1	Stochastic parameterization of convective area fractions
2	with a multicloud model inferred from observational data
3	Jesse Dorrestijn, * Daan T. Crommelin
	CWI, Amsterdam, The Netherlands
4	A. PIER SIEBESMA
	Royal Netherlands Meteorological Institute (KNMI), De Bilt, The Netherlands and
	Delft University of Technology, Delft, The Netherlands
5	HARMEN J.J. JONKER
	Delft University of Technology, Delft, The Netherlands
6	Christian Jakob

ARC Centre of Excellence for Climate System Science, Monash University, Melbourne, Australia

 $^{^{\}ast}Corresponding author address:$ Jesse Dorrestijn, CWI, P.O. Box 94079, 1090 GB, Amsterdam, The Netherlands

E-mail: J.Dorrestijn@cwi.nl

ABSTRACT

Observational data of rainfall from a rain radar in Darwin in Australia is combined with data 8 defining the large-scale dynamic and thermodynamic state of the atmosphere around Darwin 9 to develop a multicloud model based on a stochastic method using conditional Markov chains. 10 We assign the radar data to clear sky, moderate congestus, strong congestus, deep convective 11 or stratiform clouds and estimate transition probabilities used by Markov chains that switch 12 between the cloud types and yield cloud type area fractions. Cross-correlation analysis shows 13 that the mean vertical velocity is an important indicator of deep convection. We show that, 14 if conditioned on the mean vertical velocity, the Markov chains produce fractions comparable 15 to the observations. The stochastic nature of the approach turns out to be essential for the 16 correct production of area fractions. The stochastic multicloud model can easily be coupled 17 to existing moist convection parameterization schemes used in general circulation models. 18

¹⁹ 1. The cumulus parameterization problem

The representation of clouds and convection is of major importance for numerical weather 20 and climate prediction. Moist convection, also called cumulus convection, transports heat, 21 moisture and momentum vertically in the atmosphere, it influences dynamical, thermody-22 namical and radiative processes and it has an impact on the large-scale global circulation. 23 In general circulation models (GCMs), moist convection can not be explicitly resolved since 24 the scale of the involved processes is too small, therefore the sub-grid processes have to be 25 represented by parameterizations, which are formulations of the statistical effects of the un-26 resolved variables on the resolved variables. We refer to Arakawa (2004) for an overview of 27 the the cumulus parameterization problem. Formulating moist convection parameterizations 28 is a difficult problem: it introduces uncertainties in model predictions (e.g. Randall et al. 29 (2003)) and although models do agree that the cloud feedback is positive or neutral, they do 30 not agree on the strength of the cloud feedback, e.g. Flato et al. (2013). It has been shown 31 by Lin et al. (2006) that the intraseasonal variability of precipitation is generally too small 32 in models and that convectively coupled tropical waves are not well simulated. 33

An important issue considering cumulus parameterizations is that it is still not known 34 which large-scale resolved variables are most strongly related to moist convection, and on 35 which variables the closures of the parameterizations should be based. In general we have 36 the choice between dynamical (e.g. vertical velocity) or thermodynamical (e.g. the con-37 vective available potential energy (CAPE), relative humidity (RH)) variables, which have 38 been studied in a recent paper by Davies et al. (2013a). Another important issue is that 39 if parameterizations are chosen to be *deterministic* functions of the resolved variables, the 40 subgrid response of moist convection to large-scale variations can not cover the variety of 41 responses that is possible in reality, as deterministic parameterizations can only provide the 42 expected value of the response of moist convection in a grid box. In view that GCMs reso-43 lutions are getting finer and finer, this issue becomes more important, because with smaller 44 grid boxes the fluctuations around expected subgrid responses become larger. Palmer (2001) 45

⁴⁶ pointed out that neglecting subgrid variability can result in model errors and that this can ⁴⁷ be corrected by using *stochastic* parameterizations to represent subgrid processes. This has ⁴⁸ for example been shown by Buizza et al. (1999) who improved the skill of numerical weather ⁴⁹ prediction (NWP) with the European Centre for Medium-Range Weather Forecasts's sys-⁵⁰ tem by introducing stochastic elements in the physical parameterization tendency. Their ⁵¹ pioneering work gave impulse to develop more sophisticated stochastic schemes.

Instead of perturbing all subgrid processes at once, it is possible to improve GCMs by introducing stochastic elements only in the deep convection parameterization, e.g. Lin and Neelin (2000); Lin and Neelin (2003); Teixeira and Reynolds (2008); Plant and Craig (2008) and Bengtsson et al. (2013).

Rather than relying on physical intuition or deriving parameterizations from first prin-56 ciples, stochastic parameterizations can be inferred directly from data. Crommelin and 57 Vanden-Eijnden (2008) showed that Markov chains, with only a few states, for which the 58 transition probabilities had been estimated from data, could represent the subgrid terms in 59 the Lorenz '96 (Lorenz (1996)) model quite well, better than the deterministic parameter-60 izations and the stochastic parameterizations, based on autoregressive processes, of Wilks 61 (2005). The data-driven Markov chain model inspired Kwasniok (2012) to develop a similar 62 model based on cluster-weighted Markov chains. In Dorrestijn et al. (2013b) the Markov 63 chain model of Crommelin and Vanden-Eijnden (2008) was used to study stochastic param-64 eterization of shallow convection and in Dorrestijn et al. (2013a) for deep convection. 65

⁶⁶ A promising class of moist convection parameterizations based on the idea of evolving an ⁶⁷ ensemble of several (convective) cloud types, inspired by Mapes et al. (2006) and Johnson ⁶⁸ et al. (1999), is formed by *multicloud models*, e.g. Khouider and Majda (2006); Khouider ⁶⁹ et al. (2010); Majda et al. (2007); Frenkel et al. (2013); and Peters et al. (2013). The ⁷⁰ clouds follow a life cycle starting from clear sky to *congestus clouds*, to *deep cumulus* towers ⁷¹ with *stratiform* anvil clouds as a remnant of the towers spreading over large areas, finally ⁷² dissolving and come full circle at clear sky. In the multicloud model of Dorrestijn et al. ⁷³ (2013a) also *shallow cumulus* clouds are included.

In the present paper we use high-resolution ($\sim 2.5 \times 2.5 \text{ km}^2$) observational data of 74 rainfall in combination with data defining the large-scale ($\sim 150 \times 150 \text{ km}^2$) dynamical 75 and thermodynamical state of the atmosphere to infer such a stochastic multicloud model. 76 The large-scale data are NWP analysis variable estimates improved with observations. The 77 model is similar to the multicloud model of Dorrestijn et al. (2013a) in which Large-Eddy 78 Simulation data was used to infer the model, as opposed to the observational data of this 79 study. The multicloud model produces area fractions for several cloud types which can be 80 used as stochastic parameterizations in the deep convection and cloud schemes of GCMs. 81 We also determine which large-scale variables are strongly related to deep convection. 82

Our paper is organized as follows. In Section 2 we explain how we use Markov chains 83 as a foundation for our multicloud model. Then, in Section 3 we give a description of the 84 observational data, explain how we classified the data into cloud categories and how we dealt 85 with advection while estimating transition probabilities between cloud states. In Section 4 86 we assess the skill of large-scale variables as indicators for deep convection. In Section 5 we 87 construct our model, give expected area fractions and standard deviations and we discuss 88 scale adaptivity, i.e. the ability to adapt the the size of a GCM grid box. We give results in 89 Section 6 by comparing area fractions from the model with the observations and looking at 90 their autocorrelation functions. In Section 7 we discuss the possibilities of implementation of 91 the stochastic model in a convection parameterization of a GCM and make some concluding 92 remarks. 93

⁹⁴ 2. Markov chains

The multicloud model we use in this study consists of Markov chains positioned on the nodes of a 2-dimensional micro-grid. This model set-up has been used before in Khouider et al. (2010); Dorrestijn et al. (2013a); Peters et al. (2013). The state of each Markov chain at time t is denoted $Y_n(t)$, where n is the micro-grid index. Each Y_n can take on 5 different values, corresponding to the following categories: clear sky, moderate congestus, strong congestus, deep convective and stratiform. The choice of these specific categories will be discussed in Section 3. We will refer to these categories as *cloud types*. As time evolves, the Markov chains can switch, or "make a transition", between states every $\Delta t = 10$ minutes. All the Markov chains on the micro-grid together determine the area fractions σ_m for the various cloud types:

$$\sigma_m(t) = \frac{1}{N} \sum_{n=1}^N \mathbf{1}[Y_n(t) = m],$$
(1)

in which **1** is the indicator function $(\mathbf{1}[A] = 1$ if A is true, 0 otherwise), N is the number of micro-grid nodes, and $m \in \{1, ..., 5\}$ the cloud type. We use radar data to estimate the transition probabilities, needed in the Markov chain model.

¹⁰⁸ When used in a GCM, each GCM column contains N Markov chains that can switch ¹⁰⁹ to a different state every 10 minutes, resulting in time-evolving area fractions σ_m for each ¹¹⁰ cloud type and for each GCM column. These area fractions can be used in the convection ¹¹¹ and cloud schemes of a GCM. For example, the deep convective area fractions, σ_4 , can serve ¹¹² as a mass flux closure at cloud base for a deep convection parameterization scheme:

$$M_b = \rho \,\sigma_4 \, w_{cb},\tag{2}$$

¹¹³ in which ρ is the density and w_{cb} is the vertical velocity in a deep convective updraft at cloud ¹¹⁴ base (e.g. Arakawa et al. (2011); Möbis and Stevens (2012)). More examples of possible ¹¹⁵ applications in GCMs are given in Section 7.

As mentioned before, we use Markov chains with 5 possible states, so that the transition probabilities form a 5×5 transition matrix. Since these transition probabilities depend strongly on the large-scale state of the atmosphere, we make these probabilities conditional on functions of large-scale variables (i.e., the variables that are normally resolved by GCMs). These functions are called *indicators* of deep convection. In Section 4 we discuss appropriate indicators. The framework of conditional Markov chains (CMCs) for parameterization was introduced by Crommelin and Vanden-Eijnden (2008). For now, we consider a discretized indicator X, such that the possible states of X correspond to a finite number Γ of large-scale states. So, for each $\gamma \in \{1, \ldots, \Gamma\}$ we estimate a 5 × 5 transition probability matrix. The probability of CMCs switching from state α to state β given the large-scale state γ can be estimated as follows (see also Crommelin and Vanden-Eijnden (2008)):

128

$$\frac{T_{\gamma}(\alpha,\beta)}{\sum_{\beta}T_{\gamma}(\alpha,\beta)}$$

 $\operatorname{Prob}(Y_n(t + \Delta t) = \beta | Y_n(t) = \alpha, X(t) = \gamma) =$

(3)

129 where

$$T_{\gamma}(\alpha,\beta) = \sum_{t,n} \mathbf{1}[Y_n(t+\Delta t) = \beta]\mathbf{1}[Y_n(t) = \alpha]\mathbf{1}[X_n(t) = \gamma]$$

counts the number of transitions observed in the data from cloud type α to β given that the 130 large-scale state is γ . The indices n and t run over space and time covered in the training 131 *data set* which is used to estimate the transition probabilities. We remark that we do not 132 condition the Markov chains on $X(t+\Delta t)$, which reduces the number of matrices to estimate 133 significantly. For the estimation of the transition matrices we use data sets corresponding to 134 two different scales: data sets that are formed by high-resolution observations of rainfall at 135 a scale that is equal to or smaller than the micro-grid scale of the CMCs and data sets that 136 represent the large-scale atmospheric state at the grid scale of a GCM. In the next section 137 we introduce the high-resolution observation data sets. 138

¹³⁹ 3. The radar data

The microscale data consists of observational data of precipitation obtained from the Darwin C-Band Polarimetric (CPOL) Radar in Darwin, North-Australia. This data is described in detail in Kumar et al. (2013). In the same article it is explained how the radar data can be used to calculate cloud top height (CTH) and rain rates. For two time periods, 10 November 2005-15 April 2006 and 20 January 2007-18 April 2007, we have integer valued

CTH and rain rate observations at 10-minute timesteps, for a circular area with radius 150 145 km and resolution of 2.5×2.5 km². In Fig. 1 we show a snapshot of the CTH and the rain 146 rates at one time instance. The fields are rather noisy at the outer ring of the radar domain 147 and the radar does not give observations in the center of the radar domain, which is known 148 as the "cone of silence" and is due to the 42° maximum elevation angle (May and Ballinger 149 (2007)). Therefore, we only use pixels in between 25 km and 97.5 km from the center of the 150 domain. This forms an annular shaped sub-domain consisting of 4720 pixels of $2.5 \times 2.5 \; \rm km^2$ 151 corresponding to an area size of approximately 172×172 km². Fig. 2 contains histograms of 152 the CTH and the rain rates, showing the distribution of these quantities. We consider CTH 153 below 1.5 km as clear sky to avoid the influence of radar ground clutter. There is a bi-modal 154 distribution of CTH, with a minimum at around 4 km, which is close to the freezing level 155 at 5 km. To classify our cloud types, we use thresholds for CTH to distinguish high clouds, 156 low clouds and clear sky. The bi-modal distribution in the cloud top histogram suggests 157 a CTH threshold to distinguish low and high clouds (e.g. congestus and deep convective 158 clouds) of around 4 or 5 km. Congestus clouds have been observed up to 9.5 km in the 159 atmosphere (Johnson et al. (1999)). We adopt the approach of Kumar et al. (2013), who 160 developed a more objective identification of congestus and deep convective clouds, taking the 161 value 6.5 km as a threshold. Further, we employ a rain rate threshold to make a distinction 162 between clouds with intense precipitation and those with little or no precipitation. This 163 enables us to make a distinction between deep convective clouds and stratiform clouds as 164 well as a distinction between strong and moderate congestus. The rain rate histogram in 165 Fig. 2b, shows an approximately exponential distribution, so it is impossible to argue for 166 an obvious rain rate threshold. In the literature thresholds for partitioning convective and 167 stratiform precipitation vary between 10 and 25 mm h^{-1} , and there are several methods for 168 partitioning which are described in Lang et al. (2003). We choose a threshold of 12 mm h^{-1} 169 to distinguish between deep convective and stratiform clouds and a threshold of 3 mm h^{-1} 170 to distinguish between moderate and strong congestus. Combining these thresholds results 171

in the following five cloud types: (1) clear sky, (2) moderate congestus, (3) strong congestus,
(4) deep convective and (5) stratiform. In Table 1 we summarize the classification into cloud
types. Note that, although desired, shallow cumulus clouds are not included in the model,
for the obvious reason that the rain radar does not observe non-precipitating clouds.

After classification we have 2-dimensional fields with discrete values (integers from 1 to 5). 176 In Fig. 3 we give an example of a classified field, which is the classified field corresponding 177 to the CTH and rain rate fields shown in Fig. 1. After the classification the observed 178 area fractions, σ_m , can be calculated according to (1), with Y_n the observed cloud type and 179 N = 4720 the number of radar pixels in the annular domain. The observed area fractions are 180 strongly time-dependent, with σ_1 (clear sky) varying between 0% and 100%, σ_2 (moderate 181 congestus) between 0% and 55%, σ_3 (strong congestus) between 0% and 2.5%, σ_4 (deep 182 convective) ranging from 0 to about 10% and σ_5 (stratiform) ranging from 0 to about 99%. 183 The observed fractions are depicted in Fig. 9 (discussed in Section 6) for a time period of 5 184 days for all cloud types, and the deep convective area fraction also in Fig. 7a (discussed in 185 Section 6) for a longer period of 3 months. 186

Besides calculating observed area fractions for the different cloud types, the classified data are used to estimate transition probabilities between the cloud types for the CMCs, using (3). This is a key step in creating the multicloud model. To give an idea of the observed transition probabilities, not yet conditioned on the large-scale variables, we give the estimated transition matrix:

192

¹⁹³ The probability of a transition from cloud type m to cloud type n can be found in the nth

column of row m. For example, the probability that a deep convective pixel will be assigned 194 to stratiform 10 minutes later, is **0.3540**. The probability that a deep site is again a deep 195 site 10 minutes later, is 0.4295, much larger than the expected deep convective area fraction 196 (at most 0.03 as can be seen Fig. 6, discussed later in this paper). Some evidence for the 197 life cycle can be seen in this transition matrix, a deep convective cloud likely turns into 198 stratiform, which turns into clear sky. Some entries are artefacts of the estimation method, 199 for example the probability of clear sky turning into stratiform is 0.0329, but in reality the 200 stratiform cloud spreads out from the top of a deep cumulus cloud. 201

For correct estimation of cloud type transition probabilities, we have to take into account 202 that clouds are advecting horizontally through the domain. To do this, we translate the 203 advected clouds in a radar image back to their position in the previous image. In this way, 204 we minimize transitions that are only a result of advection. The advection, with zonal wind u205 and the meridional wind v, is assumed to be a function of height and time only. We calculate 206 this translation separately for every cloud type (as they are located at different heights in 207 the atmosphere). Let $Z_m(x_i, y_j, t) = \mathbf{1}[Y(x_i, y_j, t) = m]$, with $Y(x_i, y_j, t)$ the discretized 208 radar pixel at location (x_i, y_j) at time t and (x_i, y_j) running over all $N_{ij} = 4720$ pixels in 209 the annular shaped sub-domain. We calculate for every cloud type m and for every time 210 interval $[t, t + \Delta t]$ the optimal horizontal displacements $u_m \Delta t$ and $v_m \Delta t$ which minimize the 211 correlation 212

$$\frac{1}{N_{ij}}\sum_{ij}Z_m(x_i+u_m\Delta t, y_j+v_m\Delta t, t)Z_m(x_i, y_i, t+\Delta t).$$

By applying the Correlation Theorem (e.g. Press et al. (1992)), fast Fourier transforms can be used to reduce the calculation time for finding the displacements. At the boundaries at the outer edge and in the center of the radar domain, clouds flow into and out of the domain. We also have to account for this during the estimation of cloud type transition probabilities. More specifically, we do not count transitions of "clouds" (including clear sky) that are inside the radar domain at time t, but which are outside the domain at the previous time step $t - \Delta t$ or at the next time step $t + \Delta t$, due to advection. Without corrections, the estimated probability transition matrix is significantly different: for example the probability that a pixel assigned to the deep convective cloud type is deep convective 10 minutes later would be estimated at 0.29 instead of 0.43.

The focus in this paper will primarily be on the deep convective area fractions, when we determine the large-scale variable on which to condition the CMC (Section 4) and when we test the CMC (Section 7). Although the other fractions can have applications in GCMs, the deep convective area fractions are the most important. Describing the convective transport by deep convection accurately is crucial for a GCM to work properly. Conditioning each individual cloud type on different large-scale variables could improve the model, in particular for the strong congestus clouds, that precede deep convection.

²³⁰ 4. The large-scale data

We have data available that defines the large-scale dynamic and thermodynamic state of 231 the atmosphere around Darwin for the time periods November 2005-April 2006 and January 232 2007-April 2007 for which we also have the radar data. The large-scale fields are averages over 233 6 hour intervals and have a vertical resolution of 40 pressure levels, from ground level to about 234 20 km altitude. The data has been prepared by Davies et al. (2013a) who used a variational 235 analysis method to improve NWP analysis large-scale variable estimates by constraining 236 the moisture budgets with observational rain data from the CPOL radar. The large-scale 237 data is also used in Davies et al. (2013b); Peters et al. (2013) and Gottwald et al. (2014). 238 Here, we use the data to investigate which large-scale variables are suitable indicators for the 239 convective state of the atmosphere and compare our findings with the results of Davies et al. 240 (2013a). Then, we will use the large-scale data accordingly for conditioning the multicloud 241 CMC model. As in Davies et al. (2013a), we consider thermodynamical and dynamical 242 variables. In particular, we will consider the following well-known indicators: CAPE, the 243 mean vertical velocity $\langle \omega \rangle$, and RH. CAPE is a measure for the stability of the atmosphere 244

²⁴⁵ and is formally defined as follows:

CAPE :=
$$R_d \int_{p_{\rm NB}}^{p_{\rm LFC}} (T_{v,p} - \overline{T}_v) d\ln p$$
,

in which $T_{v,p}$ is the virtual temperature of an undiluted parcel, \overline{T}_v is the virtual temperature of the environment, R_d is the gas constant of dry air, p_{NB} the level of neutral buoyancy and p_{LFC} the level of free convection (e.g. Siebesma (1998)). The mean vertical velocity we define as

$$\langle \omega \rangle := \frac{1}{p_0 - p^*} \int_{p^*}^{p_0} \overline{\omega}(p) dp,$$

in which $\overline{\omega}$ is the large-scale vertical velocity in hPa h⁻¹, p_0 the pressure at the surface, and p^* 250 is pressure level 340 hPa, chosen because the resulting $\langle \omega \rangle$ gives the highest correlation with 251 deep convective area fractions (as calculated with (4) that is given below). We find that the 252 vertical integral over ω gives higher correlations than ω at a single pressure level. Further, 253 the relative humidity is chosen at pressure level 640 hPa, also because it gives the highest 254 correlation with deep convective area fractions. To assess how well an indicator correlates 255 with deep convection, we calculate the time-lagged cross-correlation function (CCF) of the 256 indicator and the deep convective area fraction. 257

Given the timeseries of the deep convective area fraction $\sigma_4(t)$ and the timeseries of the indicator X(t), the normalized CCF of X(t) and $\sigma_4(t)$ is:

$$CCF(\tau) = \int_{-\infty}^{\infty} \tilde{X}(t+\tau)\tilde{\sigma}_4(t)dt$$
(4)

with $\tilde{X}(t) = \frac{X(t) - \mu_X}{\sigma_X}$ (i.e. the indicator normalized by subtracting its mean μ_X and dividing by its standard deviation σ_X), $\tilde{\sigma}_4$ defined analogously, and τ the time lag of X w.r.t. σ_4 . As such, the CCF lies in between -1 and 1. If the maximum value of the CCF is attained at positive time lag τ , the indicator X(t) tends to *follow* rather than *precede* deep convection. In Fig. 4 we plot the CCFs of the indicators $-\langle\omega\rangle$, CAPE and RH with the observed deep convective area fraction for the 2005/2006 period. The figure for the 2007 period is similar (not included). Before calculating the CCF, we linearly interpolate X to get its

values every 10 minutes instead of every 6 hours, because the sequences X and $\tilde{\sigma}_4$ must have 267 the same length. We see that $\langle \omega \rangle$ has a larger correlation at zero time lag than CAPE and 268 RH. Moreover, also for negative time lags of a few hours this correlation is higher. In this 269 respect $\langle \omega \rangle$ is the best indicator of deep convection. We note that the maximum correlation 270 of $\langle \omega \rangle$ with σ_4 is attained at a positive time lag. This may seem to indicate that $\langle \omega \rangle$ is an 271 effect rather than a cause of deep convection. However, this is a subtle issue, as $\langle \omega \rangle$ may also 272 both be a trigger (i.e., cause) of deep convection and be reinforced by it, so that separating 273 cause and effect becomes difficult. In Peters et al. (2013) a related discussion can be found. 274 In order to use an indicator for constructing the CMC according to (3), it must be 275 discretized into a finite number of states. If only one indicator is used, which is the case is this 276 paper, a finite number (Γ) of intervals can be chosen, defined by thresholds. If a combination 277 of several indicators is used, one can choose thresholds for each indicator separately, or use 278 a clustering method as in Dorrestijn et al. (2013b,a) and Kwasniok (2012). To give an 279 example, in Fig. 5 we show a histogram of $\langle \omega \rangle$ discretized using 25 intervals. These intervals 280 have been found by using a cluster method, k-means, which minimizes the distance between 281 the $\langle \omega \rangle$ -values and the centers of the intervals. Using equidistant intervals is also an option, 282 however, since the $\langle \omega \rangle$ -values are not distributed uniformly, we prefer the non-equidistant 283 intervals found by k-means. Interval number 25, corresponds to negative $\langle \omega \rangle$ or strongly 284 positive large-scale vertical velocity (illustrated by the arrow), which is favourable for deep 285 convection, and we will later see in Fig. 6 that the averaged observed deep convective and 286 stratiform area fractions are large (around 3% and 90%, respectively) for interval number 287 25.288

²⁸⁹ 5. A description of the multicloud model

Having classified the radar data into cloud types, and having identified (and discretized) a suitable large-scale indicator, $\langle \omega \rangle$, we estimate the transition probability matrices of the ²⁹² CMC using (3). We take the period from 10 November 2005 until 15 April 2006 as the ²⁹³ training data set, and we set $\Gamma = 25$. So, we have to estimate 25 matrices each of size 5×5 , ²⁹⁴ giving 625 parameters in total. This may seem a large number, however the training data ²⁹⁵ set is very large, containing $O(10^8)$ observations of transitions (radar images at 10-minute ²⁹⁶ intervals during 157 days, with 4720 pixels in each image).

In Section 6 we will validate the CMCs with the *test data set*, but since we have estimated transition matrices, we can already get some insight into the statistical properties of the cloud type area fractions generated by the CMC as compared to the observed area fractions in the *training data set*.

In Fig. 6, we plot the expected fractions and the standard deviation for both the observations and the CMC as a function of the $\langle \omega \rangle$ -intervals seen before in Fig. 5. The expected values of the CMC correspond to the invariant distribution of the transition matrix for each $\langle \omega \rangle$ -interval. The CMC expected values are almost equal to the observational expectations for all cloud types, the small differences can be ascribed to the way we corrected for horizontal advection (as described before in Section 3).

We see in Fig. 6a that the expected deep convective area fractions increase with increas-307 ing $\langle \omega \rangle$ -interval (corresponding to increasing upward mean vertical velocities) and has its 308 maximum of around 0.03 for interval number 24. Further, the strong congestus fractions 309 in Fig. 6b, increase with increasing $\langle \omega \rangle$ -interval, however, for interval number larger than 310 22, the fraction decreases rapidly, while expected deep and stratiform cloud fractions keep 311 increasing. The expected stratiform fractions increase with increasing $\langle \omega \rangle$ -interval up to 312 very high expected values of 90%. The expected value of moderate congestus is around 313 15% for downward mean motion, increases slightly with increasing $\langle \omega \rangle$ -interval number. For 314 $\langle \omega \rangle$ -interval numbers above 22, the expected value of moderate congestus decreases which is 315 caused by the stratiform decks that are dominating the radar domain (for this $\langle \omega \rangle$ -interval 316 numbers). Expected clear sky fractions decrease rapidly as a function of the $\langle \omega \rangle$ -interval. 317 The standard deviation of the observational deep convective arac fractions tends to in-318

crease with increasing $\langle \omega \rangle$ -interval number, so it tends to increase if the expected value increases and for high values of the $\langle \omega \rangle$ -interval number the standard deviation is almost equal to the expected value. The normalized observed standard deviation, the standard deviation divided by the mean, is decreasing with increasing mean, with values decreasing from 5 down to about 1. So, we agree with the conclusion of Davies et al. (2013a) that noise (or stochastic behaviour) decreases as a function of increasing forcing.

The standard deviation of the observational strong congestus area fractions depends 325 on the expected values as well, with a normalized standard deviation ranging from 1 (for 326 relatively high fractions) up to 3 (for relatively low fractions). The standard deviation of the 327 stratiform area fractions tends to increase as a function of the $\langle \omega \rangle$ -interval, but decreases if 328 the expected values become very large because of the upper bound of 100%. For moderate 329 congestus, the normalized standard deviation ranges between 0.5 and 1. The standard 330 deviation of the clear sky area fraction is around 10-20%, independently of the $\langle \omega \rangle$ -interval 331 number, with an exception of interval number 25 for which the standard deviation is only 332 2.4%. 333

The theoretical standard deviation of the CMC can be calculated explicitly for each $\langle \omega \rangle$ -334 interval and is equal to $\sqrt{N^{-1}p(1-p)}$, in which p is the expected value of the fraction. So, 335 the theoretical standard deviation depends only on the expected value of the fraction and 336 the number of CMCs used to calculate the cloud type area fractions. We choose a value of 337 N = 100 such that the standard deviation of the deep convective area fractions is comparable 338 to the standard deviation of the observed deep convective area fractions in the training data 339 set. This implies that the standard deviation of the CMC is too small for cloud types with 340 larger standard deviations (clear sky, moderate congestus and stratiform) and too large for 341 the strong congestus cloud type (which has a small standard deviation). 342

343 Scale adaptivity

Ideally a parameterization of deep convection should be adaptive to the size of the GCM 344 grid box, see Arakawa et al. (2011). By construction of the multicloud model, our parame-345 terization of deep convection is indeed scale adaptive. The value N of the number of CMCs 346 can be adapted to the horizontal grid spacing of the GCM. For a large size of the GCM grid 347 box, a large number of clouds fit into the model column and therefore a large number of 348 CMCs should be taken to calculate the cloud type area fractions. For very large GCM grids, 349 the number of CMCs becomes very large and hence the σ_m tend to a deterministic limit 350 (equal to the expected values associated with the large-scale interval number). For smaller 351 grid box sizes, the number of CMCs is smaller and as a result, the area fractions generated 352 by the multicloud model will be "more stochastic", fluctuating significantly around their 353 expected values. It is difficult to say to which horizontal size a CMC corresponds exactly. 354 The size corresponding to a CMC is equal to the typical horizontal size of the cloud type 355 under consideration. Therefore, the horizontal size is larger than the area of a radar data 356 pixel $(2.5 \times 2.5 \text{ km}^2)$, which explains that producing area fractions with CMCs while using 357 a number smaller than the number of radar pixels in the radar domain gives better results 358 in Section 6, N = 100 versus N = 4720. We emphasize that the value of N = 100 is found 359 during the training phase and not during the the testing phase of the model. For N = 100360 the horizontal area size corresponding to a CMC is approximately 17×17 km², which is the 361 area of the 4720 radar pixels divided by 100. 362

To summarize the different length scales that are used in this paper: radar pixels of 2.5 km, clouds of length scale ~ 17 km which is also the length scale corresponding to 1 CMC, the length corresponding to the large-scale variables, the radar domain and a GCM grid box ~ 150 km. Finally, the length scale corresponding to the CMC fractions: $17\sqrt{N}$.

³⁶⁷ 6. Results

To assess how well the multicloud model reproduces the convective behaviour observed in the radar data set, we first consider the cloud type area fractions. Then, we will look at autocorrelation functions (ACFs) of the fractions and $\langle \omega \rangle$.

371 a. Area fractions

As mentioned, the radar data can be used to calculate observed area fractions of each cloud type. We use $\langle \omega \rangle$ as indicator and take N = 100 CMCs. Then, we train the CMCs as explained in Section 5 using the training data set 2005/2006. We assess the model by driving the CMCs with $\langle \omega \rangle$ as observed in the other data set (from 2007). Thus, different data sets are used for training and evaluation.

In Fig. 7a we show the deep convective area fractions as observed in the Darwin radar 377 test data set (2007). It can be seen that the deep convective events are very intermittent 378 in the radar data, with periods of enhanced deep convection and periods with less wide-379 spread convective events. In Fig. 7b and 7c we give two realizations of the deep convective 380 area fractions as reproduced by the CMCs. The CMC fractions display similar intermittent 381 behaviour, with maximum values that are slightly too high compared to the observations. 382 The CMC fractions have discrete values, namely $\sigma_4 \in \{0.01, 0.02, 0.03, \ldots\}$, because N = 100383 CMCs are used. To further assess the quality of the deep convective fractions, we calculate 384 histograms of the deep convective area fractions (Fig. 7d). Since the CMC fractions are 385 integer multiples of 0.01, we bin the Darwin observed fractions into intervals of length 0.01, 386 apart from the first interval which is [0, 0.005). Because high values of the deep convective 387 fractions are rare, we plot the histograms on a logarithmic y-axis. We observe that the 388 observational fractions decrease exponentially, as is expected since rain rates tend to decrease 389 exponentially (see Fig. 2). The CMC fractions follow the exponential decrease well and the 390 values are only slightly off. 391

We repeat the computations with CAPE as indicator instead of $\langle \omega \rangle$. In Fig. 8a we 392 show the resulting CMC deep convective area fractions (compare to Fig. 7a). We observe 393 that the fractions are also intermittent, but high fraction values are too rare. Further, 394 although periods of enhanced convection and of less convective events are visible, they are 395 not comparable with the observations. In the histograms with a logarithmic y-axis (Fig. 8b) 396 it is indeed visible that fractions larger than 0.04 are too rare, although a fraction of 8% is 397 reached in one of the 100 realizations. We conclude that in the present setting CAPE is less 398 suitable as indicator for deep convection than $\langle \omega \rangle$. 399

As our third experiment, we use $\langle \omega \rangle$ again as indicator and keep everything as in the 400 first experiment except for taking $N = 69^2 = 4761$ which is (close to) the number of radar 401 pixels used to train the CMCs. We observe (Fig. 9) that high values of the deep convective 402 area fractions are not reached anymore, values are not higher than 0.04. Because N is much 403 larger than before, the fractions are rather close to the (deterministic) expectation values. 404 This means that, although the number of CMCs is equal to the number of radar lattice sites, 405 the CMC fractions show lower maxima. We note that in our current set-up the CMCs on the 406 2D micro lattice sites are independent of their lattice neighbors, which is not the case for the 407 sites in the radar data. This is the underlying cause of the lower CMC maxima. Introducing 408 local interactions between neighboring CMCs can improve this, but it makes the estimation 409 of the CMCs much more complicated, see Dorrestijn et al. (2013a). 410

As a final experiment we take again N = 100 CMCs and $\langle \omega \rangle$ as indicator, but we interchange the roles of training data set and test data set. Thus, we train the CMCs with the 2007 data set and validate using fractions for the 2005/2006 period. The deep convective area fractions in the 2005/2006 radar data reach higher maxima than in the 2007 data set, with an overall maximum of about 10 percent (not shown). The fractions of the CMCs are less likely to attain these highest peak values. Notwithstanding this issue, the distribution of the CMC fractions is still comparable to that of the observed fractions.

For a more detailed look at the fractions, in Fig. 10 we show the area fractions of all

5 cloud types corresponding to the first experiment (with N = 100 and $\langle \omega \rangle$ as indicator) 419 for a much shorter period of 5 days. The timing of the deep convective events produced by 420 the CMCs is almost correct, there is a small time lag visible in Fig. 10a. Furthermore, it 421 is clear that the deep convective fractions of the CMC show maximum values of the peaks 422 in agreement with the observations, which is not the case for the expected values of the 423 CMC. The conclusion is that the stochastic fluctuations of the multicloud model fractions 424 are needed in order to produce the correct maximum values of the deep convection area 425 fraction peaks. The stochastic nature of the approach is essential for production of the 426 correct area fractions. A day-night cycle can be seen in the deep convective fractions, owing 427 to the presence of land in the radar domain. This cycle is also present in the CMC fractions. 428 The strong congestus fractions in Fig. 10b are small, so the CMC fractions, being integer 429 multiples of 0.01, have difficulties attaining the observational fractions. So, N = 100 seems 430 to be too small for the strong congestus area fractions. In Fig. 10c, we see stratiform area 431 fractions. The CMC fractions follow the observations correctly (in a time sense), but the 432 local maxima tend to be too low. The stochastic part of the fractions is not as prominent 433 as for the deep convective area fractions. The observational moderate congestus fractions 434 in Fig. 10d are difficult to follow for the CMCs: the value zero is never attained for the 435 CMC fractions. A conclusion is that $\langle \omega \rangle$ is not such a good indicator of moderate congestus 436 clouds. These depend probably more on boundary layer processes. The clear sky fractions 437 (Fig. 10e) of the CMC follow the observations quite well, but the minimum values of are not 438 small enough. The clear sky fractions are important, as $1 - \sigma_1$ is the cloud cover observed 439 by the radar, which is a usable quantity in GCMs, however, keep in mind that the radar is 440 not able to detect all clouds. 441

442 b. Autocorrelation functions

As a final assessment in this paper, we inspect ACFs of the cloud type area fractions and $\langle \omega \rangle$. The ACF of the cloud type area fraction σ_m is

$$ACF(\tau) = \int_{-\infty}^{\infty} \tilde{\sigma}_m(t+\tau)\tilde{\sigma}_m(t)dt,$$
(5)

which is the CCF of $\tilde{\sigma}_m$ with itself, cf. (4). Recall that $\tilde{\sigma}_m$ is the normalized σ_m . The ACF 445 of $\langle \omega \rangle$ is defined analogously. A main advantage of using Markov chains instead of drawing 446 samples that are uncorrelated in time from the observed distribution of cloud types is that 447 a Markov process should be better capable of capturing the observed ACF. In Fig. 11 we 448 show normalized ACFs of the observed area fractions (solid line with stars), the CMC area 449 fractions with N = 100 conditioned on $\langle \omega \rangle$ (solid line) and on CAPE (dashed line) and the 450 ACF corresponding to 69² CMCs conditioned on $\langle \omega \rangle$ (dotted line), for (a) deep convective 451 (b) strong congestus (c) stratiform (d) moderate congestus and (e) clear sky. Also the 452 ACF of $\langle \omega \rangle$ is shown (dash-dotted line). In (a) we see that appenently, the ACF of the 453 deep convective area fractions produced by N = 100 CMCs decreases too rapidly initially. 454 Without the correction for advection as explained in Section 3 the ACF decreases even more 455 rapidly (not shown). The rapid initial decrease indicates that the probability of a transition 456 from deep to deep is estimated too low. We see that the daily cycle is well captured in the 457 case that we conditioned on $\langle \omega \rangle$. When CAPE is used as indicator the ACF decreases more 458 rapidly than when conditioned on $\langle \omega \rangle$ and it can be seen that the daily cycle is not captured. 459 The ACF for the observational data set of 2005/2006 is similar to the ACF for the 2007 data 460 set (not shown). We note that for a large number of CMCs, close to the deterministic limit, 461 the ACF follows the ACF of $\langle \omega \rangle$ almost perfectly. In (b), we see that in order for the CMCs 462 to follow the observational strong congestus ACFs, the $N = 69 \times 69$ performs better than the 463 $N = 10^2$. In (c) and (e) we see ACFs of the CMC, that are comparable to the observational 464 ACF, only if conditioned on $\langle \omega \rangle$, not if conditioned on CAPE. The presence of a daily cycle 465 in the fractions is clearly visible if conditioned on $\langle \omega \rangle$ except for strong congestus fractions 466

⁴⁶⁷ produced with N = 100 CMCs. Considering all ACFs, we conclude that the ACFs for CMCs ⁴⁶⁸ conditioned on $\langle \omega \rangle$ are better than if conditioned on CAPE (except for moderate congestus). ⁴⁶⁹ For N = 100, the ACF of deep convection is better than for $N = 69^2$, while this is not the ⁴⁷⁰ case for strong congestus and moderate congestus. For stratiform and clear sky, the number ⁴⁷¹ of CMCs does not strongly influence the ACFs.

472 7. Discussion and conclusion

In this study we constructed a multicloud model from observational radar data in Dar-473 win, Australia, combined with large-scale data representing the atmosphere around Darwin. 474 The multicloud model consists of CMCs switching between different cloud types (moderate 475 congestus, strong congestus, deep convective and stratiform clouds and clear sky), a model 476 set-up similar to Khouider et al. (2010) and Dorrestijn et al. (2013a). The model is able to 477 reproduce cloud type area fractions comparable to the observational fractions (especially for 478 the deep convective area fractions, on which we focussed primary). The vertically averaged 479 large-scale vertical velocity $\langle \omega \rangle$ was found to be a good indicator, whereas CAPE or RH 480 were found to be less suitable indicators. This is in agreement with the findings of Davies 481 et al. (2013a). 482

The number N of CMCs used to form cloud type area fractions was shown to be an 483 important parameter of the model: for moderate values of N the model shows significant 484 stochastic fluctuations and the model is able to produce area fractions comparable with the 485 observational fractions. For large values of N the model is more deterministic and unable 486 to reproduce fractions well. The stochastic nature of the model is essential for making the 487 fractions comparable to the observations. Further, by changing N the multicloud model 488 can be adapted to the horizontal scale if implemented in a GCM, providing a way to make 489 the parameterization scale-adaptive. This makes the model suitable for GCMs using non-490 uniform grids. Further, the model can be used as a start for GCMs reaching grid sizes that 491

fall in the grey zone, i.e. for grid sizes so small that subgrid convective flux terms are of the 492 same order as the resolved flux terms. For a discussion of the grey zone we refer to Yu and 493 Lee (2010) and Dorrestijn et al. (2013b). The horizontal size to which a CMC corresponds 494 is not clearly determined. In principle it corresponds to the horizontal size of the cloud 495 type under consideration, which is different for all cloud types. Using a different number of 496 CMCs for each cloud type is an option, but it is complicated and lies out of the scope of 497 this research. During the training process, we found a value of N = 100, only taking deep 498 convective area fractions into account, which corresponds to an area size of 17×17 km² for 499 a single CMC. 500

In recent work of Gottwald et al. (2014), data-driven methods, similar to our approach, 501 are used to parameterize deep convection. Observational data of a radar located at Kwajalein 502 is used to infer two stochastic processes. A stochastic process for which samples are drawn at 503 random from the estimated distributions and a CMC. Both processes are conditioned on ω 504 at 500hPa. The states of the CMC correspond to deep convective area fractions. With both 505 approaches they are able to reproduce deep convective area fractions for the Darwin region. 506 The models are computationally less expensive than our model, because no micro-grid is used 507 and only deep convective area fractions are considered. They point out that for the training 508 process of the CMC not enough data is available, since they only use spatially averaged 509 fraction values to train the CMCs. Interestingly, they show that only a small adaptation 510 has to be performed before using the models at a different location than where they have 511 been trained. This supports that also our multicloud model could be used more globally. 512 However, since convection is (in part) location dependent, e.g. the presence of land or sea, 513 our model could be improved by using observations from multiple locations. This could 514 lead to a data-driven parameterization of convection and clouds for the usage in numerical 515 weather and climate prediction models. 516

As the multicloud model was able to reproduce the cloud type area fractions quite well, a natural step is to test this model in a GCM. In Section 2, we mentioned that the deep

convective area fractions σ_4 can be used as a closure for the mass flux at cloud base as in 519 (2). The strong congestus area fractions σ_3 , which also represents convection, can be added 520 using a different updraft velocity, and the same can be done with the moderate congestus 521 fractions σ_2 . When using the fractions only as a mass flux closure, it is assumed that the 522 GCM can further calculate the entire vertical tendency profiles for e.g. heat and moisture. 523 An alternative is to define vertical heat and moisture tendency profiles corresponding to each 524 cloud type (e.g. Khouider et al. (2010)) or explicitly inferring vertical heat and moisture 525 tendency profiles from data as in Dorrestijn et al. (2013b). Another possible application of 526 the model in a GCM is that $\sum_{m>1} \sigma_m$, or $1 - \sigma_1$, can be used in the parameterization of 527 cloud cover. 528

The main weakness of our model is that there is no spatial dependence between the CMCs 529 other than through the large-scale state, which results in too small standard deviations for 530 the CMC fractions when N is chosen to be equal to the number of radar sites. The peak 531 values of the the observational fractions of the cloud types stratiform, moderate congestus 532 and for clear sky are difficult to produce, while keeping N such that the peak values of 533 the deep convective area fractions are good. The standard deviation for the cloud types 534 stratiform, moderate congestus and for clear sky are too small and we saw that the ACFs of 535 the CMCs using N = 100 decrease too much initially (except for stratiform and clear sky). 536 To summarize the strengths of our approach: realistic observational data is used to 537 estimate the model; the CMC cloud type area fraction were shown to be comparable to 538 the observations, which is notable, because we used different data sets for training and 539 validation. Furthermore, we saw that the model can be adapted to the scale of the GCM, 540 giving larger fluctuations when a smaller number of Markov chains is used to produce area 541 fractions. Due to the conditioning, memory effects are build in that are often absent in 542 conventional stochastic convection schemes. Implementation in a GCM for assessing the 543 model in a dynamical environment is possible and it can be improved by using additional 544 data from different locations. 545

546 Acknowledgments.

We are grateful to Karsten Peters and Keith Myerscough for useful comments on the paper. This research was supported by the Division for Earth and Life Sciences (ALW) with financial aid from the Netherlands Organization for Scientific Research (NWO).

REFERENCES

- Arakawa, A., 2004: The cumulus parameterization problem: Past, present, and future. J.
 Climate, 17, 2493–2525.
- Arakawa, A., J.-H. Jung, and C.-M. Wu, 2011: Toward unification of the multiscale modeling
 of the atmosphere. *Atmos. Chem. Phys.*, **11**, 3731–3742.
- ⁵⁵⁶ Bengtsson, L., M. Steinheimer, P. Bechtold, and J.-F. Geleyn, 2013: A stochastic
 ⁵⁵⁷ parametrization for deep convection using cellular automata. Q.J.R. Meteorol. Soc., 139,
 ⁵⁵⁸ 1533–1543.
- ⁵⁵⁹ Buizza, R., M. Milleer, and T. Palmer, 1999: Stochastic representation of model uncertainties ⁵⁶⁰ in the ECMWF ensemble prediction system. *Q.J.R. Meteorol. Soc.*, **125**, 2887–2908.
- ⁵⁶¹ Crommelin, D. and E. Vanden-Eijnden, 2008: Subgrid-scale parameterization with condi⁵⁶² tional markov chains. J. Atmos. Sci., 65, 2661–2675.
- Davies, L., C. Jakob, P. May, V. Kumar, and S. Xie, 2013a: Relationships between the largescale atmosphere and the small-scale convective state for Darwin, Australia. J. Geophys. *Res. Atmos.*, 118, 534,11–545.
- Davies, L., et al., 2013b: A single column model ensemble approach applied to the TWP-ICE
 experiment. J. Geophys. Res., 118, 6544–6563.
- ⁵⁶⁸ Dorrestijn, J., D. Crommelin, J. Biello, and S. Böing, 2013a: A data-driven multi-cloud
 ⁵⁶⁹ model for stochastic parametrization of deep convection. *Phil. Trans. R. Soc. A.*, **371**,
 ⁵⁷⁰ 20120 374.

551

- Dorrestijn, J., D. Crommelin, A. Siebesma, and H. Jonker, 2013b: Stochastic parameteri-571 zation of shallow cumulus convection estimated from high-resolution model data. Theor. 572 Comput. Fluid Dyn., 27, 133–148. 573
- Flato, G., et al., 2013: Evaluation of climate models. Climate Change 2013: The Physical 574

Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the

nor, S. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex, and P. Midgley, Eds., Cambridge

- Intergovernmental Panel on Climate Change, T. Stocker, D. Qin, G.-K. Plattner, M. Tig-
- University Press, Cambridge, United Kingdom and New York, NY, USA, 741–866. 578

575

576

- Frenkel, Y., A. Majda, and B. Khouider, 2013: Stochastic and deterministic multicloud 579 parameterizations for tropical convection. Clim. Dyn., 41, 1527–1551. 580
- Gottwald, G. A., K. Peters, and L. Davies, 2014: Data-driven stochastic subgrid-scale 581 parametrization for tropical convection. Submitted. 582
- Johnson, R., T. Rickenbach, S. Rutledge, P. Ciesielski, and W. Schubert, 1999: Trimodal 583 characteristics of tropical convection. J. Climate, 12, 2397–2418. 584
- Khouider, B., J. Biello, and A.J.Majda, 2010: A stochastic multicloud model for tropical 585 convection. Comm. Math. Sci., 8, 187–216. 586
- Khouider, B. and A. Majda, 2006: A simple multicloud parameterization for convectively 587 coupled tropical waves. I. Linear Analysis. J. Atmos. Sci., 63, 1308–1323. 588
- Kumar, V., C. Jakob, A. Protat, P. May, and L. Davies, 2013: The four cumulus cloud modes 589 and their progression during rainfall events: A C-band polarimetric radar perspective. J. 590 Geophys. Res. Atmos., 118, 8375–8389. 591
- Kwasniok, F., 2012: Data-based stochastic subgrid-scale parametrisation: an approach using 592 cluster-weighted modelling. Phil. Trans. R. Soc. A., 370, 1061–1086. 593

- Lang, S., W.-K. Tao, J. Simpson, and B. Ferrier, 2003: Modeling of convective-stratiform precipitation processes: Sensitivity to partitioning methods. J. Appl. Meteor., 42, 505– 527.
- ⁵⁹⁷ Lin, J.-B. and J. Neelin, 2000: Influence of a stochastic moist convective parameterization ⁵⁹⁸ on tropical climate variability. *Geophys. Res. Lett.*, **27**, 3691–3694.
- Lin, J.-B. and J. Neelin, 2003: Toward stochastic deep convective parameterization in general circulation models. *Geophys. Res. Lett.*, **30**, 1162.
- Lin, J.-L., et al., 2006: Tropical intraseasonal variability in 14 IPCC AR4 climate models.
 Part I: Convective signals. J. Climate, 19, 2665–2690.
- Lorenz, E., 1996: Predictability a problem partly solved. *Proceedings: Seminar on Predictability*, Reading, United Kingdom, 1–18.
- Majda, A., S. Stechmann, and B. Khouider, 2007: MaddenJulian oscillation analog and
 intraseasonal variability in a multicloud model above the equator. *Proc. Natl. Acad. Sci. USA*, **104**, 9919–9924.
- Mapes, B., S. Tulich, J. Lin, and P. Zuidema, 2006: The mesoscale convective life cycle: building block or prototype for large-scale tropical waves? *Dyn. Atmos. Oceans*, **42**, 3–29.
- May, P. and A. Ballinger, 2007: The statistical characteristics of convective cells in a monsoon
 regime (Darwin, Northern Australia). *Mon. Weather Rev.*, 135, 82–92.
- Möbis, B. and B. Stevens, 2012: Factors controlling the position of the Intertropical Convergence Zone on an aquaplanet. J. Adv. Model. Earth Syst., 4, M00A04.
- Palmer, T., 2001: A nonlinear dynamical perspective on model error: a proposal for nonlocal stochastic-dynamic parametrization in weather and climate prediction models. Q.J.R.
 Meteorol. Soc., 127, 279–304.

- Peters, K., C. Jakob, L. Davies, B. Khouider, and A. Majda, 2013: Stochastic behavior of
 tropical convection in observations and a multicloud model. J. Atmos. Sci., 70, 3556–3575.
- Plant, R. and G. Craig, 2008: A stochastic parameterization for deep convection based on
 equilibrium statistics. J. Atmos. Sci., 65, 87–105.
- Press, W. H., S. Teukolsky, W. Vetterling, and B. Flannery, 1992: Numerical Recipes in C:
 The Art of Scientific Computing. Cambridge Univ. Press, 994 pp.
- Randall, D., M. Khairoutdinov, A. Arakawa, and W. Grabowski, 2003: Breaking the cloud
 parameterization deadlock. *Bull. Amer. Meteor. Soc.*, 84, 1547–1564.
- ⁶²⁵ Siebesma, A., 1998: Shallow cumulus convection. Buoyant Convection in Geophysical Flows,
- E. Plate, E. Fedorovich, X. Viegas, and J. Wyngaard, Eds., Kluwer, 441–486.
- Teixeira, J. and C. Reynolds, 2008: Stochastic nature of physical parameterizations in ensemble prediction: A stochastic convection approach. *Mon. Wea. Rev.*, **136**, 483–496.
- Wilks, D., 2005: Effects of stochastic parameterizations in the Lorenz '96 system. Q.J.R.
 Meteorol. Soc., 131, 389–407.
- ⁶³¹ Yu, X. and T.-Y. Lee, 2010: Role of convective parameterization in simulations of a convec-
- tion band at grey-zone resolutions. Tellus A, 62, 617-632.

⁶³³ List of Tables

⁶³⁴ 1 Cloud type classification using thresholds for the cloud top height and the ⁶³⁵ rain rate.

TABLE 1. Cloud type classification using thresholds for the cloud top height and the rain rate.

CTH [km]	rain rate $[mm h^{-1}]$	
	≤ 12	> 12
≥ 6.5	stratiform $(m = 5)$	deep convective $(m = 4)$
	≤ 3	> 3
$\in [1.5, 6.5)$	moderate congestus $(m = 2)$	strong congestus $(m = 3)$
< 1.5	clear (m = 1)	

List of Figures

637	1	(a) A snapshot of the cloud top height derived from Darwin radar observations	
638		and (b) the corresponding rain rate.	31
639	2	Histograms of (a) the cloud top height and (b) the rain rate observed with the	
640		Darwin radar in the periods November 2005 - April 2006 and January-April	
641		2007.	32
642	3	Example of radar data assigned to the categories clear sky, moderate conges-	
643		tus, strong congestus, deep convective and stratiform, corresponding to the	
644		CTH and rain rate snapshots of Fig. 1.	33
645	4	Cross-correlation functions (CCFs) of the deep convective area fraction with	
646		- $\langle \omega \rangle$, CAPE and RH at 640 hPa for the 2005/2006 data set.	34
647	5	Histogram of the 25 intervals of $\langle \omega \rangle$, found by clustering the linearly inter-	
648		polated $\langle \omega \rangle$ values. The first and last (25th) intervals are open on one side.	
649		Because ω is a velocity in terms of pressure, positive $\langle\omega\rangle$ corresponds to down-	
650		ward mean large-scale motion and negative $\langle\omega\rangle$ to upward mean motion (as	
651		illustrated by the arrows).	35
652	6	Observational mean cloud type area fractions as a function of the $\langle\omega\rangle$ intervals	
653		for the $2005/2006$ training period (solid line with circles) plus and minus	
654		the standard deviation (dash-dotted line) and the CMC expected cloud type	
655		area fractions (solid line) plus and minus the standard deviation while using	
656		${\cal N}=100$ CMCs (dashed line). Note the different scaling on the y-axis.	36
657	7	(a) Deep convective area fractions observed in Darwin (b,c) two realizations of	
658		deep convective area fractions produced by $N=100~{\rm CMCs}$ conditioned on $\langle\omega\rangle$	
659		and (d) the corresponding histograms comparing the CMC fractions (averaged	
660		over 100 realizations) with the observed fractions (binned into intervals) on a	
661		logarithmic y-axis.	37

662	8	Deep convective area fractions produced by $N = 100$ CMCs conditioned on	
663		CAPE and (b) the corresponding histograms in which the CMC fractions (av-	
664		eraged over 100 realizations) are compared to the observed fractions (binned	
665		into intervals) on a logarithmic y-axis.	38
666	9	Deep convective area fractions produced by $N = 69^2$ CMCs conditioned on	
667		$\langle\omega\rangle$ and (b) the corresponding histograms of the binned CMC fractions av-	
668		eraged over 100 realizations compared to the binned observed fractions on a	
669		logarithmic y-axis.	39
670	10	Area fractions of (a) deep convective, (b) strong congestus, (c) stratiform,	
671		(d) moderate congestus and (e) clear sky observed in Darwin (dashed line),	
672		produced by 100 CMCs (solid line) conditioned on $\langle\omega\rangle$ and the corresponding	
673		expected area fractions of the CMCs (dash-dotted line) for a period of 5 days.	
674		Note the different scaling on the y-axis.	40
675	11	Normalized ACFs of the observational area fractions (solid lines with stars),	
676		the CMC area fractions with $N=100$ conditioned on $\langle\omega\rangle$ (solid lines) and on	
677		CAPE (dashed lines), the ACF corresponding to $69^2~{\rm CMCs}$ conditioned on	
678		$\langle\omega\rangle$ (dotted lines) for the cloud types (a) deep convective (b) strong congestus	
679		(c) stratiform (d) moderate congestus and (e) clear sky. Also the ACF of $\langle\omega\rangle$	
680		is shown (dash-dotted lines).	41



FIG. 1. (a) A snapshot of the cloud top height derived from Darwin radar observations and (b) the corresponding rain rate.



FIG. 2. Histograms of (a) the cloud top height and (b) the rain rate observed with the Darwin radar in the periods November 2005 - April 2006 and January-April 2007.



FIG. 3. Example of radar data assigned to the categories clear sky, moderate congestus, strong congestus, deep convective and stratiform, corresponding to the CTH and rain rate snapshots of Fig. 1.



FIG. 4. Cross-correlation functions (CCFs) of the deep convective area fraction with - $\langle \omega \rangle$, CAPE and RH at 640 hPa for the 2005/2006 data set.



FIG. 5. Histogram of the 25 intervals of $\langle \omega \rangle$, found by clustering the linearly interpolated $\langle \omega \rangle$ values. The first and last (25th) intervals are open on one side. Because ω is a velocity in terms of pressure, positive $\langle \omega \rangle$ corresponds to downward mean large-scale motion and negative $\langle \omega \rangle$ to upward mean motion (as illustrated by the arrows).



FIG. 6. Observational mean cloud type area fractions as a function of the $\langle \omega \rangle$ intervals for the 2005/2006 training period (solid line with circles) plus and minus the standard deviation (dash-dotted line) and the CMC expected cloud type area fractions (solid line) plus and minus the standard deviation while using N = 100 CMCs (dashed line). Note the different scaling on the y-axis.



FIG. 7. (a) Deep convective area fractions observed in Darwin (b,c) two realizations of deep convective area fractions produced by N = 100 CMCs conditioned on $\langle \omega \rangle$ and (d) the corresponding histograms comparing the CMC fractions (averaged over 100 realizations) with the observed fractions (binned into intervals) on a logarithmic y-axis.



FIG. 8. Deep convective area fractions produced by N = 100 CMCs conditioned on CAPE and (b) the corresponding histograms in which the CMC fractions (averaged over 100 realizations) are compared to the observed fractions (binned into intervals) on a logarithmic y-axis.



FIG. 9. Deep convective area fractions produced by $N = 69^2$ CMCs conditioned on $\langle \omega \rangle$ and (b) the corresponding histograms of the binned CMC fractions averaged over 100 realizations compared to the binned observed fractions on a logarithmic y-axis.



FIG. 10. Area fractions of (a) deep convective, (b) strong congestus, (c) stratiform, (d) moderate congestus and (e) clear sky observed in Darwin (dashed line), produced by 100 CMCs (solid line) conditioned on $\langle \omega \rangle$ and the corresponding expected area fractions of the CMCs (dash-dotted line) for a period of 5 days. Note the different scaling on the y-axis.



FIG. 11. Normalized ACFs of the observational area fractions (solid lines with stars), the CMC area fractions with N = 100 conditioned on $\langle \omega \rangle$ (solid lines) and on CAPE (dashed lines), the ACF corresponding to 69² CMCs conditioned on $\langle \omega \rangle$ (dotted lines) for the cloud types (a) deep convective (b) strong congestus (c) stratiform (d) moderate congestus and (e) clear sky. Also the ACF of $\langle \omega \rangle$ is shown (dash-dotted lines).