

We thank the reviewer for their careful reading of our revised manuscript and we have responded (in blue) to their comments (in black).

The authors responded well to the reviewer's comments and considerably improved their paper. Consequently I suggest accepting the paper for publication. I only have a few minor suggestions for corrections.

P. 7, line 20-25

The content is largely redundant with previous text, starting at p.5 - line 25, to p.6 line 10. I guess it is sufficient to state that a second set of metrics is presented in Table 2b that comprise some of those metrics introduced before.

As suggested we have removed much of this material and joined the remainder to the previous paragraph.

p. 7 - line 28

... hydrologists

Corrected.

p.8 - line 5, Change to:

Table 2a indicates that only

Modified as suggested.

p. 13 - line 7:

... analysis (reported ...

Corrected.

Table 3

Why many numbers have a preceding '=' sign in front of the number?

I suggest removing the '='.

We have added a # symbol to the table and explain in the Table footnote the meaning of the = symbol.

1 **Assessment of precipitation and temperature data from** 2 **CMIP3 Global Climate Models for hydrologic simulation**

3
4 **T. A. McMahon¹, M. C. Peel¹ and D. J. Karoly²**

5 [1]{Department of Infrastructure Engineering, University of Melbourne, Victoria, 3010,
6 Australia}

7 [2]{School of Earth Sciences and ARC Centre of Excellence for Climate System Science,
8 University of Melbourne, Victoria, 3010, Australia}

9 Correspondence to: M. C. Peel (mpeel@unimelb.edu.au)

10 11 **Abstract**

12 The objective of this paper is identify better performing CMIP3 Global Climate Models
13 (GCMs) that reproduce grid-scale climatological statistics of observed precipitation and
14 temperature as input to hydrologic simulation over global land regions. Current assessments
15 are aimed mainly at examining the performance of GCMs from a climatology perspective and
16 not from a hydrology standpoint. The performance of each GCM in reproducing the
17 precipitation and temperature statistics was ranked and better performing GCMs identified for
18 later analyses. Observed global land surface precipitation and temperature data were drawn
19 from the CRU 3.10 gridded dataset and re-sampled to the resolution of each GCM for
20 comparison. Observed and GCM based estimates of mean and standard deviation of annual
21 precipitation, mean annual temperature, mean monthly precipitation and temperature and
22 Köppen-Geiger climate type were compared. The main metrics for assessing GCM
23 performance were the Nash-Sutcliffe efficiency index and RMSE between modelled and
24 observed long-term statistics. This information combined with a literature review of the
25 performance of the CMIP3 models identified the following better performing GCMs from a
26 hydrological perspective: HadCM3 (Hadley Centre for Climate Prediction and Research),
27 MIROCM (Center for Climate System Research (The University of Tokyo), National Institute
28 for Environmental Studies, and Frontier Research Center for Global Change), MIUB
29 (Meteorological Institute of the University of Bonn, Meteorological Research Institute of

1 KMA, and Model and Data group), MPI (Max Planck Institute for Meteorology) and MRI
2 (Japan Meteorological Research Institute). The future response of these GCMs was found to
3 be representative of the 44 GCM ensemble members which confirms that the selected GCMs
4 are reasonably representative of the range of future GCM projections.

5

6 **1 Introduction**

7 Our primary objective in this paper is to identify better performing GCMs from a hydrologic
8 perspective. To do this we assess how well 22 Global Climate Models (GCMs) from the
9 World Climate Research Programme's (WCRP's) Coupled Model Intercomparison Project
10 phase 3 (CMIP3) multi-model dataset (Meehl et al., 2007) are able to reproduce GCM grid-
11 scale climatological statistics of observed precipitation and temperature over global land
12 regions. We recognise that GCMs model different variables with a range of success and that
13 no single model is best for all variables and/or for all regions (Lambert and Boer, 2001;
14 Gleckler et al., 2008). The approach adopted here is not inconsistent with Dessai et al. (2005)
15 who regarded the first step in evaluating GCM projection skill is to assess how well observed
16 climatology is simulated. We also recognize there have been assessments published in peer-
17 reviewed journals, but all appear to be assessed from a climate science perspective. This
18 review concentrates on GCM variables and statistical techniques that are relevant to
19 engineering hydrologic practice.

20

21 GCM runs for the observed period do not seek to replicate the observed monthly record at any
22 point in time and space. Rather a better performing GCM is expected to produce long-term
23 mean annual statistics that are broadly similar to observed conditions across a wide range of
24 locations. Here, the assessment of CMIP3 GCMs is made by comparing their long-term mean
25 annual precipitation (MAP), standard deviation of annual precipitation (SDP), mean annual
26 temperature (MAT), mean monthly patterns of precipitation and temperature, and Köppen-
27 Geiger climate type (Peel et al., 2007) with concurrent observed data for 616 to 11886
28 terrestrial grid cells world-wide (the number of grid cells depends on the resolution of the
29 GCM under consideration). These variables were chosen to assess GCM performance because
30 they provide insight into the mean annual, inter-annual variability and seasonality of
31 precipitation and temperature, which are sufficient to estimate the mean and variability of
32 annual runoff from a traditional monthly rainfall-runoff model (Chiew and McMahon, 2002)

1 or from a top-down annual rainfall-runoff model (McMahon et al., 2011) for hydrologic
2 simulation purposes.

3

4 The GCMs included in this assessment are detailed in Table 1. (Model acronyms adopted are
5 listed in the table.) Although no quantitative assessment of the BCCR model is made, this
6 model is included in Table 1 as details of its performance are available in the literature which
7 is discussed in Section 2. Other details in the table include the originating group for model
8 development, country of origin, model name given in the CMIP3 documentation (Meehl et
9 al., 2007), the number of 20C3M runs available for analysis, the model resolution, and the
10 number of terrestrial grid cells used in the precipitation and temperature comparisons.

11

12 Readers should note that when this project began as a component of a larger study in 2010,
13 runs from the Coupled Model Intercomparison Project phase 5 (CMIP5) were not available.
14 We are of the view that the approach adopted here is equally applicable to evaluating CMIP5
15 runs for hydrologic simulations. Conclusions about better performing models drawn from this
16 analysis may prove similar to a comparable analysis of CMIP5 runs since most models in
17 CMIP5 are, according to Knutti et al. (2013), “strongly tied to their predecessors”. Analysis of
18 the CMIP5 models indicates that the CMIP3 simulations are of comparable quality to the
19 CMIP5 simulations for temperature and precipitation at regional scales (Flato et al., 2013).

20

21 This study is part of a larger research project that seeks to enhance our understanding of the
22 uncertainty of future annual river flows world-wide through catchment scale hydrologic
23 simulation, leading to more informed decision-making for the sustainable management of
24 scarce water resources, nationally and internationally. To achieve this, it is necessary to
25 determine, as a minimum, how the mean and variability of annual streamflows will be
26 affected by climate change. Other factors of less importance are changes in the auto-
27 correlation of annual streamflow, changes in net evaporation from reservoir water surfaces,
28 and changes in monthly flow patterns, with the latter being more important for relatively
29 small reservoirs. In this paper we deal with the key drivers of streamflow production namely
30 the mean and the standard deviation of annual precipitation and mean annual temperature, the
31 latter is adopted here as a surrogate for potential evapotranspiration (PET), along with

1 secondary factors, the mean monthly patterns of precipitation and temperature. Adopting
2 temperature as a surrogate for PET is contentious. We provide a detailed discussion of this
3 issue in the Supplementary Material associated with this paper. Suffice to say that a more
4 complex PET formulation requires additional GCM variables other than temperature which
5 are less reliable. This simplicity comes at the expense of potentially inadequate representation
6 of future changes in PET, which may have important negative consequences when modelling
7 streamflow in energy limited catchments. Nevertheless, in the following discussion we
8 concentrate on mean annual temperature as the GCM variable representing PET.

9

10 Computer models of most water resource systems that rely on surface reservoirs to offset
11 streamflow variability adopt a monthly time-step to ensure that seasonal patterns in demand
12 and reservoir inflows are adequately accounted for. However, in a climate change scenario it
13 is more likely that an absolute change in streamflow will have a greater impact on system
14 yield than shifts in the monthly inflow or demand patterns. This will certainly be the case for
15 reservoirs that operate as carryover systems rather than as within-year systems (for an
16 explanation see McMahon and Adeloze (2005)). Therefore, in this paper we assess the GCMs
17 in terms of annual precipitation and annual temperature, and patterns of mean monthly
18 precipitation and temperature.

19

20 Following this introduction we describe, and summarise in the next section, several previous
21 assessments of CMIP3 GCM performance. We also include some general comments on GCM
22 assessment procedures. In Section 3, data (observed and GCM based) used in the analysis are
23 described. Details and results of the subsequent analyses comparing GCM estimates of
24 present climate mean and standard deviation of annual precipitation, mean annual
25 temperature, mean monthly precipitation and temperature patterns and Köppen-Geiger
26 climate type against observed data are set out in Section 4. In Section 5, we review the results
27 and compare the literature information with our assessments of the GCMs. The final section
28 of the paper presents several conclusions.

29

1 **2 Literature**

2 As noted above, to assess the impact of climate change on surface water resources of a region
3 through hydrologic simulation, it is necessary to assess, as a minimum, the performance of the
4 mean and the standard deviation of annual precipitation and mean annual temperature, and the
5 mean monthly patterns of precipitation and temperature. Noting this background we describe
6 in the next section procedures that have been adopted in the literature to assess GCM
7 performance.

8 **2.1 Procedures to assess GCM performance**

9 Ever since the first GCM was developed by Phillips (1956) (see Xu, 1999), attempts have
10 been made to assess the adequacy of GCM modelling. Initially, these evaluations were simple
11 side-by-side comparisons of individual monthly or seasonal means or multi-year averages
12 (Chervin, 1981). To assess model performance, Chervin (1981) extended the evaluation
13 procedure by examining statistically the agreement or otherwise of the ensemble average and
14 standard deviation between the GCM modelled climate and the observed data using the
15 vertical transient heat flux in an example application. Legates and Willmott (1992) compared
16 observed with simulated average precipitation rates by 10° latitude bands. On a two-
17 dimensional plot, Taylor (2001) developed a diagram in which each point consisted of the
18 spatial correlation coefficient and the spatial root mean square (RMS) along with the ratio of
19 the variances of the modelled and the observed variables. Recently, some authors have used
20 the Taylor diagram (Covey et al., 2003; Bonsal and Prowse, 2006) or a similar approach
21 (Lambert and Boer, 2001; Boer and Lambert, 2001). Murphy et al. (2004) introduced a
22 Climate Prediction Index (CPI) which is based on a broad range of present-day climates. This
23 index was later used by Johns et al. (2006) for a different set of climate variables than those
24 used by Murphy et al. (2004). Whetton et al. (2005) introduced a demerit point system in
25 which GCMs were rejected when a specified threshold was exceeded. Min and Hense (2006)
26 introduced a Bayesian approach to evaluate GCMs and argued that a skill-weighted average
27 with Bayes factors is more informative than moments estimated by conventional statistics.
28 Shukla et al. (2006) suggested that differences in observed and GCM simulated variables
29 should be examined in terms of their probability distributions rather than individual moments.
30 They proposed the differences could be examined using relative entropy. Perkins et al. (2007)
31 also claimed that assessing the performance of a GCM through a probability density function
32 (PDF) rather than using the first or a second moment would provide more confidence in

1 model assessment. To compare the reliability of variables (in time and space) rather than
2 individual models, Johnson and Sharma (2009a, 2009b) developed the Variable Convergence
3 Score which is used to rank a variable based on the ensemble coefficient of variation. They
4 observed the variables with the highest scores were pressure, temperature, and humidity.
5 Reichler and Kim (2008) introduced a Model Performance Index by first estimating a
6 normalised error variance based on the square of the grid-point differences between simulated
7 (interpolated to the observational grid) and the observed annual climate weighted and
8 standardised with respect to the variance of the annual observations. The error variance was
9 scaled by the average error found in the reference models and, finally, averaged over all
10 climates.

11
12 It is clear from this brief review that no one procedure has been universally accepted to assess
13 GCM performance, which is consistent with the observations of Räisänen (2007). We also
14 note the comments of Smith and Chandler (2010, page 379) who said: “It is fair to say that
15 any measure of performance can be subjective, simply because it will tend to reflect the
16 priorities of the person conducting the assessment. When different studies yield different
17 measures of performance, this can be a problem when deciding on how to interpret a range of
18 results in a different context. On the other hand, there is evidence that some models
19 consistently perform poorly, irrespective of the type of assessment. This would tend to
20 indicate that these model results suffer from fundamental errors which render them
21 inappropriate.”

22
23 In 1992, Legates and Willmott (1992) assessed the adequacy of GCMs based mainly on
24 January and July precipitation fields. Although a number of GCM assessments were carried
25 out during the following one and one-half decades, it was not until 2008 that mean
26 precipitation, either absolute or bias, was included in GCM published assessments. In that
27 year, Reichler and Kim (2008, page 303) argued that the mean bias is an important
28 component of model error.

29
30 In Tables 2a and 2b we summarize the application of the numerical metrics and the ranking
31 metrics of precipitation and temperature respectively applied to CMIP 3 data sets at the global

1 or country scales. These references cover the period from 2006 to 2014. Across these 15
2 papers, we observe that for precipitation and temperature the spatial root mean square error,
3 either using raw data (RMSE) or normalised data as a percentage of the mean value
4 (RRMSE), is adopted in seven of the 15 studies. (The data are normalised by the
5 corresponding standard deviation of the reference or observed data.) This spatial root mean
6 square metric, along with the bias in the mean of the data, is relevant to hydrologists as it
7 provides an indication of the uncertainty in the climate variables of interest to them. Of more
8 relevance to hydrologists is the uncertainty in temporal mean and variance of climatic
9 variables, which for precipitation are only reported in four of the 15 studies. Although spatial
10 correlation is not used directly in general hydrologic investigations, in GCM assessments it is
11 often combined with the variance and spatial RMSE through the Taylor diagram (Taylor,
12 2001) which is an excellent summary of the performance of a GCM projected variable. As
13 noted in Table 2, three papers utilize this approach. Lambert and Boer (2001, page 89)
14 extended the Taylor diagram to display the relative mean square differences, the pattern
15 correlations and the ratio of variances for modelled and observed data. This approach to
16 displaying the second order statistics appears not to have been widely adopted. It is noted in
17 Table 2a that only four papers include the mean or bias of the raw precipitation data in the
18 GCM assessments which is important from a hydrological perspective. **The second set of
19 metrics listed in Table 2b is used essentially for ranking GCMs by performance. Several other
20 assessment tools not included in Table 2b are the climate prediction index (Murphy et al.,
21 2004) and Bayesian approaches (Min and Hense, 2006).**

22

23 Specific climate features like the preservation of the ENSO signal (van Oldenborgh et al.,
24 2005) would also be considered to be a non-numerical measure of GCM performance, but in
25 some regions to be no less important to **hydrologists** than the numerical measures. Most of
26 these ranking metrics have been developed for specific purposes with respect to GCMs and
27 several have little utility to the practicing hydrologist who is primarily interested in bias,
28 variance and uncertainty in projected estimates of precipitation and temperature (plus net
29 radiation, wind speed and humidity to derive potential ET) as input to drive stand-alone
30 global and catchment hydrologic models.

31

1 2.2 Results of CMIP3 GCMs assessments

2 Table 2a indicates that only two papers (Räisänen, 2007; Gleckler et al., 2008) detail
3 numerical measures for both mean annual precipitation and temperature for 21 and 22 CMIP3
4 GCMs respectively at a global scale. Reifen and Toumi (2009) (17 GCMs) and Knutti et al.
5 (2010) (23 GCMs) address, inter alia, only mean annual temperature. Hagemann et al. (2011)
6 used three GCMs to estimate precipitation and temperature characteristics, but the paper
7 includes only precipitation results.

8

9 Räisänen (2007) results illustrate the wide range of model performances that exist: for
10 precipitation, RMSE = 1.35 mm day⁻¹ with a range of 0.97 – 1.86 and for temperature RMSE
11 = 2.32 °C with a range of 1.58 – 4.56. Reichler and Kim (2008) considered 14 variables
12 covering mainly the period 1979–1999 to assess the performance of CMIP3 models using
13 their Model Performance Index. They concluded that there was a continuous improvement in
14 model performance from the CMIP1 models compared to those available in CMIP3 but there
15 are still large differences in the CMIP3 models' ability to match observed climates. Gleckler
16 et al. (2008) normalised the data in Taylor diagrams for a range of climate variables and
17 concluded that some models performed substantially better than others. However, they also
18 concluded that it is not yet possible to answer the question: What is the best model?

19

20 Reifen and Toumi (2009) (Table 2b) using temperature anomalies observed that "... there is
21 no evidence that any subset of models delivers significant improvement in prediction
22 accuracy compared to the total ensemble". On the other hand, Macadam et al. (2010) (Table
23 2a) assessed the performance of 17 CMIP3 GCMs comparing the observed and modelled
24 temperatures over five 20-year periods and concluded that GCM rankings based on anomalies
25 can be inconsistent over time whereas rankings based on actual temperatures can be
26 consistent over time.

27

28 In summary, Gleckler et al. (2008) stated that the best GCM will depend on the intended
29 application. In the overarching project of which this study is a component, we are interested
30 in the uncertainty in annual streamflow estimated through hydrologic simulation using GCM
31 precipitation and temperature and how that uncertainty will affect estimates of future yield

1 from surface water reservoir systems. Consequently, we are interested in which GCMs
2 reproduce precipitation and temperature satisfactory. Based on the references of Reichler and
3 Kim (2008), Gleckler et al. (2008) and Macadam et al. (2010), the performance of 23 CMIP3
4 GCMs assessed at a global scale are ranked in Table 3. In the table eight models that meet the
5 Reichler and Kim (2008) criterion are also ranked in the upper 50% based on the Macadam et
6 al. (2010) and Gleckler et al. (2008) references. These models are CCCMA-t47, CCSM,
7 GFDL2.0, GFDL2.1, HadCM3, MIROCM, MPI and MRI.

8

9 **3 Data**

10 Two data sets are used in the GCM assessment that follows in Section 4. One is based on
11 observed data and the other on GCM simulations of present climate (20C3M). It should be
12 noted that of the 22 GCMs examined herein, multiple runs or projections were available for
13 nine models. The resulting 46 runs are identified in the tables summarising the results.

14

15 The first data set is based on monthly observed precipitation and temperature gridded at $0.5^\circ \times$
16 0.5° resolution over the global land surface from CRU 3.10 (New et al., 2002) for the period
17 January 1950 to December 1999. For grid cells where monthly observations are not available,
18 the CRU 3.10 dataset is based on interpolation of observed values within a ‘correlation decay
19 distance’ of 450 km for precipitation and 1,200 km for temperature. The CRU 3.10 data set
20 provides information about the number of observations within the correlation decay distance
21 of each grid cell for each month. In this analysis we defined a grid cell as ‘observed’ if $\geq 90\%$
22 of months at that grid cell has at least one observation within the correlation decay distance
23 for the period January 1950 to December 1999. Only ‘observed’ grid cells are used to
24 compute summary statistics in the following analysis.

25

26 The second data set is monthly precipitation and temperature data for the present climate
27 (20C3M) from 22 of the 23 GCMs listed in Table 1 and consists of 46 GCM runs. The
28 20C3M monthly data for precipitation and temperature were extracted from the CMIP3
29 dataset. As shown in Table 1 the GCMs have a wide range of spatial resolutions, all of which
30 are coarser than the observed CRU data. In order to make comparisons between observed and
31 GCM data either the CRU and/or GCM data must be re-sampled to the same resolution. To

1 avoid re-sampling coarse resolution data to a finer resolution we only re-sampled the CRU
2 data here. Thus in the following analysis the performance of each GCM is assessed at the
3 resolution of the GCM and the CRU data are re-sampled to match the GCM resolution.
4 Therefore, the number of grid cells in each comparison varies with the GCM resolution and
5 ranged from 616 to 11886 for the temperature comparisons and 425 to 8291 for the
6 precipitation comparisons. The difference in number of grid cells between temperature and
7 precipitation is due to more terrestrial grid cells having observed temperature data than
8 precipitation data over the period 1950-1999.

9

10 In the following analysis comparisons are made between observed and GCM values of mean
11 and standard deviation of annual precipitation and mean annual temperature. The GCM
12 values are based on *concurrent raw* (that is, not downscaled nor bias corrected) data from the
13 20C3M simulation. For example, if a grid cell has observed calendar-year data from 1953 –
14 1994, then the comparison is made with GCM values from the 20C3M run for the concurrent
15 calendar years 1953 – 1994. Although the aim of a 20C3M run from a given GCM is not to
16 strictly replicate the observed monthly record, we expect better performing GCMs to
17 reproduce mean annual statistics that are broadly similar to observed conditions. Average
18 monthly precipitation and temperature patterns are also compared to assess how well GCM
19 runs reproduce observed seasonality. Finally, we assess how well the Köppen-Geiger climate
20 classification (Peel et al., 2007) estimated from the CMIP3 data compares with present-day
21 gridded observed climate classification.

22

23 **4 Comparison of present climate GCM data with observed data**

24 In the analyses that follow, GCM estimates of mean annual precipitation and temperature and
25 the standard deviation of annual precipitation are compared against observed estimates for
26 terrestrial grid cells with $\geq 90\%$ observed data during the period 1950–1999.

27

28 Eight standard statistics – Nash-Sutcliffe efficiency (NSE) (Nash and Sutcliffe, 1970),
29 product moment coefficient of determination (R^2) (MacLean, 2005), standard error of
30 regression (Maidment, 1992), bias (MacLean, 2005), percentage bias (Maidment, 1992),
31 absolute percentage bias (MacLean, 2005), root mean square error (RMSE) (MacLean, 2005)

1 and mean absolute error (MacLean, 2005) – were computed as the basis of comparison but we
2 report only the NSE, R^2 and RMSE in the following discussion. For our analysis, the NSE is
3 the most useful statistic as it shows the proportion of explained variance relative to the 1:1
4 line in a comparison of two estimates of the same variable. R^2 is included because many
5 analysts are familiar with its interpretation. Both NSE and R^2 were computed in arithmetic
6 (untransformed) and natural log space. We have also included RMSE values (computed from
7 the untransformed values) as many GCM analyses include this measure.

8

9 In the following sub-sections comparisons between the concurrent raw GCM data and
10 observed values for MAP, SDP, MAT, long-term average monthly precipitation and
11 temperature patterns, and Köppen-Geiger climate classification at the grid cell scale are
12 presented and discussed. Although we rank the models by each selection criteria and combine
13 the ranks by addition, we note the warning of Stainforth et al. (2007) who argue that model
14 response should not be weighted but ruled in or out. We follow this approach in this paper by
15 identifying better performing GCMs to be used for hydrologic simulations reported in a
16 companion paper (Peel et al., 2014). This approach is consistent with the concept recognised
17 by Randall et al. (2007, page 608) that “... for models to predict future climatic conditions
18 reliably, they must simulate the current climatic state with some as yet unknown degree of
19 fidelity. Poor model skill in simulating present climate could indicate that certain physical or
20 dynamical processes have been misrepresented”. It is noted that our comparisons are
21 conducted over the global terrestrial land surface rather than focussing on a single catchment,
22 region or continent. This allows us to assess whether a GCM performs consistently well
23 across a large area and reduces the chance of a GCM being selected due to a random high
24 performance over a small area.

25 **4.1 Mean annual precipitation**

26 Comparisons of mean annual precipitation and the standard deviation of annual precipitation
27 between GCM estimates and observed data for the grid cells across the 46 runs are presented
28 in Table 4. For MAP, the Nash-Sutcliffe efficiency varied from a maximum of 0.68 ($R^2 =$
29 0.69) with a RMSE value of 335 mm year⁻¹ for model MIUB(3) to -0.54 for GISS-EH(3).
30 (GCM run number is enclosed by parenthesis, for example MIUB(3) is run 3 for the GCM
31 MIUB.) The MAP values for MIUB(3) are compared with the observed CRU MAP values in

1 Figure 1. Each data point in this figure represents a MAP comparison at one of the 632
2 MIUB(3) terrestrial grid cells where observed CRU 3.10 data were available for the period
3 January 1950 to December 1999. The relationship between GCM and observed MAP shown
4 in this figure is representative of the other GCMs where high MAPs are underestimated and
5 low MAPs are overestimated. GISS-EH(3), shown in Figure 2, is an example of a poorly
6 performing GCM in terms of mean annual precipitation. Here, based on untransformed data,
7 the Nash-Sutcliffe efficiency is -0.54 ($R^2 = 0.37$) with a RMSE value of 697 mm year^{-1} .

8

9 The range of NSE values for the MAP comparisons across the 46 GCM runs is plotted in
10 Figure 3. The results may be classified into four groups: 5 runs exhibiting $\text{NSE} > 0.6$, 27 runs
11 $0.4 < \text{NSE} \leq 0.6$, 6 runs $0 < \text{NSE} \leq 0.4$ and 8 runs ≤ 0 where the predictive power of the GCM
12 is less than using the average observed MAP across all grid cells (Gupta et al., 2009).

13 **4.2 Standard deviation of annual precipitation**

14 For the standard deviation of annual precipitation, HadCM3 was the best performing model
15 with a NSE of 0.57, R^2 of 0.62 and a RMSE of 51 mm year^{-1} . MIROCH also yielded a NSE of
16 0.57 and an R^2 of 0.58 but with a RMSE of 63 mm year^{-1} . These results along with other
17 standard deviation values are listed in Table 4. Figure 4 is a plot for MIUB(3), which is
18 representative (rank 4, that is the 4th best performance of the 46 runs) of the relationship
19 between GCM and observed SDP, and shows the model underestimates the standard deviation
20 of annual precipitation for high values and overestimates at low values of standard deviation
21 compared with observed values.

22 **4.3 Mean annual temperature**

23 The comparison of the GCM mean annual temperatures with concurrent observed data for the
24 grid cells are listed for each model run in Table 4. In contrast to the precipitation modelling,
25 the mean annual temperatures are simulated satisfactorily by most of the GCMs. Except for
26 the IAP and the GFLD2.0 models (NSE = ~ 0.90 and 0.93 respectively), all model runs exhibit
27 NSE values ≥ 0.94 with 17 of the 46 GCM runs having a NSE value ≥ 0.97 . A comparison
28 between MIUB(3) estimates of mean annual temperature (NSE = 0.96, rank 33) and observed
29 values from the CRU data set is presented in Figure 5. Also shown in Figure 5 is a linear fit

1 between GCM and observed MAT. The average fit for the 46 GCM runs (not shown)
2 exhibited a small negative bias of -1.03°C and a slope of 1.01.

3 **4.4 Average monthly precipitation and temperature patterns**

4 Because a monthly rainfall-runoff model is applied in the next phase of our analysis (reported
5 in a companion paper) it is considered appropriate to assess how well the GCMs simulate the
6 observed mean monthly patterns of precipitation and temperature (see also the argument of
7 Charles et al., 2007). The NSE was used for the assessment by comparing the 12 long-term
8 average monthly values. For each GCM run the average precipitation and temperature values
9 for each month were calculated for each grid cell. Nash-Sutcliffe efficiencies were computed
10 between the equivalent 12 GCM and 12 CRU based monthly averages. The median NSE
11 values across terrestrial grid cells where observed CRU 3.10 data were available for the
12 period January 1950 to December 1999 for each GCM run are summarised in Table 4. As
13 shown in Table 4 average monthly patterns of precipitation are poorly modelled. In fact, 57%
14 of the 46 model runs have a median NSE value of < 0 . For these GCMs their predictive power
15 for the monthly precipitation pattern is less than using the average of the 12 monthly values at
16 each of the terrestrial grid cells. Only two GCMs have NSE values > 0.25 . In contrast, the
17 median NSEs of all monthly temperature patterns are > 0.75 , with 41% > 0.90 . The NSE
18 metric reflects how well the GCM replicates both the monthly pattern and the overall average
19 monthly value (bias). Thus the monthly pattern of temperature is generally well reproduced
20 by the GCMs, whereas the monthly pattern of precipitation is not, which is mainly due to the
21 bias in the GCM average monthly precipitation.

22 **4.5 Köppen-Geiger classification**

23 The Köppen-Geiger climate classification (Peel et al., 2007) (see Table 5) provides an
24 alternate way to assess the adequacy of how well a GCM represents climate because the
25 classification is based on a combination of annual and monthly precipitation and temperature
26 data. Two comparisons between the MPI(3) model and CRU observed data are presented in
27 Table 6. The MPI(3) was chosen as an example here as over the three levels of climate classes
28 it estimated the observed climate correctly more often than the other model runs. In Table 6(a)
29 a comparison at the first letter level of the Köppen-Geiger climate classification is shown.
30 This comparison reveals how well the GCM reproduces the distribution of broad climate
31 types: tropical, arid, temperate, cold and polar over the terrestrial surface. In Table 6(b) the

1 comparison shown is for the second letter level of the Köppen-Geiger climate classification,
2 which assesses how well the GCM reproduces finer detail within the broad climate types; for
3 example, the seasonal distribution of precipitation or whether a region is semi-arid or arid.
4 The diagonal values shown in Tables 6(a) and 6(b) represent the number of grid cells
5 correctly classified by the GCM whereas the off-diagonal values are the number of grid cells
6 incorrectly classified by the GCM for the one- and two-letter level respectively. At the first
7 letter level MPI(3) reproduces the correct climate type at 81% of the terrestrial grid cells.
8 Within this good performance the MPI(3) produces more polar climate and less tropical and
9 cold grids cells than observed. At the second letter level, MPI(3) reproduces the correct
10 climate type at 67% of the terrestrial grid cells. The model produces less grid cells of tropical
11 rainforest, cold with a dry winter and cold without a dry season than expected and more cold
12 with a dry summer and polar tundra than expected.

13

14 Table 7 summarises the overall proportion of GCM grid cells that were classified correctly for
15 each GCM run across the three levels of classification. As we wish to have a ranking of the
16 comparisons we adopted this simple measure as it is regarded as “... one of the most basic and
17 widely used measures of accuracy...” for comparing thematic maps (Foody, 2004, page 632).
18 From Table 7 we observe that GCM accuracy in reproducing the climate classification
19 decreases as one moves from coarse to fine detail climate classification. The average accuracy
20 (and range) for the three classes are: 0.48 (0.36-0.60) for the three-letter classification, 0.57
21 (0.47-0.68) for the two-letter classification, and for one-letter 0.77 (0.66-0.82). In other
22 words, at the three-letter scale nearly 50% of GCM Köppen-Geiger estimates are correct,
23 increasing to nearly 60% at the two-letter level and, finally, at the one-letter aggregation more
24 than 75% are correct across the 46 GCM runs. Using these average values across the three
25 classes, the following seven models performed satisfactorily in identifying Köppen-Geiger
26 climate class correctly: CNRM, CSIRO, HadCM3, HadGEM, MIUB, MPI, and MRI. Of
27 these models the least successful run was for CSIRO with the percentage correct for each
28 class being: three-letter 51%, two-letter 60% and one-letter 78%.

29

1 **5 Discussion**

2 **5.1 Relating GCM resolution to performance**

3 In the analysis presented in the previous section each GCM's performance in reproducing
4 observed climatological statistics was assessed at the resolution of the individual GCM. The
5 question of whether GCMs with a finer resolution outperform GCMs with a coarser resolution
6 is addressed in Figure 6 where GCM performance in reproducing observed terrestrial MAP
7 and MAT, based on the NSE, is related to GCM resolution, defined as the number of grid
8 cells used in the comparison. The plot suggests there is no significant relationship between
9 GCM resolution and GCM performance beyond 1500 grid cells for either MAP or MAT.
10 Interestingly, some lower resolution GCMs, <1500 grid cells, perform as well as higher
11 resolution GCMs for MAP and MAT, yet for others, they perform poorly. While it is
12 sometimes assumed that higher resolution should normally lead to improved performance,
13 there are many other factors that affect performance. These include the sophistication of the
14 parameterisation schemes for different sub-grid scale processes, the time spent in developing
15 and testing the individual schemes and their interactions. Our purpose here is to report this
16 observation rather than speculate what it might mean for GCM model development. Our
17 observation is consistent with Masson and Knutti (2011) who comment that "... model
18 resolution in CMIP3 seems to only affect performance in simulating present-day temperature
19 for small scales over land" (page 2691) and for precipitation they comment that "...no clear
20 relation seems to exist at least within the relatively narrow range of resolutions covered by
21 CMIP3" (page 2686).

22 **5.2 Joint comparison of precipitation and temperature**

23 In using GCM climate scenarios in a water resources study, it is appropriate to ensure
24 consistency between precipitation and temperature by adopting projections of these variables
25 from the same GCM run. Grid cell based NSEs for mean annual temperature and mean annual
26 precipitation from each GCM are compared in Figure 7, which illustrates the performance of
27 each GCM for both variables. Models that have relatively high NSEs for precipitation do not
28 necessarily have relatively high values for temperature. It is interesting to note that the rank of
29 the models based on NSE of the MAP is unrelated to the ranking of the models based on
30 MAT. Fortunately, however, most of the NSEs for MAT are relatively high and the

1 acceptance or rejection of a GCM as a better performing model is largely dependent on its
2 precipitation characteristics.

3 **5.3 Identifying better performing GCMs**

4 To identify the better performing GCMs across the different variables assessed, the results in
5 Table 4 are ranked by Nash-Sutcliffe efficiency and summarised in Table 8. The monthly
6 patterns of precipitation and temperature are combined by ranking the average of their
7 respective NSE values. The overall rank for each GCM run is based on combining, by
8 addition, the ranks for the individual variables and, finally, identifying the best performing
9 run from each GCM. Selection of the better performing GCMs using these rankings is not
10 inconsistent with Stainforth et al. (2007) who argued that model response should not be
11 weighted but ruled in or out. From Table 8 we identify several GCMs, listed in Table 9, as
12 better performing models. These selected GCMs were based on the assumption that
13 performance across the four variables (MAP, SDP, MAT and combined monthly pattern) is
14 equally weighted. GCMs that achieved MAP NSE > 0.50, SDP NSE > 0.45, MAT NSE >
15 0.95 and mean monthly pattern of precipitation NSE > 0.0 (Table 4) were identified as better
16 performing. (Because nearly all the GCM runs modelled mean monthly patterns of
17 temperature satisfactorily, this measure was not considered in the selection of models listed in
18 column 1, Table 9.) The following GCMs were selected (Table 9): HadCM3, INGV,
19 MIROCM, MIUB, MPI, and MRI. INGV was included although it failed the monthly
20 precipitation pattern criterion. The above criteria were selected to identify a small number of
21 GCMs that would require less bias correction to produce annual precipitation and temperature
22 consistent with observations.

23

24 In Table 9, we summarise our observations from the literature review in Section 2 and the
25 results from our analyses in Tables 4 and 8, where we identified six GCMs that satisfied our
26 selection criteria (Table 9, column 1). From the literature review (Table 3), eight GCMs were
27 identified as being satisfactory. We have added MIUB because in the literature review it
28 ranked first overall, although no guidance was available from Reichler and Kim (2008). We
29 also added MIROCH to this list as it performed better according to Gleckler et al. (2008) than
30 several models in the above list and met the performance index of Reichler and Kim (2008).
31 Columns (1) and (2) of Table 9 suggest there is some consistency between our analyses from

1 a hydrological perspective and that reported in the literature from a climatological
2 perspective. From the table, we identify that, in terms of our objective to assess how well the
3 CMIP3 GCMs are able to reproduce observed annual precipitation and temperature statistics
4 and the mean monthly patterns of precipitation and temperature, the following models are
5 deemed acceptable for the next phase of our project: HadCM3, MIROCM, MIUB, MPI and
6 MRI. Although not used in the selection criteria we observe our selected GCMs performed
7 well in the Köppen-Geiger climate assessment. We note here that INGV also performed
8 satisfactorily but it was not included in our adopted GCMs as it was not reviewed in the
9 papers of Gleckler et al. (2008), Reichler and Kim (2008) and Macadam et al. (2010).

10 **5.4 Comparing future responses of selected GCMs**

11 In order to confirm that the selected GCM runs are representative of the range of future
12 responses to climate change in the CMIP3 ensemble, we plot in Figure 8 the ratio of mean
13 annual precipitation for the period 2015-2034 (from the A1B scenario) to 1965-1994 against
14 the mean annual temperature difference between 2015-2034 and 1965-1994 for the global
15 land surface. The five selected GCM runs are well distributed amongst the 44 GCM ensemble
16 members, which indicate that the selected GCMs are reasonably representative of the range of
17 future GCM projections if all the runs were considered. We observe that most GCM runs are
18 clustered around the median response, except for the seven CCSM runs in the top right
19 quadrant with a precipitation ratio $> \sim 1.04$.

20

21 **6 Conclusions**

22 Our primary objective in this paper is to identify better performing GCMs from a hydrologic
23 perspective over global land regions. The better performing GCMs were identified by their
24 ability to reproduce observed climatological statistics (mean and the standard deviation of
25 annual precipitation and mean annual temperature, and the mean monthly patterns of
26 precipitation and temperature) for hydrologic simulation. The GCM selection process was
27 informed by our results presented here and by a literature review of CMIP3 GCM
28 performance. In terms of the Nash-Sutcliffe efficiency there was a large spread in values for
29 mean annual precipitation and the standard deviation of annual precipitation over concurrent
30 periods. The highest NSE for mean annual precipitation was 0.68 and 0.57 for the standard
31 deviation of annual precipitation. On the other hand, for mean annual temperatures, the NSEs

1 between modelled and observed data were very high, with median NSE being 0.97. Overall,
2 all GCMs reproduced the Köppen-Geiger climate satisfactorily at the broad first letter level.
3 From the literature, the following GCMs were identified as being suitable to simulate annual
4 precipitation and temperature statistics: CCCMA-T47, CCSM, GFDL2.0, GFDL2.1,
5 HadCM3, MIROCH, MIROCM, MIUB, MPI and MRI. After combining our results with the
6 literature the following GCMs were considered the better performing models from a
7 hydrologic perspective: HadCM3, MIROCM, MIUB, MPI and MRI. The future response of
8 the better performing GCMs was found to be representative of the 44 GCM ensemble
9 members which confirms that the selected GCMs are reasonably representative of the range of
10 future GCM projections. Our approach for evaluating GCM performance for hydrologic
11 simulation could be applied to CMIP5 runs.

12

13 **Acknowledgements**

14 This research was financially supported by Australian Research Council Grant LP100100756
15 and FT120100130, Melbourne Water and the Australian Bureau of Meteorology. Lionel
16 Siriwardena, Sugata Narsey and Dr Ian Smith assisted with extraction and analysis of CMIP3
17 GCM data. Lionel Siriwardena also assisted with extraction and analysis of the CRU 3.10
18 data. We acknowledge the modeling groups, the Program for Climate Model Diagnosis and
19 Intercomparison (PCMDI) and the WCRP's Working Group on Coupled Modelling (WGCM)
20 for their roles in making available the WCRP CMIP3 multi-model dataset. Support of this
21 dataset is provided by the Office of Science, U.S. Department of Energy. The authors thank
22 two anonymous reviewers who provided stimulating comments on the discussion paper.

23

1 **References**

- 2 Boer, G. J., and Lambert, S. J.: Second order space–time climate difference statistics, *Clim.*
3 *Dynam.*, 17, 213–218, 2001.
- 4 Bonsal, B. T., and Prowse, T. D.: Regional assessment of GCM-simulated current climate
5 over Northern Canada, *Arctic*, 59(2), 115-128, 2006.
- 6 Charles, S. P., Bari, M. A., Kitsios, A., and Bates, B. C.: Effect of GCM bias on downscaled
7 precipitation and runoff projections for the Serpentine catchment, Western Australia, *Int. J.*
8 *Climatol.*, 27, 1673-1690, 2007.
- 9 Chervin, R. M.: On the Comparison of Observed and GCM Simulated Climate Ensembles,
10 *Journal of Atmospheric Sciences*, 38(5), 885-901, 1981.
- 11 Chiew, F. H. S., and McMahon, T. A.: Modelling the impacts of climate change on Australian
12 streamflow, *Hydrol. Process.*, 16, 1235–1245, 2002.
- 13 Covey, C., Achutarao, K. M., Cubasch, U., Jones, P., Lambert S. J., Mann, M. E., Phillips, T.
14 J., and Taylor, K. E.: An overview of results from the Coupled Model Intercomparison
15 Project, *Global and Planetary Change*, 37, 103–133, 2003.
- 16 Dessai, S., Lu, X., and Hulme, M.: Limited sensitivity analysis of regional climate change
17 probabilities for the 21st century, *J. Geophys. Res.*, 110, D19108,
18 doi:10.1029/2005JD005919, 2005.
- 19 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S. C., Collins, W., Cox, P.,
20 Driouech, F., Emori, S., Eyring, V., Forest, C., Gleckler, P., Guilyardi, E., Jakob, C., Kattsov,
21 V., Reason, C., and Rummukainen, M.: Evaluation of Climate Models. In: *Climate Change*
22 *2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment*
23 *Report of the Intergovernmental Panel on Climate Change [Stocker, T. F., D. Qin, G.-K.*
24 *Plattner, M. Tignor, S. K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P. M. Midgley*
25 *(eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA,*
26 *2013.*
- 27 Foody, G. M.: Thematic map comparison: Evaluating the statistical significance of
28 differences in classification accuracy, *Photogramm. Eng. Rem. S.*, 70(5), 627-633, 2004.
- 29 Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, *J.*
30 *Geophys. Res.-Atmos.*, 113, D06104, doi:10.1029/2007JD008972, 2008.

1 Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean
2 squared error and NSE performance criteria: Implications for improving hydrological
3 modelling, *J. Hydrol.*, 377, 80-91, 2009.

4 Hagemann, S., Chen, C., Haerter, J.O., Heinke, J., Gerten, D. and Piani, C.: Impact of a
5 statistical bias correction on the projected hydrological changes obtained from three GCMs
6 and two hydrology models. *J. Hydrometeor.* 12, 556-578, 2011.

7 Heo, K.-Y., Ha, K.-J., Yun, K.-S., Lee, S.-S., Kim, H.-J. and Wang, B.: Methods for
8 uncertainty assessment of climate models and model predictions over East Asia. *International*
9 *Journal of Climatology*, 34:10.1002/joc.2014.34.issue-2, 377-390, 2014.

10 Johns, T. C., Durman, C. F., Banks, H. T., Roberts, M. J., McLaren, A. J., Ridley, J. K.,
11 Senior, C. A., Williams, K. D., Jones, A., Rickard, G. J., Cusack, S., Ingram, W. J., Crucifix,
12 M., Sexton, D. M. H., Joshi, M. M., Dong, B.-W., Spencer, H., Hill, R. S. R., Gregory, J. M.,
13 Keen, A. B., Pardaens, A. K., Lowe, J. A., Bodas-Salcedo, A., Stark, S., and Searl, Y.: The
14 new Hadley Centre climate model (HadGEM1): evaluation of coupled simulations, *J.*
15 *Climate*, 19, 1327–1353, 2006.

16 Johnson, F. M., and Sharma, A.: GCM simulations of a future climate: How does the skill of
17 GCM precipitation simulations compare to temperature simulations, 18th World
18 IMACS/MODSIM Congress, Cairns, Australia, 2009a.

19 Johnson, F. and Sharma, A.: Measurement of GCM skill in predicting variables relevant for
20 hydroclimatological assessments, *J. Climate*, 22(16), 4373-4382, 2009b.

21 Knutti, R., Furrer, R., Tebaldi, C., Cermak, J. And Meehl, G.A.: Challenges in combining
22 projections from multiple climate models. *Journal of Climate*, 23, 2739-2758.

23 Knutti, R., Masson, D., and Gettelman, A.: Climate model genealogy: Generation CMIP5 and
24 how we got there, *Geophys. Res. Lett.*, 40, 1194–1199, 2013.

25 Lambert, S. J., and Boer, G. J.: CMIP1 evaluation and intercomparison of coupled climate
26 models, *Clim. Dynam.*, 17, 83–106, 2001.

27 Legates, D. R., and Willmott, C. J.: A comparison of GCM-simulated and observed mean
28 January and July precipitation, *Global and Planetary Change*, 5(4), 345-363, 1992.

- 1 Macadam, I., Pitman, A. J., Whetton, P. H., and Abramowitz, G.: Ranking climate models by
2 performance using actual values and anomalies: Implications for climate change impact
3 assessments, *Geophys. Res. Lett.*, 37, L16704, 2010.
- 4 MacLean, A.: Statistical evaluation of WATFLOOD (Ms), 2005.
- 5 Maidment, D. R.: *Handbook of Hydrology*. McGraw-Hill Inc., 1992.
- 6 Masson, D. and Knutti, R.: Spatial-scale dependence of climate model performance in the
7 CMIP3 ensemble, *J. Clim.* 24, 2680-2692, 2011.
- 8 McMahon, T. A., and Adedoye, A. J.: *Water Resources Yield*. Water Resources Publications,
9 CO, USA, 220pp, 2005.
- 10 McMahon, T. A., Peel, M. C., Pegram, G. G. S., and Smith, I. N.: A simple methodology for
11 estimating mean and variability of annual runoff and reservoir yield under present and future
12 climates, *J. Hydrometeorol.*, 12(1), 135-146, 2011.
- 13 Meehl, G. A., Covey, C., Delworth, T., Latif, M., McAvaney, B., Mitchell, J. F. B., Stouffer,
14 R. J., and Taylor, K. E.: The WCRP CMIP3 multi-model dataset: A new era in climate change
15 research, *B. Am. Meteorol. Soc.*, 88, 1383-1394, 2007.
- 16 Min, S.-K., and Hense, A.: A Bayesian approach to climate model evaluation and multi-
17 model averaging with an application to global mean surface temperatures from IPCC AR4
18 coupled climate models, *Geophys. Res. Lett.*, 33, doi:10.1029/2006GL025779, 2006.
- 19 Murphy, J. M., Sexton, D. M. H., Barnett, D. N., Jones, G. S., Webb, M. J., Collins, M. J., and
20 Stainforth, D. A.: Quantification of modelling uncertainties in a large ensemble of climate
21 change simulations, *Nature*, 430, 768–772, 2004.
- 22 Nash, J. E., and Sutcliffe, J. V.: River flow forecasting through conceptual models Part 1 – A
23 discussion of principles, *J. Hydrol.*, 10, 282-290, 1970.
- 24 New, M., Lister, D., Hulme, M., and Makin, I.: A high-resolution data set of surface climate
25 over global land areas, *Clim. Res.*, 21, 1–25, 2002.
- 26 Peel, M. C., Finlayson, B. L., and McMahon, T. A.: Updated world map of the Köppen-
27 Geiger climate classification, *Hydrol. Earth Syst. Sc.*, 11, 1633-1644, 2007.

- 1 Peel, M. C., Srikanthan, R., McMahon, T. A., and Karoly, D. J.: Approximating uncertainty
2 of annual runoff and reservoir yield using stochastic replicates of Global Climate Model data,
3 Submitted to Hydrol. Earth Syst. Sc., 2014.
- 4 Perkins, S. E., Pitman, A. J., Holbrook, N. J., and McAneney, J.: Evaluation of the AR4
5 climate models simulated daily maximum temperature, minimum temperature and
6 precipitation over Australia using probability density functions, *J. Climate*, 20, 4356-4376,
7 2007.
- 8 Phillips, N. A.: The general circulation of atmosphere: a numerical experiment, *Q. J. R.*
9 *Meteorol. Soc.*, 82, 123–164, 1956.
- 10 Räisänen, J.: How reliable are climate models? *Tellus A*, 59, 2-29, 2007.
- 11 Raju, K.S. and Kumar, D.N.: Ranking of global climate models for India using multicriterion
12 analysis. *Climate Research*, 60, 103-117, 2014.
- 13 Randall, R.A. and Wood, R.A. (Coordinating lead authors): Climate models and their
14 evaluation. Contribution of Working Group I to the Fourth Assessment Report of the
15 Intergovernmental Panel on Climate Change AR4, Chapter 8, 589-662, 2007.
- 16 Reichler, T., and Kim, J.: How well do coupled models simulate today's climate? *B. Am.*
17 *Meteorol. Soc.*, 89, 303–311, 2008.
- 18 Reifen, C., and Toumi, R.: Climate projections: Past performance no guarantee of future skill?
19 *Geophys. Res. Lett.*, 36, L13704, doi:10.1029/2009GL038082, 2009.
- 20 Shukla, J., DelSole, T., Fennessy, M., Kinter, J., and Paolino, D.: Climate model fidelity and
21 projections of climate change, *Geophys. Res. Lett.*, 33, L07702, DOI:
22 10.1029/2005GL025579, 2006.
- 23 Smith, I., and Chandler, E.: Refining rainfall projections for the Murray Darling Basin of
24 south-east Australia – the effect of sampling model results based on performance, *Climatic*
25 *Change*, 102, 377-393, 2010.
- 26 Stainforth, D. A., Allen, M. R., Tredger, E. R., and Smith, L. A.: Confidence, uncertainty and
27 decision-support relevance in climate predictions, *Philos. T. R. Soc. A*, 365, 2145-2161, 2007.
- 28 Suppiah, R., Hennessy, K. L., Whetton, P. H., McInnes, K., Macadam, I., Bathols, J.,
29 Ricketts, J. and Page, C. M.: Australian climate change projections derived from simulations
30 performed for IPCC 4th Assessment Report. *Aust. Met. Mag*, 56, 131-152, 2007.

- 1 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, J.
2 Geophys. Res., 106, 7183–7192, 2001.
- 3 van Oldenborgh, G. J., Philip, S. Y. and Collins, M.: El Niño in a changing climate: a multi-
4 model study. Ocean Science, 1, 81-95, 2005.
- 5 Watterson, I. G.: Calculation of probability density functions for temperature and
6 precipitation change under global warming. Journal of Geophysical Research 113, D12106,
7 doi:10.1029/2007JD009254, 2008.
- 8 Whetton, P., McInnes, K. L., Jones, R. J., Hennessy, K. J., Suppiah, R., Page, C. M., and
9 Durack, P. J.: Australian Climate Change Projections for Impact Assessment and Policy
10 Application: A Review. CSIRO Marine and Atmospheric Research Paper 001,
11 www.cmar.csiro.au/e-print/open/whettonph_2005a.pdf, 2005.
- 12 Xu, C. Y.: Climate change and hydrologic models: A review of existing gaps and recent
13 research developments, Water Resour. Manag., 13, 369-382, 1999.

Table 1. Details of 23 GCMs considered in this paper.

Acronym	Originating group	Country	Model name in CMIP3	Number of 20C3M runs available	Resolution		Number of prec. grid cells*	Number of temp. grid cells‡
					Lat (°)	Long (°)		
BCCR	Bjerkness Centre for Climate Research	Norway	bccr-bcm2.0	na†	1.9	1.9	na	na
CCCMA-t47	Canadian Centre for Climate Modeling and Analysis	Canada	cccma_cg3_1_t47	1	~3.75	3.75	631	916
CCCMA-t63	Canadian Centre for Climate Modeling and Analysis	Canada	cccma_cg3_1_t63	1	~2.8	2.8125	1169	1706
CCSM	National Centre for Atmospheric Research	USA	ccsm	8	~1.4	1.40625	5184	7453
CNRM	Météo-France / Centre National de Recherches Météorologiques	France	cnrm	1	~2.8	2.8125	1169	1706
CSIRO	Australia CSIRO	Australia	csiro_mk3_0	1	~1.87	1.875	2820	4068
GFDL2.0	NOAA Geophysical Fluid Dynamics Laboratory	USA	gfdl2_cm2_0	1	2	2.5	1937	2828
GFDL2.1	NOAA Geophysical Fluid Dynamics Laboratory	USA	gfdl2_cm2_1	1	~2	2.5	1911	2758
GISS-AOM	NASA Goddard Institute of Space Studies	USA	giss_aom_r1, 2	2	3	4	754	1076
GISS-EH	NASA Goddard Institute of Space Studies	USA	giss_eh1, 2,3	3	3 and 4	5	425	616
GISS-ER	NASA Goddard Institute of Space Studies	USA	giss_model_e_r	3	3 and 4	5	425	616
HadCM3	Hadley Centre for Climate Prediction and Research	UK	hadcm3	1	2.5	3.75	982	1421
HadGEM	Hadley Centre for Climate Prediction and Research	UK	HadGem	1	1.25	1.875	4316	6239
IAP	Institute of Atmospheric Physics, Chinese Acad. Sciences	China	iap_fgoals1.0_g	3	6.1 ~2.8	2.8125	1159	1664
INGV	National Institute of Geophysics and Vulcanology, Italy	Italy	ingv20c ECHAM4.6	1	~1.1	1.125	8291	11886
INM	Institute for Numerical Mathematics, Russia	Russia	inmcm3.0	1	4	5	420	620
IPSL	Institute Pierre Simon Laplace	France	ipsl_cm4	1	~2.5	3.75	980	1403
MIROCH	Center for Climate System Research (The	Japan	miroc3_2_hires	1	~1.1	1.125	8291	11886

	University of Tokyo), National Institute for Environmental Studies, and Frontier Research Center for Global Change		(mirochi)						
MIROCM	Center for Climate System Research (The University of Tokyo), National Institute for Environmental Studies, and Frontier Research Center for Global Change	Japan	miroc3_2_medres (mirocmedr)	3	~2.8	2.8125	1169	1706	
MIUB	Meteorological Institute of the University of Bonn, Meteorological Research Institute of KMA, and Model and Data group	Germany South Korea	miub_echo_g	3	~3.7	3.75	631	916	
MPI	Max Planck Institute for Meteorology	Germany	mpi_echam5 (mpi)	3	~1.8	1.875	2820	4068	
MRI	Japan Meteorological Research Institute	Japan	mri_cgcm2_3_2a (mri)	5	~2.8	2.8125	1169	1706	
PCM	National Center for Atmospheric Research	USA	pcm	1	~2.8	2.8125	1169	1706	

† na: Not available. * Based on mean annual precipitation comparison between GCM and CRU. ‡ Based on mean annual temperature comparison between GCM and CRU.

Table 2a: Numerical measures of performance assessment of CMIP3 GCMs

Reference	Global, country, large region	GCMs	Precipitation							Temperature					
			Reference data sets	Mean of raw data	Bias in mean of raw data	Variance of raw data	RMS or similar metric	Spat. correl.	Taylor plots	Reference data sets	Mean of raw data	Bias in mean of raw data	RMS or similar metric	Spat. correl.	Taylor plots
Bonsal & Prowse (2006)	Northern Canada	7 GCMs	CRU & other data		yes					yes (abs)		yes			yes (abs)
Perkins et al. (2007)	Australia	16 GCMs	Bureau of Met., Australia								Bureau of Met., Australia				
Suppiah et al. (2007)	Australia	23 GCMs	Bureau of Met., Australia				yes (abs)^	yes			Bureau of Met., Australia		yes	yes	
Räsänen (2007)	Global	21 GCMs	CRU, GPCPv2		yes as figure	yes as figure	yes	yes			CRU TS2.0, NCEP-NCAR		yes as figure	yes (abs)	yes
Gleckler et al. (2008)	Global	22 GCMs	GPCP/ CMAP				yes (norm) #		yes (norm)		ERA40/NCEP-NCAR			yes (norm)	yes (norm)
Reifen & Toumi (2009)	Global	17 GCMs									HadCRUT3 5° x 5°			yes (abs)	
Knutti et al. (2010)	Global	23 GCMs									ERA40		yes	yes (abs)	yes
Macadam et al. (2010)	Global	17 GCMs	HadCRUT3 dataset									yes as figure*			
Hagemann et al. (2011)	Global	MPI CNRM IPSL	WFD (ERA-40)	yes as figure	yes as figure	yes as figure									

Heo et al. (2014)	East Asia	21GCMs	CMAP			yes (norm)	yes	yes (norm)	NCEP-I			yes (norm)	yes	yes (norm)
Raju and Kumar (2014)	India	11 GCMs	NCAP/NCAR 2.5° × 2.5°		yes (norm)	yes (abs)			NCAP/NCAR 2.5° × 2.5°			yes (norm)		
Number of references		11		1	4	2	5	3	3	2	2	7	4	3

^ (abs): based on absolute data; # (norm): based on normalized data; ^ (ratio): as ratio of reference data; * also as an anomaly

1 Table 2b Ranking measures of performance assessment of CMIP3 GCMs

Reference	Global, country, large region	GCMs	PDF & related measures	Performance index based on variance	Entropy	Skill score	Variance convergence score	Signal noise ratio
Shukla et al. (2006)	Global	13 GCMs			yes			
Perkins et al. (2007)	Australia	16 GCMs	yes			yes		
Gleckler et al. (2008)	Global	22GCMs				yes		
Reichler & Kim (2008)	Global	21 GCMs		yes				
Watterson (2008)	Australia	23 GCMs	yes					
Johnson & Sharma (2009b)	Australia	9 GCMs					yes	
Knutti et al. (2010)	Global	23 GCMs	yes					
Heo et al. (2014)	East Asia	21 GCMs			yes			yes
Number of references		8	3	1	2	2	1	1

2

1 Table 3. Summary of performance of 23 CMIP3 GCMs in simulating present climate based
 2 on literature review

	Source	Macadam et al. (2010)	Gleckler et al. (2008)		Reichler and Kim (2008)
	Variables	Temperature	Precipitation		Many
GCM	Method	Ranking [‡]	Rel error ranking	Overall rank	Perf index*
BCCR		15	=13 [#]	14	No
CCCMA-t47		9	=1	=3	Yes
CCCMA-t63		na [†]	=3	na	Yes
CCSM		6	=13	=10	Yes
CNRM		8	=19	13	No
CSIRO		12	=3	7	No
GFDL2.0		16	=3	=10	Yes
GFDL2.1		5	=3	2	Yes
GISS-AOM		na	=13	na	No
GISS-EH		na	=19	na	No
GISS-ER		1	=17	9	No
HadCM3		2	=9	5	Yes
HadGEM		14	=17	15	Yes
IAP		na	=9	na	No
INGV		na	na	na	Yes
INM		13	=19	16	No
IPSL		10	=9	=10	No
MIROCH		na	=9	na	Yes
MIROCM		7	=3	=3	Yes
MIUB		4	=1	1	na
MPI		3	=13	8	Yes
MRI		11	=3	6	Yes
PCM		17	22	17	No

3 * As summarised in Smith and Chandler (2010) (The performance index is based on the error
 4 variance between modelled and observed climate for 14 climate and ocean variables. ‘Yes’

- 1 indicates the variance error is less than the median across the GCMs.) † na: Not available or
- 2 not applicable. ‡Rank 1 is best rank. # = more than one GCM with this rank.

1 Table 4. Performance statistics comparing CMIP3 GCM mean and standard deviation of annual precipitation,
 2 mean annual temperature, and mean monthly patterns of precipitation and temperature with concurrent
 3 observed data. (Analysis based on untransformed data.).

GCM Name	MAP			SDP			MAT			Monthly pattern	
	R ²	NSE	RMSE	R ²	NSE	RMSE	R ²	NSE	RMSE	NSE	NSE
										Prec	Temp
CCCMA-t47	0.498	0.457	435	0.342	0.252	63	0.984	0.953	3.14	0.409	0.838
CCCMA-t63	0.519	0.458	447	0.397	0.328	65	0.984	0.940	3.59	0.364	0.797
CCSM(1)*	0.496	0.483	460	0.426	0.413	71	0.982	0.981	2.06	-0.178	0.910
CCSM(2)	0.488	0.473	464	0.423	0.411	71	0.982	0.981	2.03	-0.210	0.912
CCSM(3)	0.493	0.479	462	0.418	0.403	71	0.981	0.980	2.08	-0.195	0.908
CCSM(4)	0.500	0.488	457	0.426	0.410	71	0.982	0.980	2.08	-0.174	0.911
CCSM(5)	0.493	0.480	461	0.423	0.410	71	0.983	0.981	2.02	-0.210	0.909
CCSM(6)	0.494	0.480	461	0.437	0.426	70	0.982	0.981	2.04	-0.181	0.909
CCSM(7)	0.496	0.483	460	0.429	0.420	71	0.982	0.981	2.06	-0.173	0.907
CCSM(9)	0.500	0.488	457	0.400	0.393	72	0.982	0.980	2.08	-0.157	0.910
CNRM	0.445	0.246	527	0.479	0.321	65	0.979	0.967	2.67	-0.631	0.879
CSIRO	0.387	0.363	503	0.462	0.452	65	0.971	0.959	2.99	0.034	0.825
GFDL2.0	0.544	0.528	434	0.588	0.460	63	0.980	0.934	3.79	-0.092	0.760
GFDL2.1	0.534	0.518	436	0.570	0.196	77	0.979	0.970	2.54	0.071	0.884
GISS-AOM(1)	0.330	-0.093	624	0.142	0.039	73	0.972	0.969	2.55	-0.325	0.873
GISS-AOM(2)	0.330	-0.087	623	0.132	0.027	74	0.972	0.970	2.54	-0.306	0.876
GISS-EH(1)	0.373	-0.510	692	0.210	-0.397	78	0.963	0.956	3.03	-0.856	0.858
GISS-EH(2)	0.375	-0.502	690	0.176	-0.589	83	0.962	0.955	3.07	-0.920	0.852
GISS-EH(3)	0.368	-0.535	697	0.181	-0.521	81	0.962	0.955	3.06	-0.858	0.856
GISS-ER(1)	0.386	-0.347	653	0.254	-0.115	70	0.970	0.960	2.87	-0.819	0.854

GISS-ER(2)	0.381	-0.357	656	0.203	-0.372	77	0.970	0.959	2.90	-0.739	0.850
GISS-ER(4)	0.386	-0.340	652	0.223	-0.214	72	0.970	0.960	2.88	-0.742	0.854
HadCM3	0.662	0.630	363	0.618	0.572	51	0.988	0.973	2.43	0.227	0.893
HadGEM	0.571	0.302	531	0.457	0.178	82	0.977	0.953	3.22	0.046	0.824
IAP(1)	0.496	0.438	456	0.191	0.096	75	0.963	0.894	4.64	-0.910	0.777
IAP(2)	0.493	0.433	458	0.188	0.041	77	0.962	0.895	4.61	-0.989	0.779
IAP(3)	0.499	0.440	455	0.186	0.048	77	0.963	0.896	4.60	-0.922	0.781
INGV	0.681	0.672	371	0.492	0.468	70	0.983	0.973	2.45	-0.263	0.882
INM	0.450	0.439	431	0.287	0.099	65	0.969	0.952	3.21	-0.247	0.833
IPSL	0.394	0.116	563	0.421	0.223	68	0.967	0.957	3.05	-0.147	0.846
MIROCH	0.588	0.370	514	0.583	0.570	63	0.974	0.971	2.54	0.107	0.906
MIROCM(1)	0.555	0.512	424	0.477	0.454	58	0.970	0.969	2.58	0.061	0.899
MIROCM(2)	0.552	0.508	425	0.525	0.501	56	0.970	0.969	2.58	0.054	0.900
MIROCM(3)	0.549	0.505	427	0.459	0.428	60	0.971	0.970	2.52	0.041	0.902
MIUB(1)	0.689	0.676	336	0.527	0.510	51	0.979	0.960	2.92	0.166	0.870
MIUB(2)	0.684	0.671	338	0.529	0.513	51	0.979	0.962	2.85	0.155	0.867
MIUB(3)	0.691	0.678	335	0.524	0.515	51	0.979	0.958	2.99	0.167	0.860
MPI(1)	0.543	0.538	429	0.464	0.437	66	0.985	0.984	1.88	0.014	0.939
MPI(2)	0.541	0.536	430	0.462	0.415	67	0.985	0.983	1.90	-0.002	0.939
MPI(3)	0.542	0.536	430	0.507	0.479	63	0.986	0.984	1.87	0.007	0.940
MRI(1)	0.617	0.535	414	0.507	0.499	56	0.977	0.969	2.57	0.217	0.912
MRI(2)	0.615	0.537	413	0.513	0.491	56	0.976	0.968	2.64	0.216	0.907
MRI(3)	0.617	0.541	411	0.523	0.505	55	0.977	0.969	2.57	0.222	0.911
MRI(4)	0.619	0.539	412	0.532	0.523	54	0.977	0.969	2.60	0.195	0.911
MRI(5)	0.615	0.538	412	0.503	0.487	56	0.977	0.968	2.62	0.211	0.907
PCM	0.360	0.190	546	0.336	0.135	73	0.975	0.943	3.49	-0.415	0.798

1

* In the paper a parenthesis after a GCM name indicates run number.

1 Table 5. Köppen-Geiger climate classification (adapted from Peel et al., 2007)

Köppen-Geiger class	Description of climate
Af	Tropical, rainforest
Am	Tropical, monsoon
Aw	Tropical, savannah
BWh	Arid, desert hot
BWk	Arid, desert cold
BSh	Arid, steppe hot
BSk	Arid, steppe cold
Csa	Temperate, dry and hot summer
Csb	Temperate, dry and warm summer
Csc	Temperate, dry and cold summer
Cwa	Temperate, dry winter and hot summer
Cwb	Temperate, dry winter and warm summer
Cwc	Temperate, dry winter and cold summer
Cfa	Temperate, without dry season and hot summer
Cfb	Temperate, without dry season and warm summer
Cfc	Temperate, without dry season and cold summer
Dsa	Cold, dry and hot summer
Dsb	Cold, dry and warm summer
Dsc	Cold, dry and cool summer
Dsd	Cold, dry summer and very cold winter
Dwa	Cold, dry winter and hot summer
Dwb	Cold, dry winter and warm summer
Dwc	Cold, dry winter and cool summer
Dwd	Cold, dry winter and very cold winter
Dfa	Cold, without dry season and hot summer
Dfb	Cold, without dry season and warm summer
Dfc	Cold, without dry season and cool summer
Dfd	Cold, without dry season and very cold winter
ET	Polar, tundra
EF	Polar, frost

1 Table 6. Köppen-Geiger climate estimated by MPI(3) compared with the observed Köppen-Geiger climate for (a) the
 2 one-letter and (b) the two-letter climate classification.

(a)		CRU						
	Land Surface	A	B	C	D	E	Sum	
GCM	A	414	19	8	0	0	441	
	B	68	339	52	17	0	476	
	C	24	62	319	27	0	432	
	D	0	76	16	1085	17	1194	
	E	0	6	7	143	121	277	
	Sum	506	502	402	1272	138	2820	

3

(b)		CRU													
	Land Surface	Af	Am	Aw	BW	BS	Cs	Cw	Cf	Ds	Dw	Df	ET	EF	Sum
GCM	Af	57	0	2	0	0	0	0	0	0	0	0	0	0	59
	Am	24	19	13	0	0	0	0	0	0	0	0	0	0	56
	Aw	25	49	225	0	19	0	4	4	0	0	0	0	0	326
	BW	2	1	2	134	50	3	4	0	0	0	2	0	0	198
	BS	4	11	48	50	105	13	19	13	4	0	11	0	0	278
	Cs	0	0	0	10	18	35	9	20	1	0	6	0	0	99
	Cw	0	1	17	0	5	0	62	1	0	1	0	0	0	87
	Cf	2	2	2	3	26	1	35	156	0	0	19	0	0	246
	Ds	0	0	0	0	33	2	1	1	38	1	40	0	0	116
	Dw	0	0	0	0	5	0	1	0	0	102	2	0	0	110
	Df	0	0	0	3	35	0	4	7	2	57	843	17	0	968
	ET	0	0	0	0	6	2	2	3	8	22	113	93	0	249
	EF	0	0	0	0	0	0	0	0	0	0	0	11	17	28
	Sum	114	83	309	200	302	56	141	205	53	183	1036	121	17	2820

1 Table 7. Proportion of CMIP3 GCM grid cells (20C3M) that
 2 reproduce observed CRU Köppen-Geiger climate classification over
 3 the period 1/1950 – 12/1999.

GCM Name	Köppen-Geiger climate class*		
	Three-letter	Two-letter	One-letter
CCCMA-t47	0.498	0.620	0.753
CCCMA-t63	0.429	0.558	0.709
CCSM(1)	0.488	0.558	0.749
CCSM(2)	0.489	0.563	0.748
CCSM(3)	0.424	0.545	0.744
CCSM(4)	0.466	0.549	0.749
CCSM(5)	0.444	0.519	0.727
CCSM(6)	0.490	0.563	0.757
CCSM(7)	0.488	0.556	0.749
CCSM(9)	0.489	0.560	0.755
CNRM	0.539	0.602	0.775
CSIRO	0.506	0.601	0.775
GFDL2.0	0.430	0.563	0.726
GFDL2.1	0.508	0.590	0.781
GISS-AOM(1)	0.460	0.559	0.773
GISS-AOM(2)	0.456	0.561	0.773
GISS-EH(1)	0.407	0.487	0.751
GISS-EH(2)	0.402	0.482	0.741
GISS-EH(3)	0.400	0.473	0.744
GISS-ER(1)	0.426	0.478	0.732
GISS-ER(2)	0.424	0.468	0.722
GISS-ER(4)	0.426	0.478	0.732
HadCM3	0.549	0.624	0.797
HadGEM	0.563	0.676	0.818
IAP(1)	0.362	0.484	0.790
IAP(2)	0.368	0.480	0.784
IAP(3)	0.369	0.490	0.784
INGV	0.495	0.616	0.815
INM	0.452	0.526	0.731
IPSL	0.459	0.544	0.749
MIROCH	0.496	0.631	0.806
MIROCM(1)	0.477	0.597	0.749
MIROCM(2)	0.477	0.594	0.759
MIROCM(3)	0.469	0.583	0.748
MIUB(1)	0.528	0.604	0.783
MIUB(2)	0.528	0.604	0.783
MIUB(3)	0.520	0.610	0.778

MPI(1)	0.599	0.666	0.801
MPI(2)	0.593	0.657	0.805
MPI(3)	0.602	0.669	0.808
MRI(1)	0.534	0.644	0.808
MRI(2)	0.521	0.625	0.798
MRI(3)	0.527	0.632	0.798
MRI(4)	0.528	0.634	0.799
MRI(5)	0.532	0.641	0.803
PCM	0.397	0.481	0.660

1

*The three-, two- and one-letter climate classes are listed in Table 5.

2

1 Table 8. CMIP3 GCM run rank (rank 1 = best) based on Nash-Sutcliffe Efficiency
 2 values from comparison of 20C3M and concurrent observed grid cell data.

GCM Name	MAP rank	SDP rank	MAT rank	Monthly pattern rank*	Rank sum	Overall GCM rank
CCCMA-t47	28	30	38	19	115	12
CCCMA-t63	27	28	42	22	119	13
CCSM(1)	21	22	7	18	68	8
CCSM(2)	26	23	5	17	71	
CCSM(3)	25	26	10	21	82	
CCSM(4)	20	25	11	16	72	
CCSM(5)	24	24	4	21	73	
CCSM(6)	23	19	6	20	68	
CCSM(7)	22	20	8	19.5	69.5	
CCSM(9)	19	27	9	17	72	
CNRM	36	29	26	30.5	121.5	14
CSIRO	34	16	32	28.5	110.5	11
GFDL2.0	14	14	43	34	105	10
GFDL2.1	15	32	15	17.5	79.5	9
GISS-AOM(1)	40	39	20	30.5	129.5	
GISS-AOM(2)	39	40	17	29.5	125.5	15
GISS-EH(1)	45	44	35	35.5	159.5	22
GISS-EH(2)	44	46	37	39	166	
GISS-EH(3)	46	45	36	36.5	163.5	
GISS-ER(1)	42	41	28	36	147	19
GISS-ER(2)	43	43	31	36.5	153.5	
GISS-ER(4)	41	42	29	36	148	
HadCM3	5	1	13	12	31	1
HadGEM	35	33	39	28	135	17
IAP(1)	31	36	46	44	157	
IAP(2)	32	38	45	45	160	
IAP(3)	29	37	44	44	154	21
INGV	3	13	12	28	56	5
INM	30	35	40	35	140	18
IPSL	38	31	34	29.5	132.5	16
MIROCH	33	2	14	14.5	63.5	7
MIROCM(1)	16	15	22	17	70	
MIROCM(2)	17	8	21	17	63	6
MIROCM(3)	18	18	16	17.5	69.5	
MIUB(1)	2	6	30	18.5	56.5	
MIUB(2)	4	5	27	19.5	55.5	4
MIUB(3)	1	4	33	19	57	
MPI(1)	9	17	2	10.5	38.5	
MPI(2)	12	21	3	12	48	

MPI(3)	11	12	1	10.5	34.5	2
MRI(1)	13	9	18	5	45	
MRI(2)	10	10	25	10.5	55.5	
MRI(3)	6	7	19	5.5	37.5	3
MRI(4)	7	3	23	8	41	
MRI(5)	8	11	24	11.5	54.5	
PCM	37	34	41	38.5	150.5	20

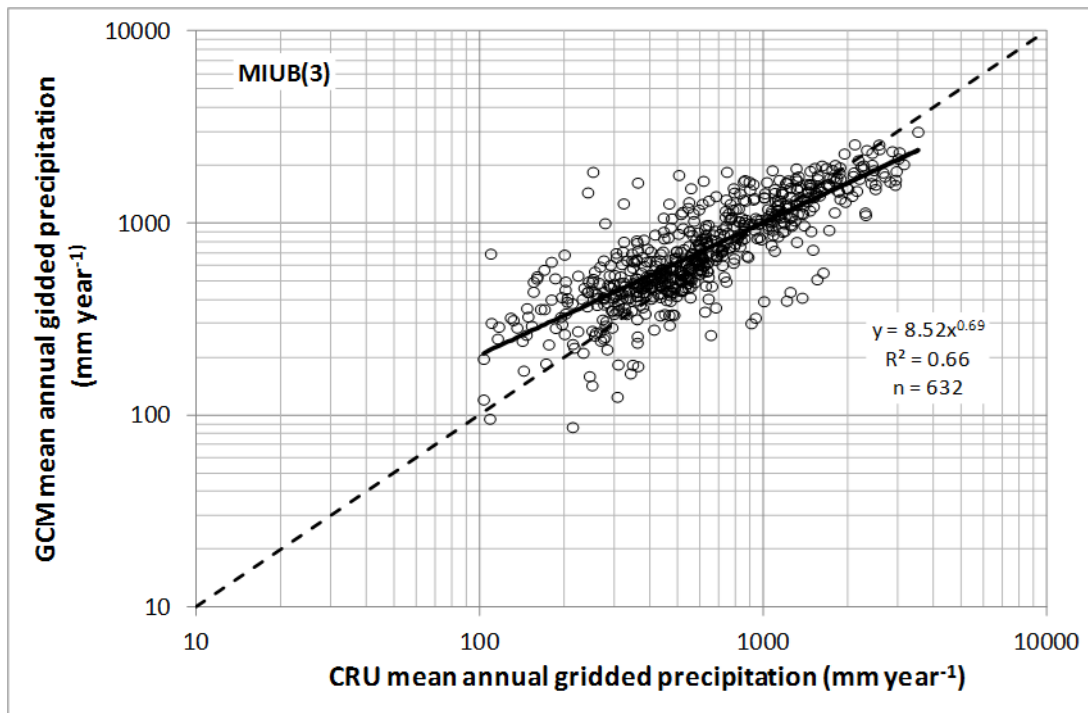
- 1 * Monthly pattern rank is the rank of the average of the monthly pattern NSEs for
- 2 precipitation and temperature
- 3

1 Table 9. Better performing CMIP3 GCMs identified from the literature and
 2 our analyses.

Grid cells (Tables 4 & 8) (Col. 1)	Literature (Table 3) (Col. 2)	Better performing GCMs (Col. 3)
	CCCMA-t47 CCSM GFDL2.0 GFDL2.1	
HadCM3 INGV	HadCM3	HadCM3
	MIROCH*	
MIROCM MIUB MPI MRI	MIROCM MIUB* MPI MRI	MIROCM MIUB MPI MRI

3 *Added to list – see Section 5.3 for explanation

4

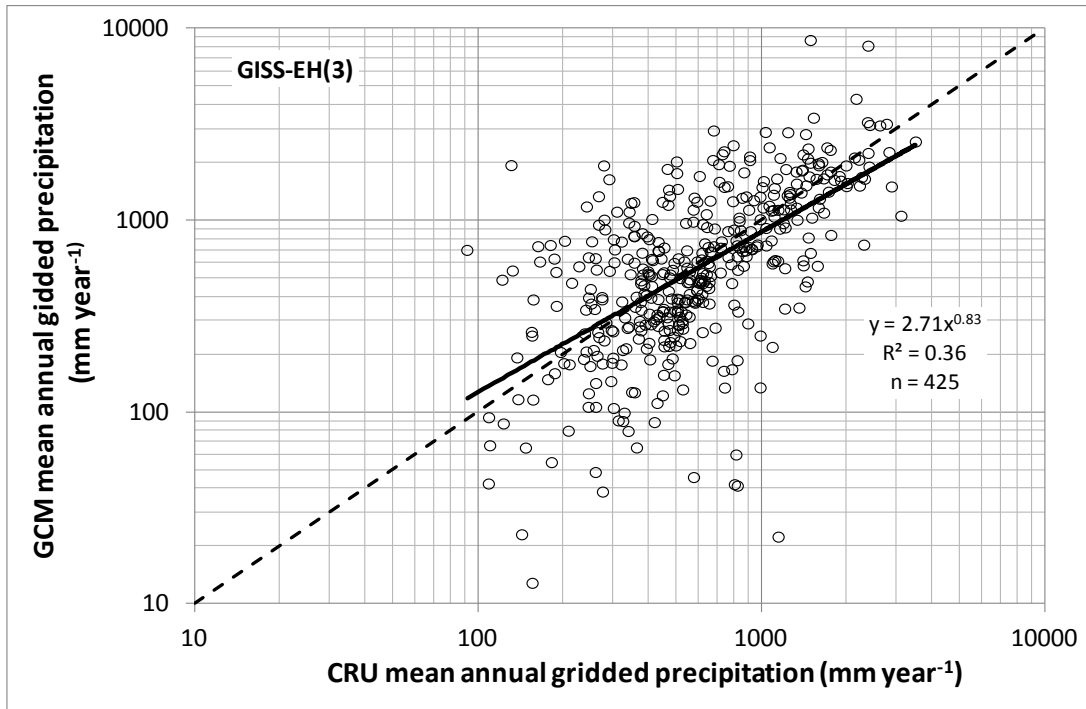


1

2

3 Figure 1. Comparison of MIUB(3) model estimates of observed mean annual precipitation
 4 with CRU estimates. (Based on untransformed precipitation NSE = 0.678, rank 1 of 46 runs,
 5 and $R^2 = 0.691$).

6

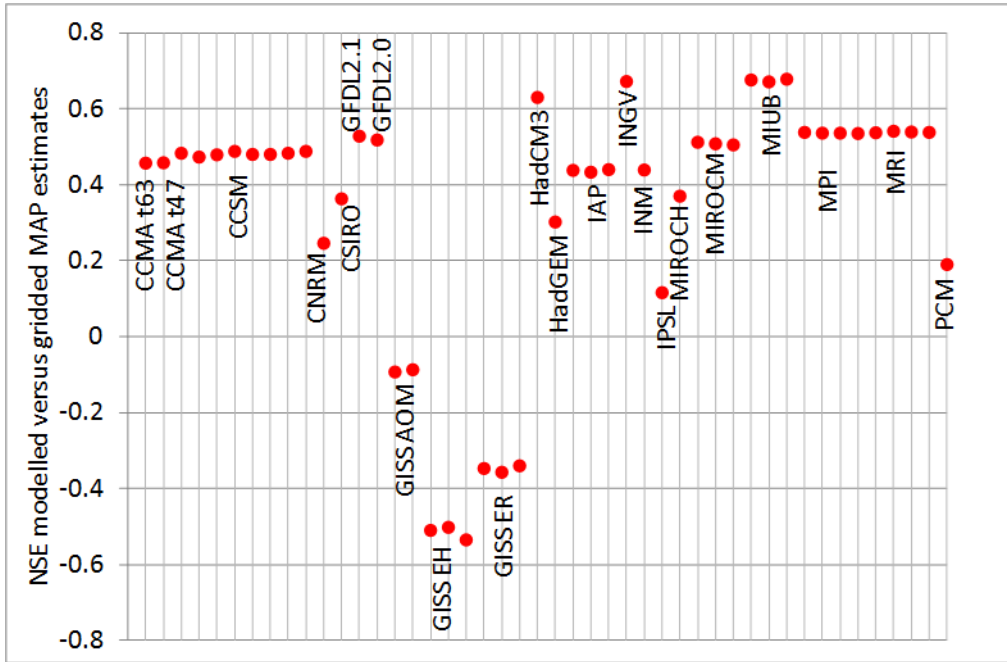


1

2

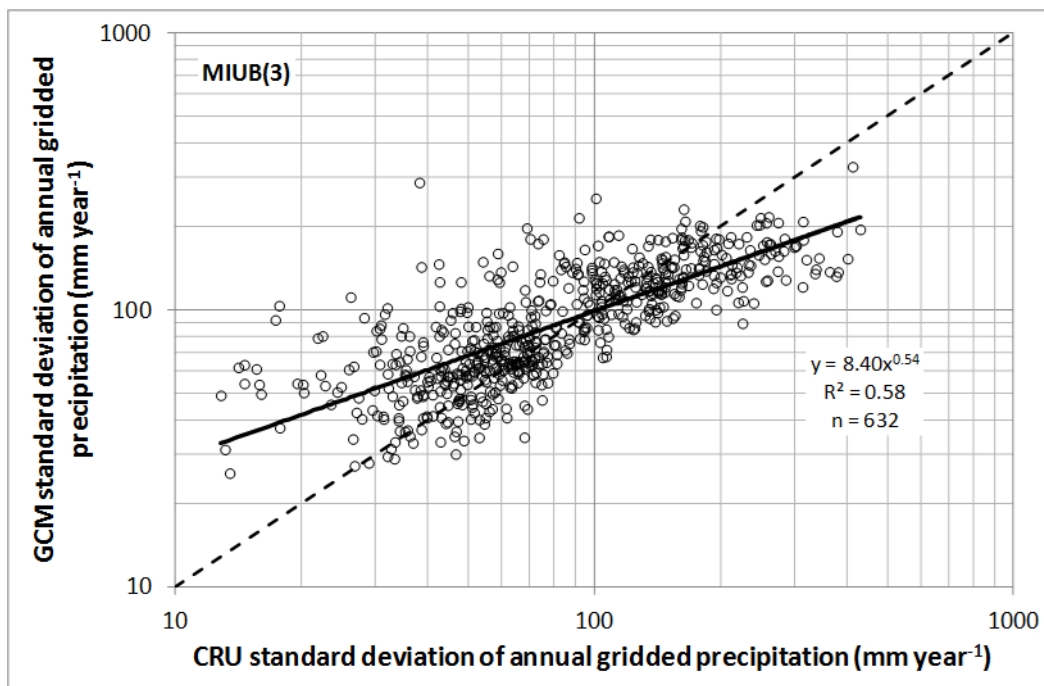
3 Figure 2. Comparison of GISS-EH(3) model estimates of observed mean annual precipitation
4 with CRU estimates. (Based on untransformed precipitation NSE = -0.535, rank 46 of 46
5 runs, and $R^2 = 0.368$).

6



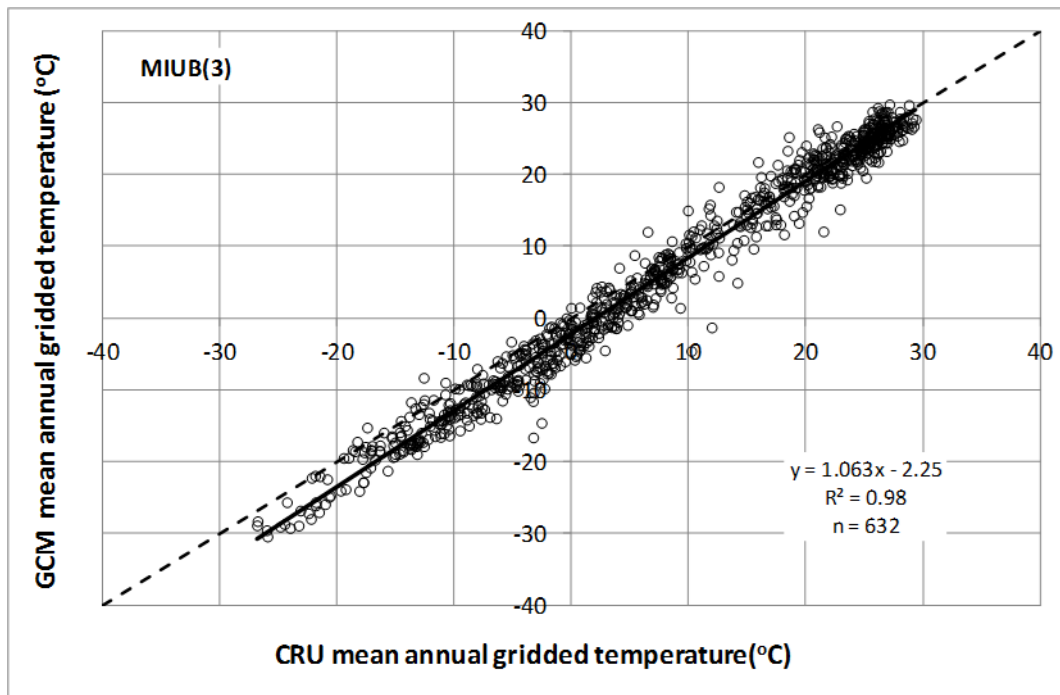
1
2
3
4
5

Figure 3. Nash-Sutcliffe Efficiency (NSE) values for modelled versus observed MAP untransformed estimates for 46 CMIP3 GCM runs



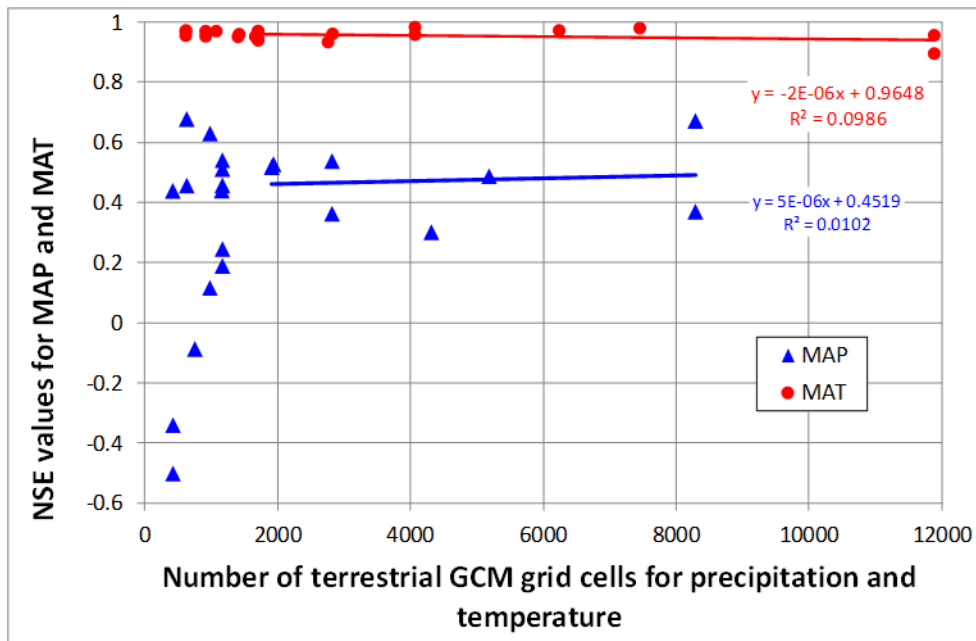
1
 2
 3
 4
 5
 6

Figure 4. Comparison of MIUB(3) model estimates of the standard deviation of annual precipitation with CRU observed estimates. (Based on untransformed precipitation NSE = 0.515, rank 4 of 46 runs, and $R^2 = 0.524$).



1
2
3
4
5
6

Figure 5. Comparison of MIUB(3) model estimates of mean annual temperature with CRU estimates. (Based on untransformed temperature NSE = 0.958, rank 33 of 46 runs, and $R^2 = 0.979$).

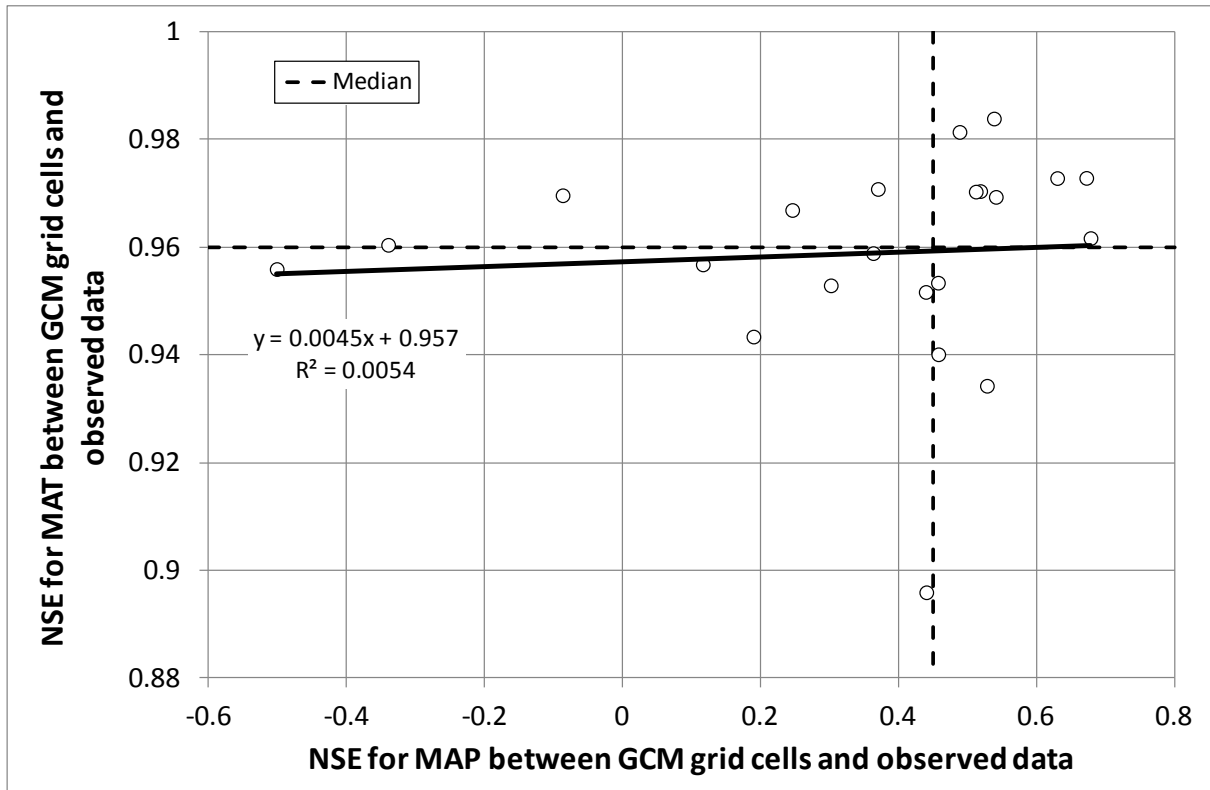


1

2

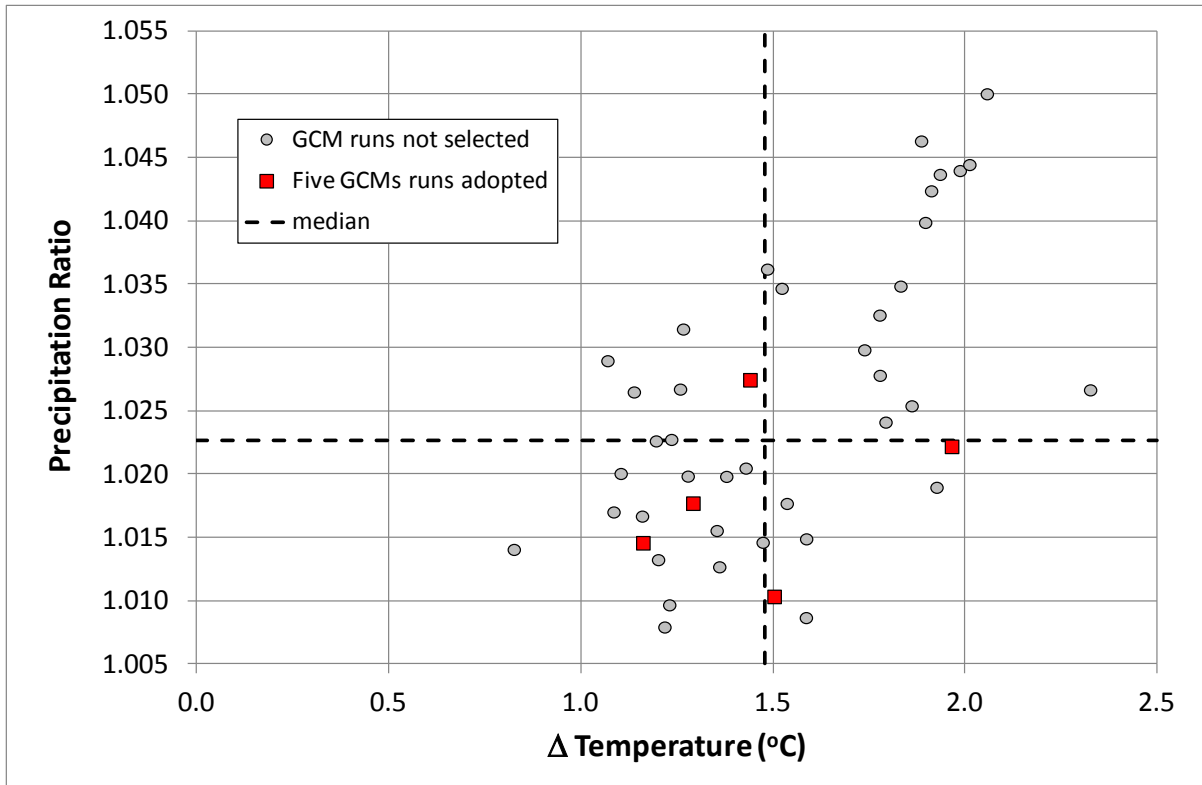
3 Figure 6. Relating 22 CMIP3 GCM resolutions (as the number of terrestrial grid cells for
 4 MAP) to model performance based on Nash-Sutcliffe efficiency for mean annual precipitation
 5 and mean annual temperature. (The trend lines are fitted to data with >1500 grid cells.).

6



1
2
3
4
5
6

Figure 7. Comparison of Nash-Sutcliffe efficiency values between CMIP3 GCM and observed mean annual temperatures with Nash-Sutcliffe efficiency values between CMIP3 GCM and observed mean annual precipitation.



1
 2 Figure 8. Ratio of 2015-2034 to 1965-1994 mean annual precipitation compared with the
 3 change in mean annual temperature ($[2015-2034] - [1965-1994]$) for the selected five CMIP3
 4 GCMs runs compared with the 23 CMIP3 GCMs including all ensemble members for the
 5 global land surface

6
 7

1 **Supplementary Material**

2 **Appendix A: Estimating potential evapotranspiration for climate change impact** 3 **assessments**

4 Projected changes in water and energy at the catchment scale are the fundamental basis of all
5 hydrologic climate change impact assessments. Hydrologic models require time-series of
6 precipitation and, usually, potential evapotranspiration to represent the interaction of water
7 and energy within a catchment. Therefore, for hydrologic climate change impact assessments,
8 an estimate of potential evapotranspiration (PET) is required. For the practitioner the question
9 is which PET method to adopt? Here we briefly review three questions that influence the
10 choice of PET equation: (1) does the equation represent all relevant processes; (2) what PET
11 information does a hydrologic model actually use; and (3) are future projections of variables
12 used to estimate PET reliable?

13 **A.1 Does the PET equation represent all relevant processes?**

14 McMahon et al. (2013) discuss a range of PET equations used in rainfall-runoff modelling.
15 Frequently adopted methods to represent PET include Penman (Penman, 1948), Penman-
16 Monteith (Monteith, 1965), FAO reference crop (Allen et al., 1998), Morton (Morton, 1983)
17 and pan evaporation data. Ideally to represent future PET conditions the method adopted
18 should adequately capture all changes in the energy and aerodynamic components of the
19 evaporative process.

20 The potential danger of using a PET equation that does not adequately represent all relevant
21 processes is highlighted by recent trends in pan evaporation data. Over the past several
22 decades the magnitude of evaporation from Class-A pans has decreased (between -1 to -4 mm
23 year⁻²) while at the same time annual temperatures have risen (Roderick et al., 2009a).
24 Roderick et al. (2009b) warn against using temperature only PET estimates for climate change
25 studies as they would suggest that rising temperature would lead to rising evaporative
26 demand; the opposite of what has been observed from pan data recently. Roderick et al.
27 (2009b) attribute much of the observed decline in pan evaporation to declines in radiation
28 and/or wind speed. Donohue et al. (2010), using the Penman formulation and gridded
29 Australian data (1981-2006), attributed increasing surface temperature with contributing +1.5

1 mm year⁻² toward evaporative demand. However, the temperature contribution was more than
2 offset by negative contributions from changes in wind speed (-1.3 mm year⁻²), net radiation (-
3 0.6 mm year⁻²) and actual vapour pressure (-0.4 mm year⁻²) to give an overall decrease in
4 evaporative demand of -0.8 mm year⁻². Donohue et al. (2010) also compared the performance
5 of five formulations of differing complexity namely Thornthwaite (Thornthwaite (1948),
6 Priestly-Taylor (Priestley and Taylor, 1972), Morton point and areal (Wang et al., 2001) and
7 Penman (1948)) and preferred Penman, the most complex form, based on its ability to best
8 capture the dynamics of evaporative demand. Overall, Roderick et al. (2009a, 2009b), Chen et
9 al. (2005) and Hobbins et al. (2008) conclude that PET estimates based only on T are
10 problematic, particularly in energy limited environments (cold and polar climates), for climate
11 change studies.

12 **A.2 What PET information does a conceptual hydrologic model actually use?**

13 Whether conceptual hydrologic models require, or make use of, detailed PET data was
14 assessed by Andréassian et al. (2004) and Oudin et al. (2005a, 2005b). They found that
15 hydrologic models perform as well (if not better) when calibrated with mean monthly
16 estimates of PET, or with temperature based estimates of PET, rather than time varying
17 estimates of PET or more complex Penman based PET (Penman, 1948, Allen et al., 1998).
18 Catchments used in their studies were located in France (Andréassian et al., 2004; Oudin et
19 al., 2005a, 2005b), USA (Oudin et al., 2005a, 2005b) and Australia (Oudin et al., 2005a,
20 2005b). The vast majority of their catchments have a temperate climate (not strongly water or
21 energy limited on an annual basis). Under these conditions the hydrologic models appear to
22 be largely insensitive to the complexity of the PET data used to drive them. During calibration
23 conceptual hydrologic models are flexible enough to extract the PET information they need
24 from whichever PET data (simple or complex) are used (see Chapman, 2003). Thus, as long
25 as PET estimates are broadly correct in terms of seasonal pattern and annual mean and the
26 hydrologic model was calibrated on that PET data then model performance is likely to be
27 acceptable. For example, Oudin et al. (2005b) tested 27 PET formulations, of varying
28 complexity, over 308 catchments using four daily conceptual models and proposed a simple
29 temperature (mean daily temperature for a given Day-of-Year) and extra-terrestrial radiation
30 (estimated from latitude and Day-of-Year) method that performed as well as the daily Penman

1 method. In summary, a complex estimate of PET is not necessary for successful hydrologic
2 modelling in catchments that are not strongly water or energy limited on an annual basis.

3 **A.3 Are future projections of variables used to estimate PET reliable?**

4 In the previous two sections we have seen that a simple PET formulation may be good enough
5 for hydrologic modelling, but not good enough to represent projected changes in PET. The
6 final question relates to whether GCMs are able to provide reliable outputs on which to base a
7 complex estimate of PET? Kay and Davies (2008) used IPCC third assessment report runs for
8 5 GCMs and 8 regional climate models nested within the Hadley Centre GCM to calculate
9 PET using Penman-Monteith and the temperature/radiation (T/R) method of Oudin et al.
10 (2005b). They compared their two PET estimates derived from GCM data against observation
11 based gridded values of Penman-Monteith PET for Britain. Overall, the GCM estimate of
12 PET using T/R performed better than GCM Penman-Monteith at reproducing observed
13 Penman-Monteith for all climate models. Future values of PET based on Penman-Monteith
14 were also more variable than those based on T/R, which they suggest may reflect reliability
15 issues with GCM variables, other than temperature, used to estimate Penman-Monteith.
16 Kingston et al. (2009) also highlight reliability issues with GCM inputs to the Penman-
17 Monteith equation. Although confidence in GCM-simulated temperature is generally high,
18 Kingston et al. (2009, page 4) note “less confidence can be placed in cloud cover and vapour
19 pressure”, which influence GCM-simulation estimates of net radiation at the evaporating
20 surface and relative humidity. Overall, Kay and Davies (2008) suggest hydrologic modellers
21 should be pragmatic and use as many GCMs as possible and estimate PET in a consistent way
22 for any impact analysis.

23 **A.4 Discussion and summary**

24 Ideally, estimates of PET should be based on methodologies that include all key evaporative
25 processes to ensure future changes in PET are accurately represented. A Penman based
26 equation is thus an ideal methodology to adopt. However, the reliability of future PET
27 estimates is dependent on the reliability of GCM projections of input variables. For example,
28 the Penman equation requires inputs of air temperature, net radiation at the evaporating
29 surface, wind speed and relative humidity. In this paper we have found that mean monthly

1 and mean annual temperature are well reproduced by CMIP3 GCMs. However, reported
2 confidence in GCM estimates of net radiation at the evaporating surface, wind speed and
3 relative humidity is much lower. For example, Johnson and Sharma (2009) have shown that in
4 terms of their Variable Convergence Score (VCS, scaled between 0 and 100, where 100 is
5 perfect convergence between GCMs) the predictions of the surface wind and specific
6 humidity have VCS scores of approximately 40, net longwave radiation about 20 compared
7 with surface temperature and net shortwave radiation of about 70 and precipitation at 10.
8 Therefore, although Penman based methodologies have the capacity to represent future trends
9 due to changes in all key evaporative processes, GCM projections of those process variables,
10 other than temperature, may be unrealistic. Thus at this time PET based on Penman may
11 actually increase uncertainty in future PET, as seen in Kay and Davies (2008). PET based on
12 Penman will be preferable once GCM projections of net radiation at the evaporating surface,
13 wind speed and relative humidity become more reliable.

14 As GCM projections of temperature are considered reliable, here we adopt temperature as a
15 surrogate for PET. Such an approach is likely to provide sufficient PET information for
16 successful hydrologic modelling if the model is calibrated on that data. However, by adopting
17 this approach we acknowledge that the projected trend in PET will be an increase, when in
18 reality the trend may increase or decrease due to changes in temperature, net radiation at the
19 evaporating surface, wind speed and/or relative humidity. We note the error in PET trend is
20 unlikely to be important for hydrologic modelling of water limited catchments, where changes
21 in precipitation are the main driver of changes in runoff. However, in energy limited
22 catchments, PET is a key driver of runoff and errors in PET trend will result in errors in
23 runoff trend.

24 **A.5 References**

25 Allen, R. G., Pereira, L. S., Raes, D., and Smith, M.: Crop evapotranspiration Guidelines for
26 computing crop water requirements FAO Irrigation and Drainage Paper 56. Food and
27 Agriculture Organization of the United Nations, 1998.

1 Andréassian, V., Perrin, C., and Michel, C.: Impact of imperfect potential evapotranspiration
2 knowledge on the efficiency and parameters of watershed models, *J. Hydrol.*, 286, 19-35,
3 2004.

4 Chapman, T. G.: Estimation of evaporation in rainfall-runoff models. MODSIM 2003
5 International Congress on Modelling and Simulation, Townsville, Australia, 2003.

6 Chen, D., Gao G., Xu C-Y., Guo, J., and Ren, G.: Comparison of the Thornthwaite method
7 and pan data with the standard Penman-Monteith estimates of reference evapotranspiration in
8 China, *Clim. Res.*, 28, 123-132, 2005.

9 Donohue, R. J., McVicar, T. R., and Roderick, M. L.: Assessing the ability of potential
10 evaporation formulations to capture the dynamics in evaporative demand within a changing
11 climate, *J. Hydrol.*, 386, 186-197, 2010.

12 Hobbins, M. T., Dai, A., Roderick, M. L., and Farquhar, G. D.: Revisiting the
13 parameterization of potential evaporation as a driver of long-term water balance trends,
14 *Geophys. Res. Lett.*, 35, L12403, doi:10.1029/2008GL033840, 2008.

15 Johnson, F. M. and Sharma, A.: Measurement of GCM skill in predicting variables relevant
16 for hydroclimatological assessments, *Journal of Climate*, 22, 4373-4382, 2009b.

17 Kay, A. L., and Davies, H. N.: Calculating potential evaporation from climate model data: A
18 source of uncertainty for hydrological climate change impacts, *J. Hydrol.*, 358, 221-239,
19 2008.

20 Kingston, D. G, Todd, M. C., Taylor, R. G., Thompson, J. R., and Arnell, N. W.: Uncertainty
21 in the estimation of potential evapotranspiration under climate change, *Geophys. Res. Lett.*,
22 36, L20403, 2009.

23 McMahon, T. A., Peel, M. C., Lowe, L., Srikanthan, R., and McVicar, T. R.: Estimating
24 actual, potential, reference crop and pan evaporation using standard meteorological data: a
25 pragmatic synthesis, *Hydrol. Earth Syst. Sci.*, 17, 1331-1363, 2013.

26 Monteith, J. L.: Evaporation and environment, In Fogg, G.E. (ed), *The state and movement of*
27 *water in living organisms*, Symposium Society Experimental Biology, 19, 205-234.,
28 Cambridge University Press, London, 1965.

- 1 Morton, F. I.: Operational estimates of areal evapotranspiration and their significance to the
2 science and practice of hydrology, *J. Hydrol.*, 66, 1–76, 1983.
- 3 Oudin, L., Michel, C., and Anctil, F.: Which potential evapotranspiration input for a lumped
4 rainfall-runoff model? Part 1 – Can rainfall-runoff models effectively handle detailed
5 potential evapotranspiration inputs? *J. Hydrol.*, 303, 275-289, 2005a.
- 6 Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andréassian, V., Anctil, F., and Loumagne, C.:
7 Which potential evapotranspiration input for a lumped rainfall-runoff model? Part 2 –
8 Towards a simple and efficient potential evapotranspiration model for rainfall–runoff
9 modelling, *J. Hydrol.*, 303, 290-306, 2005b.
- 10 Penman, H. L.: Natural evaporation from open water, bare soil and grass, *Proceedings Royal*
11 *Society London, A*, 193, 120-145, 1948.
- 12 Priestley, C. H. B., and Taylor, R.J.: On the assessment of surface heat flux and evaporation
13 using large scale parameters, *Mon. Weather Rev.*, 100, 81-92, 1972.
- 14 Roderick, M. L., Hobbins, M. T., and Farquhar, G.D.: Pan Evaporation Trends and the
15 Terrestrial Water Balance. I. Principles and Observations, *Geography Compass*, 3(2), 746-
16 760, 2009a.
- 17 Roderick, M. L., Hobbins, M. T., and Farquhar, G. D.: Pan Evaporation Trends and the
18 Terrestrial Water Balance. II. Energy Balance and Interpretation, *Geography Compass*, 3(2),
19 761-780, 2009b.
- 20 Thornthwaite, C. W.: An approach toward a rational classification of climate, *Geogr. Rev.*,
21 38, 55-94, 1948.
- 22 Wang, Q. J., Chiew, F. H. S., McConachy, F. L. N., James, R., de Hoedt, G. C., and Wright,
23 W. J.: *Climatic Atlas of Australia Evapotranspiration*. Bureau of Meteorology,
24 Commonwealth of Australia, 2001.

25

26