, and for the above- and below-normal rainfall categories. ROC scores
23 ≥ 0.5 are shaded.

1

2 Figure 9. As in Figure 5, but for rainfall predictions during El Niño and La Niña

3 seasons using the MMcca multi-model.

4