1
2
3
4
5
6
7
8
9 **On the Effective Number of Climate Models**
10

11 *Christopher Pennell and Thomas Reichler*
12
13
14 Department of Meteorology, University of Utah, Salt Lake City
15

16

17 Correspondence: *Thomas Reichler* (thomas.reichler@utah.edu)

18

19 Department of Meteorology, University of Utah
20 135 S 1460 E, Rm 819 (WBB)
21 Salt Lake City, UT 84112-0110
22 801-585- 0040 Fax: 801-581-4362
23

24

25

26

27

28

29

30

31 **Abstract**

32 Climate models are essential tools for assessing future climate change. In making predictions, it

33 is beneficial to examine simulations from various models that are developed at centers around the

34 world. The simple average over an ensemble of such models is often taken as the optimal

35 prediction, which in studies of current climate is demonstrated as being more accurate than

36 relying on any one individual model realization. However, this is only true to the extent that

37 different models provide statistically independent information. Here, we examine the ability of

38 current-generation models in simulating the observed present-day mean climate and show that

39 similarities in model implementation play an important role in ensemble estimation. We

40 demonstrate that the effective number of models is considerably smaller than the actual number

41 comprising the ensemble. Our results suggest that the common practice of taking simple

42 ensemble averages needs to be reconsidered.

43

## 1. Introduction

Over two dozen different climate models contribute to the ongoing mission of the Intergovernmental Panel on Climate Change (IPCC), whose aim is to provide reliable estimates of future climate change. Most findings of the IPCC's most recent 4th assessment report are based on simple averages over individual simulations produced by these models [*IPCC*, 2007]. This type of multi-model averaging improves a prediction only to the extent that different model outcomes are randomly distributed around the true future state, or in other words, independent from each other [*Abramowitz and Gupta*, 2008]. If this assumption is not met, resulting predictions are likely to be systematically biased and consequently, inaccurate.

There is reason to believe that the current generation of climate models considered by the IPCC violates the assumption of independence because these models have similar weaknesses [*Reichler and Kim*, 2007]. The probable reason is that models often share physical parameterization schemes and, at times, even large parts of the same code [*Pincus et al.*, 2008]. On the other hand, we may expect that effects from inter-model similarity could potentially be nullified over lengthy run-times. Even minor differences, for example in how small-scale processes are treated and sensible forcings are chosen [*Knutti*, 2008], could vastly amplify due to systemic non-linearities. Although this issue likely hampers the accuracy of individual simulations, its chaotic nature could lead to ideal simulation diversity within an ensemble.

With this in mind, the goal of this study is to explore the impact of model similarity in the context of a multi-model ensemble. We accomplish this by quantifying how well the current generation of models simulate present-day mean climate (section 2). Then, we examine the similarities in model deficiencies amongst the ensemble and clarify the source of these commonalities (section 3). Next, we discuss the potential impacts on current strategies for

67     ensemble prediction (section 4). And finally, we discuss intricate details relating to our

68     methodology (section 5).

69     **2. The Effective Number of Models**

70     We analyze deficiencies in 24 current-generation climate models from the 20$^{th}$ century

71     experiment of the WCRP CMIP3 archive [*Meehl et al.*, 2007]. One model (BCC-CM1) is not

72     included in our analysis since many of the atmospheric quantities used in this study are not

73     provided for this particular model. We proceed by comparing mean climatologies for simulations

74     as well as observations by calculating normalized RMS errors over the northern hemisphere for

75     37 physical and dynamical quantities (Table 1). These quantities are chosen based on the

76     availability of suitable observations, as well as standard practices in climate model validation.

77     Further data processing provides error distributions that are largely symmetric about their

78     respective means, giving us confidence in the correct interpretation of our subsequent analysis of

79     correlation coefficients. The results from these procedures provide for each model, quantity-

80     specific scores relative to that model's mean performance. These scores collectively define a

81     model's error structure (see Detailed Methodology).

82                                          <Table 1 about here>

83     These error structures together with our concept of an effective number of models are used to

84     assess the amount of shared bias contained in the ensemble. The effective number of models

85     ($N_{eff}$) is defined in the following way: $N_{eff}$ equals to one if an ensemble consists of completely

86     correlated error structures since the model members have identical deficiencies; alternatively, if

87     all error structures are uncorrelated, then $N_{eff}$ equals the actual number of models ($N$) in the

88    ensemble. Although our concept may appear somewhat novel, it is similar to ideas such as the

89    effective degrees of freedom or the effective sample size explored in the literature [*Wang and*

90    *Shen*, 1999].

91    We estimate $N_{eff}$ using two different methods: Method I incorporates an inverse technique based

92    on the probability distribution of correlation coefficients [*van den Dool*, 2007] and Method II is

93    based on an eigenanalysis of model correlations [*Bretherton*, 1999] (see Detailed Methodology).

94    Testing our two methods on two artificial experiments produces reliable results.

95    We next apply these methods to the error structures determined from the 24 CMIP3 models. In

96    order to quantify the similarity within the ensemble, we calculate $N_{eff}$ for an increasing number of

97    $N$ models, ranging between 3 and 24. More specifically, we make robust estimates of $N_{eff}$ by

98    applying our two methods to 10,000 ensembles consisting of $N$ randomly selected models

99    (bootstrap without replacement).

100   Both methods indicate that as the number of models increases within an ensemble, the amount of

101   shared bias also increases (Fig. 1). When all 24 models are eventually collected, the decrease in

102   $N_{eff}$ suggests that effectively only 12 to 16 models actually exist in the ensemble. This

103   corresponds to a 33 to 50% reduction, which demonstrates that the error structures of the

104   ensemble's members are correlated at a level beyond what would be expected purely by chance.

105                                            <Figure 1 about here>

106   **3. Examining Model Commonalities**

107    In order to explore the basis behind inter-model similarities in our ensemble, we group models at

108    different levels according to the strength of the relationships between their error structures.

109    Specifically, we use an agglomerative clustering method which operates on the pair-wise

110    distances calculated among the 24 models used in our analysis. Here, the distance between two

111    models is simply defined as ($1-r$), where $r$ represents the correlation coefficient between the

112    models' respective error structures. The outcome of our cluster analysis is depicted by the

113    dendrogram shown in Fig. 2. Other distance metrics and methodologies (not shown) produce

114    similar inter-model relationships.


115                                                         &lt;Figure 2 about here&gt;


116    Since models from the same center tend to differ little in terms of their implementation

117    [*Delworth et al.*, 2006; *Schmidt et al.*, 2005; *Hasumi and Emori*, 2004], it is reasonable to

118    assume that the large amount of bias seen in the ensemble is due to similarities between these

119    same center models. Fig. 2 demonstrates that error structures calculated from models developed

120    at the same center do indeed tend to be quite similar. For instance, the two CGCM3.1 models

121    developed at the Canadian Centre for Climate Modeling and Analysis have error structures that

122    are highly correlated ($r = 86\%$). Similarly close relationships are seen in GISS-ER and GISS-EH,

123    MIROC3.2(medres) and MIROC3.2(hires), as well as GFDL-CM2.0 and GFDL-CM2.1. Also, it

124    is worth noting that the two GFDL models appear the most distinct from the other remaining

125    models. This is shown in Fig. 2 by how these two models collectively merge at a rather long

126    distance with other models.


127    Similarities between same center models alone, however, can only partially explain the reduction

128    of $N_{eff}$ seen in Fig. 1. For example, if we remove seven specific models, leaving each center

129   represented only once in the ensemble ($N = 17$ in this case), we still find that $N_{eff}$ is between 24

130   and 35% smaller than the full ensemble. The actual values of $N_{eff}$, in this instance, are indicated

131   by the symbols corresponding to the two methods in Fig. 1. Given the remaining disparity

132   between $N$ and $N_{eff}$, we conclude that there must be similarities inherent in models across

133   different centers as well.


134   As outlined before, we arrived at the above results by examining model error structures for the

135   northern hemisphere. Examining error structures for the tropics (30°S – 30°N) and southern

136   hemisphere (90°S – 30°S) leads to very similar conclusions. Over these two regions, the decrease

137   in $N_{eff}$ even exceeds that seen over the northern hemisphere by about one model (not shown). As

138   before, same-center models tend to exhibit strong commonalities in the two regions, except for

139   the two GFDL models over the southern hemisphere. This somewhat surprising outcome may be

140   related to the large differences between the two models in simulating temperature and salinity of

141   the southern ocean [*Gnanadesikan et al.*, 2006], which in turn may feedback into the atmospheric

142   simulations over that region.


143   **4. Conclusion**


144   To summarize, for each of 24 CMIP3 models, we calculate errors in simulating present-day

145   climatological mean-fields for 37 different atmospheric quantities. We use two methods that

146   quantify the amount of inter-model similarities in these errors as it relates to the number of

147   models in a current-generation climate ensemble. As the number of models in an ensemble

148   increase, we see that the disparity between the number of models and effective number of models

149   increases as well. In a full 24 member ensemble, we find that there only effectively exist about

150    12 to 16 models. To explore the reasoning behind this reduction, we use clustering analysis to

151    group models based on similarities in their error characteristics. We see that a portion of inter-

152    model similarity can readily be explained by models developed at the same center being included

153    in the CMIP3 archive. This may not be too surprising since models from the same center often

154    share a considerable amount of code. Commonalities in model implementation, however, are

155    also seen to exist across all models. This can partially be explained by the CMIP3 archive being

156    an 'ensemble of opportunity' [*Tebaldi and Knutti*, 2007]. As opposed to utilizing sound sampling

157    design for model selection, results are essentially accepted from any center willing to participate

158    in the archive. The distribution of model simulations belonging to such an ensemble is

159    unpredictable and likely includes shared biases.


160    That the number of effective climate models is considerably less than the actual size of the

161    CMIP3 ensemble suggests that simple arithmetic averages over different models simulations can

162    give spurious confidence in a prediction. In order to produce more reliable estimates of future

163    climate change it may be necessary to refine strategies for selecting and weighting the members

164    of multi-model ensembles. Concerns about the effectiveness of simple multi-model averaging

165    have led to some recent alternatives. Perturbed physics ensembles, for instance, sample a broad

166    range of parametric uncertainty usually not explored by modeling centers and weight individual

167    "model versions" based on their skill [*Murphy et al.*, 2004]. This approach has currently been

168    attempted using only individual models, however, as incorporating multiple models is

169    computationally prohibitive. Even modern probabilistic approaches, which consider different

170    models simultaneously, typically require an assumption of model independence in order to

171    produce tractable results [*Furrer et al.*, 2006; *Tebaldi et al.*, 2005]. Still, recent evidence

172     suggests that weighted averages based on model skill show promise in improving ensemble

173     prediction [*Min and Hense*, 2006; *Murphy et al.*, 2004].


174     Simply constructing unbiased estimates does not, of course, guarantee predictive accuracy. And

175     given the relatively modest number of models included in the CMIP3 ensemble, defining

176     reasonable sampling strategies seems difficult at best. In light of these findings, it is apparent that

177     the underlying processes which give rise to multi-model bias should be better understood.

178     Quantifying the amount of inter-model similarities, in terms of an effective number of models, is

179     a step toward intelligently weighting redundant ensemble members and may benefit future work

180     in multi-model climate prediction.


181     **5. Detailed Methodology**


182     We evaluate a model's performance in simulating present-day climate by first calculating

183     normalized RMS errors for each climate quantity as


184

$$E^2 = \frac{1}{K} \sum_{n=1}^{K} w_n (o_n - s_n)^2 / \sigma_n^2 \tag{1}.$$

186


187     Here, $(o_n - s_n)$ represents the difference between an observational and model simulated field at

188     grid-point $n$, with $K$ total grid-points pertaining to specific large regions. $w_n$ provides proper

189     spatial and vertical mass weighting, while $\sigma_n^2$ denotes the interannual variance taken from

190     observations at $n$. Identical methodology has recently been applied in the literature [*Reichler and*

191     *Kim*, 2007].

192    By conducting a logarithmic transformation of these errors, we ensure symmetric numerical

193    distributions. For each model, we subtract its mean error so that errors are relative only to a

194    model's overall performance. Examining statistical moments via testing of the null multivariate

195    normal hypothesis [*Wilks*, 2006] provides acceptable evidence that errors are now essentially

196    normally distributed (p-values of 0.995 and 0.527 for skew and kurtosis respectively).

197    Method I, for calculating the number of effective models, employs an inverse procedure based on

198    analytical properties of the correlation coefficient distribution [*van den Dool*, 2007]. If two

199    independent variables are Gaussian distributed, then their correlation coefficient $r$ is Gaussian

200    distributed with zero mean and variance $1/(N - 1)$ [*Bain and Engelhardt*, 1992]. $N_{eff}$ is then

201    estimated by equating the sample variance of quantity correlations $S_r^2$ with the expected

202    population variance as

203

204
$$S_r^2 = \frac{1}{N_{eff} - 1}$$
(2).

205

206    Given correlation coefficients are symmetric about zero, larger similarities amongst models will

207    subsequently result in larger sample variability thereby reducing $N_{eff}$.

208    For Method II, we consider the eigenvalues that result from an eigenanalysis of the model error

209    structure correlation matrix [*Bretherton et al.*, 1999]. $N_{eff}$ can then be calculated as

210

211
$$N_{eff} = \frac{(\sum_{i=1}^{N} \lambda_i)^2}{\sum_{i=1}^{N} \lambda_i^2}$$
(3).

212

213     Here, $\lambda_i$ is the $i^{th}$ eigenvalue and $N$ is the actual number of models. If error structures are

214     independent, then all eigenvalues will have the same value and $N_{eff} = N$. However, if all error

215     structures are identical, then there will exist only one non-zero eigenvalue and $N_{eff} = 1$. Here, $N_{eff}$

216     is bounded inclusively between one and the number of models $N$.

217

**References**

Abramowitz, G., Gupta, H. Towards a model space and independence metric. *Geophys. Res.*

*Lett.* 35, L05705 (2008)

Bain, L. J., Engelhardt, M. *Introduction to Probability and Mathematical Statistics* 2nd edn

(Brooks/Cole, 1992)

Bretherton, C. S. et al. The effective number of spatial degrees of freedom of a time-varying

field. *J. Clim.* 12, 1990 (1999)

Delworth, T.L. et al. GFDL's CM2 global coupled climate models -- Part 1:  Formulation and

simulation characteristics. J. Climate. 19, 643 (2006)

Furrer, R. et al. Spatial patterns of probabilistic temperature change projections from a

multivariate Bayesian analysis. *Geophys. Res. Lett.* 34, L06711 (2007)

Gnanadesikan, A. et al. GFDL's CM2 Global Coupled Climate Models. Part II: The Baseline

Ocean Simulation. *J. Clim.* 19, 675 (2006)

Hasumi, H., Emori, S. *K-1 coupled GCM (MIROC) description, K-1 Technical Report No. 1*

(Center for Climate System Research, University of Tokyo, 2004) [Available online at:

http://www.ccsr.u-tokyo.ac.jp/kyosei/hasumi/MIROC/tech-repo.pdf as of October 17, 2008]

Knutti, R. Why are climate models reproducing the observed global surface warming so well?

*Geophys. Res. Lett.* 35, L18704 (2008)

Meehl, G. A. et al. *Global Climate Projections. In: Climate Change 2007: The Physical Science*

*Basis. Contribution of Working Group I to the Fourth Assessment Report of the*

*Intergovernmental Panel on Climate Change* (Cambridge Univ. Press, Cambridge, New York,

2007)

246    Meehl, G. A. et al. The WCRP CMIP3 multimodel dataset: A new era in climate change

247    research. *Bull. Amer. Meteor. Soc.* 88, 1383 (2007)

248    Min, S.-K., Hense, A. A Bayesian approach to climate model evaluation and multi-model

249    averaging with an application to global mean surface temperatures from IPCC AR4 coupled

250    climate models. *Geophys. Res. Lett.* 33, L08708 (2006)

251    Murphy, J. et al. Quantification of modelling uncertainties in a large ensemble of climate change

252    simulations. *Nature* 429, 768 (2004)

253    Pincus, R. et al. Evaluating the present-day simulation of clouds, precipitation, and radiation in

254    climate models. *J. Geophys. Res.* 113, D14209 (2008)

255    Reichler, T., Kim, J. How well do coupled models simulate today's climate? *Bull. Am. Meteorol.*

256    *Soc.* 89, 303 (2007)

257    Schmidt, G. A. et al. Present day atmospheric simulations using GISS ModelE: Comparison to

258    in-situ, satellite and reanalysis data. *J. Climate* 19, 153 (2005)

259    Tebaldi, C. et al. Quantifying uncertainty in projections of regional climate change: a Bayesian

260    approach to the analysis of multimodel ensembles. *J. Climate* 18(10), 1524 (2005)

261    Tebaldi, C., Knutti, R. The Use of the Multi-Model Ensemble in Probabilistic Climate

262    Projections. *Phil. Trans. R. Soc. A.* 2053, 365 (2007)

263    van den Dool, H. *Emperical Methods in Short-Term Climate Prediction* (Oxford Univ. Press,

264    New York, 2007)

265    Wang, X., Shen, S. S. Estimation of spatial degrees of freedom of a climate field. *J. Clim.* 12,

266    1280 (1999)

267    Wilks, D. S. *Statistical Methods in the Atmospheric Sciences* 2nd edn (Elsevier, 2006)
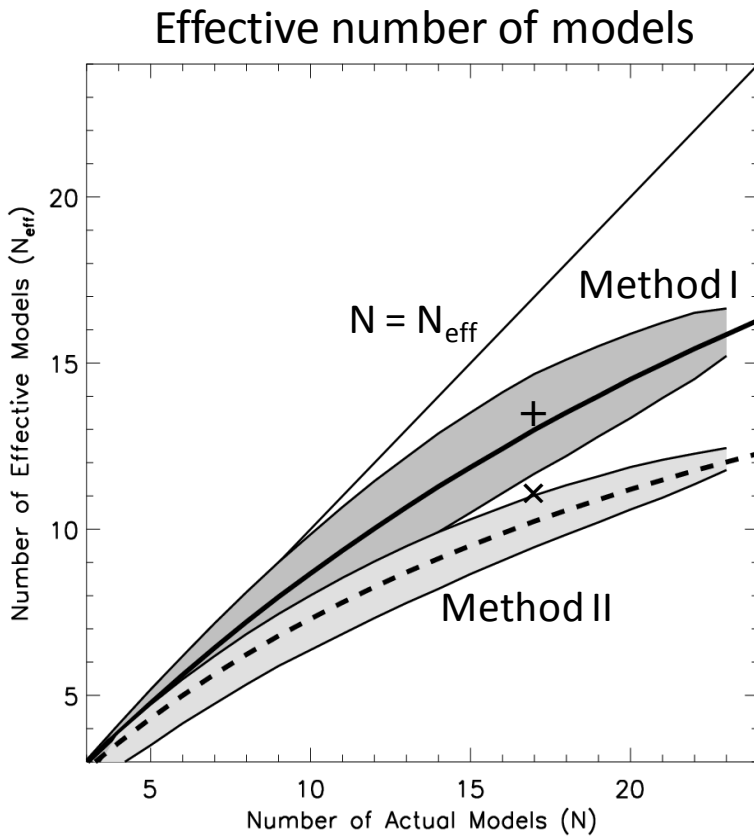
268

269 **Figure Captions**

270 **Figure 1.** Effective number of models ($N_{eff}$) and model similarities for the northern hemisphere

271 (30°-90°N). $N_{eff}$ as a function of $N$ for Method I (solid) and Method II (dashed) (see text). Grey

272 shading indicates 95% confidence intervals. + and x symbol denote $N_{eff}$ after excluding seven

273 specific companion models [GISS-ER, GISS-EH, UKMO-HadCM3, MIROC3.2(medres),

274 GFDL-CM2.0, CGCM3.1(T47), and CSIRO-Mk3.0].

275 **Figure 2.** Hierarchical clustering scheme based on model correlations. Similar (dissimilar)

276 models merge closer to the right (left).

277

**Table 1.** Climate quantities used in this study. Acronyms listed under 'Validating observations' (5[th] column) are commonly used in the literature to denote specific observational data-sets. The average is taken of validating observation sets where more than one is given for a particular quantity.

| | Quantity | Domain | Acronym | Units | Validating observations |
|---|---|---|---|---|---|
| **physics** | surface air temperature | global | TAS | K | CRU, ICOADS, NOAA |
| | surface skin temperature | land | TS | K | ISCCP |
| | zonal/meridional surface wind stress | ocean | TAUU, TAUV | $10^{-2}\,Nm^{-2}$ | GSSTF2, ICOADS |
| | sea level pressure | ocean | PSL | hPa | ERSLP, HADSLP, ICOADS |
| | surface sensible/latent heat fluxes | ocean | HFSS, HFLS | $Wm^{-2}$ | GSSTF2, HOAPS2, ICOADS, JOFURO, OAFLUX |
| | total cloudiness | global | CLT | % | CERES, ISCCP |
| | surface radiation (up/down, short-/longwave) | global | RSDS, RSUS, RLDS, RLUS | $Wm^{-2}$ | BSRN, CERES, GEBA, ISCCP |
| | TOA outgoing shortwave radiation | global | RSUT | $Wm^{-2}$ | CERES, ERBE, ISCCP |
| | TOA outgoing longwave radiation | global | RLUT | $Wm^{-2}$ | CERES, ERBE, ISCCP, NOAA |
| | TOA cloud radiative forcing | global | CFLT, CFST | $Wm^{-2}$ | CERES, ERBE, ISCCP |
| | precipitation | global | PR | mm/day | CMAP, GPCP |
| | precipitable water | global | PRW | mm | HOAPS2, NVAP |
| | snow coverage | global | SNW | % | NSIDC |
| | air temperature | zonal mean | TA | K | AIRS |
| **dynamics** | specific humidity | zonal mean | HUS | g/kg | ERA |
| | zonal/meridional wind 200 hPa | global | U200, V200 | m/s | ERA |
| | stream function 200 hPa | global | ψ200 | $10^6\,m^2s^{-1}$ | ERA |
| | velocity potential 200 hPa | global | χ200 | $10^6\,m^2s^{-1}$ | ERA |
| | temperature 200 hPa | global | T200 | K | ERA |
| | geopotential 500 hPa | global | Z500 | gpm | ERA |
| | stationary waves 500 hPa | global | SW500 | gpm | ERA |
| | zonal/meridional wind 850 hPa | global | U850, V850 | m/s | ERA |
| | zonal mean zonal/meridional wind | zonal mean | UA, VA | m/s | ERA |
| | mean meridional mass streamfunction | zonal mean | MMC | $10^9\,kg/s$ | ERA |
| **oceans** | sea surface height | ocean | ZOS | m | GRACE-DOT |
| | sea ice content | ocean | SIC | % | GICE |
| | sea surface salinity | ocean | SO | ‰ | NODC |
| | sea surface temperature | ocean | TOS | K | GISST |

## Effective number of models
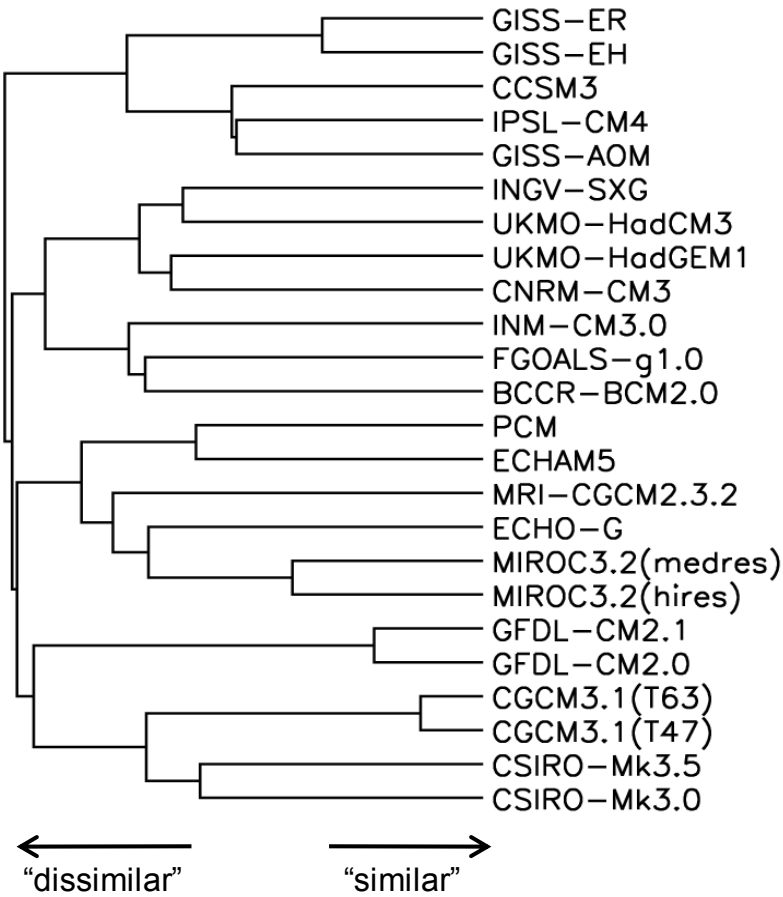
**Figure 1.** Effective number of models ($N_{eff}$) and model similarities for the northern hemisphere

(30°-90°N). $N_{eff}$ as a function of $N$ for Method I (solid) and Method II (dashed) (see text). Grey

shading indicates 95% confidence intervals. + and x symbol denote $N_{eff}$ after excluding seven

specific companion models [GISS-ER, GISS-EH, UKMO-HadCM3, MIROC3.2(medres),

GFDL-CM2.0, CGCM3.1(T47), and CSIRO-Mk3.0].

289

**Figure 2.** Hierarchical clustering scheme based on model correlations. Similar (dissimilar) models merge closer to the right (left).