



1 **Multi-Model Ensemble Forecast of Precipitation Based on**
2 **an Object-Based Diagnostic Evaluation**

3 Luying Ji ^{1,2}, Xiefei Zhi ^{*1,3}, Clemens Simmer², Shoupeng Zhu ¹, Yan Ji ¹,

4 ¹ *Key Laboratory of Meteorological Disasters, Ministry of Education (KLME) /*

5 *Collaborative Innovation Center on Forecast and Evaluation of Meteorological*

6 *Disasters (CIC-FEMD), Nanjing University of Information Science & Technology,*

7 *Nanjing, China*

8 ² *Institute for Geosciences – Section Meteorology, University of Bonn, Bonn,*

9 *Germany*

10 ³ *Nanjing Joint Center for Atmospheric Research (NJCAR), Nanjing, China*

*Corresponding author: Xiefei ZHI

Email: zhi@nuist.edu.cn

Early Online Release: This preliminary version has been accepted for publication in *Monthly Weather Review*, may be fully cited, and has been assigned DOI 10.1175/MWR-D-19-0266.1. The final typeset copyedited article will replace the EOR at the above DOI when it is published.

ABSTRACT

11

12 We analyzed 24-hour accumulated precipitation forecasts over the four-months
13 period from 1 May to 31 August 2013 over an area located in East Asia covering the
14 region 70.15°E – 139.95°E , 15.05°N – 58.95°N generated with the Ensemble Prediction
15 Systems (EPSs) from ECMWF, NCEP, UKMO, JMA and CMA contained in the
16 TIGGE dataset. The forecasts are first evaluated with the Method for Object-based
17 Diagnostic Evaluation (MODE). Then a multi-model ensemble (MME) forecast
18 technique based on weights derived from object-based scores is investigated and
19 compared with the equally-weighted MME and the traditional gridpoint-based MME
20 forecast using weights derived from the point-to-point metric, mean absolute error
21 (MAE).

22 The object-based evaluation revealed that attributes of objects derived from the
23 ensemble members of the five individual EPS forecasts and the observations differ
24 consistently. For instance, their predicted centroid location is more southwestward,
25 their shape is more circular, and their orientation is more meridional than in the
26 observations. The sensitivity of the number of objects and their attributes to
27 methodological parameters is also investigated.

28 A MME prediction technique based on weights computed from the object-based
29 scores, Median of Maximum Interest (MMI) and Object-based Threat Score (OTS), is
30 explored and the results compared to the ensemble forecasts of the individual EPS, the
31 equally-weighted MME forecast, and the traditional super-ensemble forecast. When

32 using MODE statistics for the forecast evaluation, the object-based MME prediction
33 outperforms all other predictions. This is mainly because of a better prediction of the
34 objects' centroid locations. When using the precipitation-based fractions skill score
35 (FSS), which is not used in either of the weighted MME forecasts, the object-based
36 MME forecasts are slightly better than the equally-weighted MME forecasts but inferior
37 to the traditional super-ensemble forecast based on weights derived from the point-to-
38 point metric, MAE.

39 **Key words:** The Method for Object-based Diagnostic Evaluation (MODE), 24-h
40 accumulated precipitation, multi-model ensemble forecasts

41 **1. Introduction**

42 In the past two decades, numerical weather forecasting rapidly developed and -
43 besides model improvements - evolved from traditional single deterministic forecasts
44 to ensemble forecasting (Gneiting and Raftery 2005; Bauer et al. 2015). Different
45 forecast systems differ in their overall architecture, spatial resolution, choice of initial
46 conditions, data assimilation technology, and physical parameterization schemes used
47 in the numerical models. Multi-model ensemble (MME) forecasting is an effective way
48 to make use of the forecasts from different Ensemble Prediction Systems (EPS) with
49 the goal being to reduce systematic deviations from observations and thus improve the
50 overall prediction skill. Based on The Observing System Research and Predictability
51 Experiment (THORPEX) program, which provides forecasts from different operational
52 numerical weather prediction (NWP) centers, MME forecasting is currently already
53 widely used. Many studies have shown that the MME forecast performance is superior
54 to the forecast of an individual (one-model-based) EPS (Krishnamurti et al. 1999;
55 Fraley et al. 2010; Zhi et al. 2012; Zhang et al. 2015; He et al. 2015; Ji et al. 2019).

56 Besides the equally-weighted MME, more complex MME methods, such as linear
57 regression (Krishnamurti et al. 1999; 2000), Bayesian model averaging (BMA; Raftery
58 et al. 2005; Vrugt et al. 2006), ensemble MOS (EMOS; Scheuerer 2014; Scheuerer and
59 Hamill 2015) and artificial neural networks (Yuan et al. 2007; Bakhshaii and Stull
60 2009) have been proposed and are already widely used for precipitation forecasting
61 (Tebaldi et al. 2004; Ke et al. 2008), typhoon forecasts (Kumar et al. 2003; Jordan et

62 al. 2008), and regional climate predictions (Kharin and Zwiers 2002; Yun et al. 2005).
63 Many studies suggest that unequally-weighted MME forecasts can achieve better skill
64 than equally-weighted ones (Chen et al. 2010; Zhang and Zhi 2015; Kim and Chan
65 2018). Peng et al. (2002) and Ke et al. (2009) show, however, that they are not always
66 better and may even be worse than the best individual EPS forecast.

67 Unequally-weighted MME methods often determine the weight of each contributing
68 EPS by their relative performance during a training period, which assumes a certain
69 temporal stability of their forecast performance. Most methods use scores derived from
70 point-to-point comparisons between forecasts and observations, e.g. the weighted
71 ensemble mean (WEMN, Nohara et al. 2006), the bias-removed ensemble mean
72 (BREM, Kharin and Zwiers 2002), and the super-ensemble (SUP, Krishnamurti et al.
73 2000), which uses the mean absolute error (MAE) during a training period.

74 However, point-to-point verification scores (e.g. MAE or Equitable Threat Score
75 (ETS)) provide only limited information about the quality of a precipitation forecast
76 because they only compare the observations and predictions point by point without
77 taking, for example, the resemblance of spatial patterns into account (Mass et al. 2002;
78 Baldwin and Kain 2006; Gilleland et al. 2009). Precipitation is highly discontinuous in
79 space and time. Thus, even almost perfect forecasts of e.g. the shapes and sizes of
80 precipitation systems may lead to poor point-to-point scores because of many false
81 alarms and misses known as “double penalty” already at small spatial deviations.
82 However, the correct prediction of spatial features like shape, size and approximate

83 location of extended precipitation fields are important because they can be used as a
84 valuable guidance to improve forecasts, especially of extreme weather.

85 Several methods have been used to get around the limitations of point by point
86 verification, which categorize into filtering and displacement methods. Filtering
87 methods generally apply smoothing or scale separation to evaluate the forecast for
88 different spatial scales (Marsigli et al. 2006; Roberts and Lean 2008; Ebert 2008; 2009;
89 Casati et al. 2004; 2009; Zepeda-Arce et al. 2000; Harris et al. 2001; Mittermaier 2006;
90 Marzban and Sandgathe 2009), while displacement methods identify discrete features
91 or objects in the forecast and the observations and quantify their respective
92 displacements in terms of location or other attributes (Ebert and McBride 2000;
93 Baldwin and Lakshmiarahan 2003; Keil and Craig 2007; Marzban and Sandgathe
94 2008; Gilleland et al. 2010).

95 The Method for Objected-based Diagnostic Evaluation (MODE) developed by Davis
96 et al. (2006a) is adopted for calculating verification scores in this study. MODE is a
97 typical feature-based displacement approach and an example for a spatial diagnostic
98 technique. MODE attempts to mimic the way a human would subjectively evaluate a
99 forecast via setting a precipitation threshold and spatially convoluting (scale-dependent
100 averaging) the precipitation field. The median of maximum interest (MMI; Davis et al.
101 2009) and the object-based threat score (OTS; Johnson et al. 2011a) are two scores
102 calculated from the attributes of the detected objects in the forecast and in the

103 observations (for more detail see Section 2.3), which are sensitive to different aspects
104 of forecast accuracy (Johnson and Wang 2013).

105 Although MODE has been most commonly used for the verification of high-
106 resolution model forecasts of convective storms, it can also be applied to lower-
107 resolution numerical model weather forecasts, regional climate simulations, or
108 chemistry model simulations (e.g. Brown et al. 2007; Wolff et al. 2014; Li et al. 2015).
109 Twenty-four hour accumulated precipitation over areas of hundreds of kilometers
110 exhibits characteristic spatial patterns that should be reproduced by model forecasts.
111 Object-based methods allow us to evaluate whether this is indeed the case. In this study,
112 we focus on daily precipitation forecasts with lead times of 1–7 days and venture to
113 improve their prediction skills of shape, size and/or location by a new MME approach,
114 which employs weights derived from object-based scores. We will compare its quality
115 with the predictions achieved by the individual EPS, by an equally-weighted MME
116 forecast, and by an MME forecast using weights based on point-to-point metric, MAE,
117 derived from the precipitation forecasts and the observations during a training period.

118 First, we evaluate ensemble forecasts of 24-h accumulated precipitation produced by
119 the five ensemble prediction systems (EPSs, i.e. the European Centre for Medium-
120 Range Weather Forecasts (ECMWF), the National Center for Environment Prediction
121 (NCEP), the UK Met Office (UKMO), the Japan Meteorological Agency (JMA), and
122 the China Meteorological Administration (CMA)). For each ensemble member forecast
123 of each individual EPS, MODE is used to obtain the attributes of every identified object.

124 Various attributes of one identified object are compared to the attributes of the best
125 corresponding object in the observations; then the performance of each EPS is
126 represented by the object attribute differences averaged over all objects identified for
127 each ensemble member of an individual contributing EPS.

128 Second, three MME predictions are computed and their forecast accuracy evaluated
129 using the spatial object-based measures MMI and OTS, but also using the FSS as an
130 independent skill score that was not used for calculating weights in any of the three
131 MME forecasts. We compare the three MME techniques in order to investigate if the
132 MME forecast with object-based weights provides more accurate spatial information in
133 the precipitation forecast.

134 The remainder of this paper is structured as follows. Section 2 briefly describes the
135 datasets that were used and introduces MODE. In section 3, we present the performance
136 evaluation of the five individual EPSs and of the three MME precipitation forecasting
137 methods. A discussion and major conclusions are provided in section 4.

138 **2. Data and methods**

139 ***2.1 Data***

140 We used 24-h accumulated precipitation ensemble forecasts produced by ECMWF,
141 NCEP, UKMO, JMA and CMA at $0.5^\circ \times 0.5^\circ$ resolution initialized daily at 1200 UTC
142 for lead times of 1–7 days (Table 1). The data is available from the TIGGE-ECMWF
143 portal (<http://apps.ecmwf.int/datasets/data/tigge>). TIGGE (The THORPEX Interactive

144 Grand Global Ensemble) is a key component of the THORPEX program; it contains
145 ensemble forecast data from 10 global model prediction centers and has been widely
146 used for scientific research on ensemble forecasting, predictability and the development
147 of products to improve the prediction of severe weather (Breivik et al. 2014; Loeser et
148 al. 2017; Parsons et al. 2017). We analyzed the data for a four-months period from 1
149 May to 31 August 2013 and over an area located in East Asia covering the region
150 70.15°E–139.95°E, 15.05°N–58.95°N.

151 For forecast validation we selected a high-resolution gridded dataset of hourly
152 precipitation, which merged precipitation analyses of the U.S. National Oceanic and
153 Atmospheric Administration Climate Prediction Center morphing technique
154 (CMORPH) given at a spatial resolution of 8 km, with the Chinese gauge-based
155 precipitation analysis based on about 30,000 automatic weather stations. This merged
156 gauge–satellite precipitation product (available at
157 http://data.cma.cn/data/detail/dataCode/SEVP_CLI_CHN_MERGE_CMP_PRE_HO
158 [UR_GRID_0.10/](http://data.cma.cn/data/detail/dataCode/SEVP_CLI_CHN_MERGE_CMP_PRE_HO)) with a resolution of 0.1° x 0.1° used optimal interpolation and the
159 probability density functions of both products, and has been proved to be superior to
160 other similar international products over China (Xie and Xiong 2011; Pan et al. 2012).
161 The verification data were interpolated to 0.5° x 0.5° resolution by bilinear interpolation
162 (Rauscher et al. 2010; Kopparla et al. 2013; Ahmed et al. 2019).

163 ***2.2 Method for Object-Based Diagnostic Evaluation (MODE)***

164 MODE sets weight and confidence coefficients for predefined precipitation object
165 attributes and calculates a total interest function based on a fuzzy logic approach, which
166 quantifies the similarity between any two objects (Davis et al. 2006a; Johnson et al.
167 2013). The predefined attributes are chosen by a particular user for a particular
168 application. In general, MODE consists of four steps: identifying objects, calculating
169 object attributes, finding matching objects between observations and predictions, and
170 assessing the similarity of their attributes.

171 *2.2.1 Identifying Objects and object attributes*

172 In order to extract the spatial boundary of an object, the original precipitation field is
173 spatially smoothed with a convolution radius R (unit: grid points). Then an intensity
174 threshold T (unit: mm (24 h)^{-1}) is used to define the boundaries of precipitation objects
175 (Davis et al. 2006a). The original precipitation field within these boundaries then
176 defines the precipitation objects, which are solely determined by the selection of the
177 convolution radius R , which is related to the precipitation scale, and the threshold T ,
178 which is related to the precipitation intensity. These two parameters can be chosen
179 based on the scales of interest. The result of each step is demonstrated in Fig. 1.

180 We usually pay attention to the overall location of a precipitating system, its size and
181 its shape, especially when dealing with more extreme weather (Johnson and Wang
182 2013). Therefore, the specific attributes used in our study are the area coverage of
183 precipitation objects, their aspect ratio (the ratio of minor axis to major axis; i.e. 1.0 for
184 a circular object and <1 otherwise) and orientation angle (the orientation of the major

185 axis in degrees counterclockwise starting at zonal orientation), and their centroid
186 location. For matched object pairs (introduced in Section 2.2.2), attribute differences in
187 the four mentioned object attributes (Table 2) are calculated.

188 2.2.2 Object matching

189 Object matching creates a pair consisting of one object in the forecasted field and
190 one object in the observed field. Here, we followed Davis et al. (2006a), who
191 determined paired objects solely based on their centroid distance D and their areas. If
192 $D < (Area_o^{1/2} + Area_f^{1/2})/2$ with $Area_o$ and $Area_f$ the areas of the observed
193 object and the forecasted object, respectively, both objects create a matched pair.
194 Thus, a matching object pair requires the centroid distance between both to be less than
195 their average size.

196 2.3 Quantification of similarity of matched object pairs

197 For a matched pair, its total interest I is computed via

$$198 \quad I = \frac{\sum_{i=1}^n \omega_i c_i G_i}{\sum_{i=1}^n \omega_i c_i} \quad (1)$$

199 c_i and ω_i are the confidence value and the weight of the attribute i , respectively, and
200 n is the number of attributes used. While the weight depends only on the particular
201 attribute, the confidence value varies with the sizes and distances of the paired objects
202 (Table 2). G_i is the interest value of the matched objects in terms of attribute i ; it
203 quantifies the degree of similarity between the objects for that attribute as a monotonic
204 function decreasing from 1 to 0 as the attribute dissimilarity increases (Fig. 2).

205 2.4 Quantification of object-based forecast accuracy

206 The Median of Maximum Interest (MMI; Davis et al. 2009) and the fuzzy Object-
 207 based Threat Score (OTS; Johnson et al. 2011a) are two metrics used to quantify the
 208 similarity of the objects in the forecasted and observed fields. The MMI proposed by
 209 Davis et al. (2009), which is called the standard MMI in the following, is the median of
 210 the maximum total interests in the forecasted and observed fields to which all objects
 211 contribute equally regardless of size. The MMI calculated in our study will be slightly
 212 larger than the standard MMI, because we first determine the matched objects by their
 213 centroid distance and areas, and then the total interest I is only calculated for the
 214 matched pairs. Thus, unmatched objects are not considered.

215 The OTS is the fraction of the area of all objects that is contained in matched objects,
 216 multiplied by their total interests:

$$217 \quad OTS = \frac{\sum_{p=1}^P I^p (a_f^p + a_o^p)}{A_f + A_o} \quad (2)$$

218 with P the total number of objects pairs, A_f and A_o the total area of all objects in the
 219 forecasted and the observed field, respectively, and a_f^p and a_o^p ($p=1,2,\dots,P$) the areas
 220 of the p -th paired objects in the forecasted and observed field, respectively. I^p is the
 221 total interest value for p -th matched pair. According to Eq. (2) the OTS takes the object
 222 area and the number of matched objects into account. Thus, larger objects will
 223 contribute more to the OTS than smaller objects, while over-forecasting or under-
 224 forecasting the number of objects will decrease the OTS due to more unmatched objects.
 225 Both indices range between 0 and 1 and have a value of 1.0 for perfect forecasts. Both

226 scores are used in the present study to quantify two complementary aspects of forecast
227 accuracy.

228 The MMI of each EPS is the median computed from the 51 members of ECMWF,
229 21 members of NCEP, 24 members of UKMO, 51 members of JMA, and 15 members
230 of CMA, respectively. The OTS of each EPS is the average OTS of its ensemble
231 members. Both EPS-scores computed over a training period are used as weights to
232 construct the object-based MME prediction as an alternative to point-to-point metrics
233 as used in the classical approach (see next section).

234 **2.5 Different multi-model ensemble types**

235 *2.5.1 Traditional grid point-based multi-model ensemble*

236 Super-ensembles have the potential to improve weather and climate forecast skills
237 above individual ensemble forecasts (Kim et al. 2010; Johnson et al. 2014;
238 Krishnamurti et al. 2016). They automatically remove the bias between the observations
239 and model forecasts estimated during a training period, which contributes to the
240 improved prediction skill of multi-model forecasting. In this study, the point-to-point
241 weighted multi-ensemble forecast is defined as:

$$242 \quad SUP_j = \bar{O} + \sum_{i=1}^N \delta_i (Y_{ij} - \bar{Y}_i) \quad (3)$$

$$243 \quad \delta_i = \left(\frac{1}{T} \sum_{t=1}^T |Y_{it} - O_t| \right)^{-1} / \sum_{i=1}^N \left(\frac{1}{T} \sum_{t=1}^T |Y_{it} - O_t| \right)^{-1} \quad (4)$$

244 with \bar{O} and \bar{Y}_i respectively the average observed and forecasted value by EPS i
245 ($i=1,2,\dots,5$) computed over a training period T in days, and Y_{ij} the forecast of EPS i

246 on day j of the forecast period. δ_i is the individual contributing EPS weight, with Y_{it}
 247 the forecast value of the i -th EPS on day t ($t=1,2,\dots,T$), O_t the respective observation,
 248 and N the total number of used EPS (in our case $N=5$).

249 2.5.2 Object-based multi-model ensemble

250 In this study, the weights for MME forecasts are also calculated by object-scores (i.e.
 251 MMI and OTS). As described before, first, the precipitation object with its several
 252 object attributes is identified by MODE. Second, one object in the observed field will
 253 be matched to one object in the forecast field by satisfying the matching criteria. Third,
 254 the similarity between the two matched objects is determined on the basis of the
 255 differences in their attributes. Fourth, the similarity values are used to calculate the
 256 object-based metrics MMI and/or OTS from which the object-based scores for each
 257 EPS are obtained. The performance of each EPS during a training period determines its
 258 weight. Tests identified a sliding window of 30 days before the forecast period as the
 259 optimal training period. During the training period, for a certain EPS, each ensemble
 260 member is evaluated by MODE, and then MMI and/or OTS of this EPS is calculated
 261 by the median and/or mean of all ensemble members. Finally, the multi-model
 262 ensemble forecasts MME_{MMI} or MME_{OTS} are determined by multiplying the ensemble
 263 mean of each contributing EPS by the weight calculated for the training period as
 264 follows:

$$265 \quad MME_{MMI/OTS} = \sum_{i=1}^N \delta_i^{MMI/OTS} Y_i \quad (5)$$

$$266 \quad \delta_i^{MMI} = \frac{1}{T} \sum_{t=1}^T MMI_{i,t} / \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T MMI_{i,t} \quad (6)$$

267
$$\delta_i^{OTS} = \frac{1}{T} \sum_{t=1}^T OTS_{i,t} / \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T OTS_{i,t} \quad (7)$$

268 with N the number of EPS, Y_i the ensemble mean for the i -th EPS, and $\delta_i^{MMI/OTS}$ the
 269 weight of each contributing EPS calculated by MMI or OTS. T is the length of the
 270 training period in days, and $MMI_{i,t}$ or $OTS_{i,t}$ is the MMI or OTS value for i -th EPS
 271 on day t during the training period.

272 **2.6 Fractions skill score**

273 Besides the MMI and OTS, the Fractions Skill Score, FSS (Roberts and Lean 2007;
 274 Roberts 2008), which is not used to generate the weights in the tested MMEs, is applied
 275 to evaluate the forecast skill for individual EPS and the MME forecasts. This spatial
 276 verification score quantifies forecast skill over different spatial scales. FSS is calculated
 277 based on the fractional coverage within a square neighborhood centered on each grid
 278 point. FSS requires for a given spatial scale, s , the forecasted and the observed areal
 279 fractions M_i and O_i at each grid point, respectively, with precipitation above a given
 280 threshold, and is calculated for an area divided into N sub-areas of size $s \times s$ as follows:

281
$$FSS = 1 - \frac{FBS}{FBS_{worst}} \quad (7)$$

282
$$FBS = \frac{1}{N} \sum_{i=1}^N (O_i - M_i)^2 \quad (8)$$

283
$$FBS_{worst} = \frac{1}{N} (\sum_{i=1}^N O_i^2 + \sum_{i=1}^N M_i^2) \quad (9)$$

284 FBS_{worst} is the largest Fractions Brier Score (FBS), which indicates the case when
 285 there are no common non-zero fractions between predictions and observations. The FSS
 286 ranges between 0 and 1; 0 stands for a totally mismatched forecast and 1 for a perfect
 287 forecast.

288 **3. Results**

289 **3.1 Individual objects**

290 Convolution radius R and threshold value T are the only two parameter that influence
291 object recognition and thus affect the values of object attributes such as object number,
292 area, and centroid location. In this section, we analyze the effect of the choice of R and
293 T on the object attributes. Since the effective resolution of the model dynamics is about
294 seven grid points (Skamarock 2004) and precipitation is generated grid point-wise in
295 the model by the action of parameterizations, we have chosen 3 grid points for the
296 minimum convolution radius R as compromise. Since larger R values will smooth out
297 especially the interesting heavy precipitation areas, we analyze the impact of different
298 R in the intermediate range between 3 and 6 grid points. Only the results for 24-h
299 forecasts are shown; results for other lead times are qualitatively similar.

300 The variation of the number of objects and their areas with precipitation thresholds
301 and convolution radii in the observations and the forecasts of the ECMWF EPS is
302 displayed in Fig. 3. Observations and forecast exhibit similar behavior; e.g. the number
303 of objects first increases with the precipitation threshold T until 5 mm is reached, and
304 then gradually decreases until 25 mm is reached from where the number of objects
305 strongly decrease. Generally, the forecast produces a lower number of objects than the
306 observations, suggesting that the model is more inclined to predict larger continuous
307 precipitation areas. This bias also leads to a large number of false alarms in point-to-
308 point statistics (not shown). The number of objects understandably decreases with

309 increasing convolution radius R in both observations and forecasts. Also the average
310 object areas have similar dependencies on T and R for observations and forecasts, but
311 there are also differences. The average object areas are smaller for the forecasts than
312 for the observations for precipitation less than 5mm. The forecast areas are slightly
313 larger than or equal to the observed areas for higher precipitation thresholds. The
314 average precipitation area is – different from the number of objects - relatively
315 insensitive to the choice of the convolution radius. This is maybe because averaging
316 makes them bigger – but also flatter, so the chosen threshold more or less compensates
317 for that. For a given precipitation threshold, T , the number of objects decreases with
318 increasing object area for observations and the ECMWF EPS (Fig. 4). The decrease in
319 object number with increasing area gets larger for higher precipitation thresholds. For
320 larger precipitation thresholds the forecasts produce significantly less objects with
321 larger areas than the observations (be aware, that only the forecast may produce object
322 numbers below 1, because these values are averages over the ensemble members). The
323 effect of R and T on object number and area is qualitatively similar for the other four
324 models (not shown).

325 In Figure 5 we compare the distributions of several object attributes between
326 observations and 24-hour predictions by all ensemble members of all five EPSs for a
327 convolution radius $R=4$ grid points ($\sim 220\text{km}$) and a precipitation threshold $T=10\text{mm}$ as
328 an example. This qualitative analysis of the observed and forecasted daily precipitation
329 distributions is performed in order to investigate if the different numerical models

330 behind the different EPS do capture the observed spatial features. The number of
331 objects decrease rapidly with increasing area (Fig. 5a) with the models having a lower
332 number of very small areas. The object aspect ratio distributions are broad and peak
333 around 0.6 for the observations and at 0.7 for the model forecasts (Fig. 5b). Most objects
334 have an orientation angle between -30 and 30 degrees with the largest number of objects
335 found around 15 degrees especially for the forecasts, which also have secondary peaks
336 at 90 degrees (Fig. 5c). More objects are found in the southern part of the domain, which
337 is also more pronounced in the forecasted field, while the east-west distribution is more
338 even for both observations and model forecasts (Fig. 5d,e). In general, the forecasted
339 and observed distributions are qualitatively similar, which demonstrates that the spatial
340 features of 24-h accumulated precipitation are captured reasonably well by the
341 numerical models.

342 *3.2 Comparison of matched object attributes*

343 We compare now the object attributes of centroid location, aspect ratio, and
344 orientation angle in the matched object pairs. Figures 6 and 7, respectively, show the
345 mean objects' zonal (i.e. east-west) and meridional (i.e. north-south) centroid
346 differences for the five EPSs compared to the observations for different convolution
347 radii and precipitation thresholds. The mean zonal or meridional centroids of the
348 forecasted objects are generally within 0 to 2 grid points of the observed ones for all
349 EPSs. The forecasted objects from all EPSs are located west to the observed objects for
350 thresholds less than 10mm. But for larger thresholds, the objects of NCEP, UKMO and

351 JMA EPSs are eastward from the observed ones. The predicted objects are consistently
352 southward compared to the observed ones for all thresholds and convolution radii. The
353 aspect ratio deviations between model predictions and observations are always positive
354 but small (Fig. 8) indicating that the shape of the forecasted objects is more circular
355 than for the observed objects. The orientation angle differences are on average within
356 0 to 10 degrees except for the large convolution radii at a threshold of 50mm (Fig. 9).
357 The average positive deviations between forecasted and observed objects indicate a
358 more meridional orientation of the former. In summary, the forecasted objects are more
359 circular, more southwest and have a more meridional orientation than the observed
360 objects. A most likely hypothesis is that these characteristics of forecasted objects are
361 attributed to model dynamics and physics (Johnson et al. 2011b; Johnson and Wang
362 2013).

363 Since the four main attributes are not very sensitive to the choice of R (Fig. 6-9) and
364 the difference of object numbers between the forecast and the observation becomes
365 small when R is larger than 3 grid points (Fig. 3), we choose $R=4$ grid points to
366 investigate the performance of object-based MME forecasting. We have chosen a
367 threshold value of 10 mm for 24-h accumulated precipitation in order to focus on
368 moderate to strong precipitation and exclude light precipitation, which is usually
369 overpredicted in frequency especially by non-convection-permitting model simulations
370 (Giorgi et al. 1992; Golding 2000; Dravitzki and McGregor 2011).

371 ***3.3 Multi-model ensemble forecasting***

372 Many studies have shown advantages of MME prediction over predictions from a
373 single EPS (Candille 2009; Beck et al. 2016; Wanders and Wood 2016; Samuels et al.
374 2018). We first calculate the weights based on the MAE by point-to-point statistics and
375 the MMI or OTS based on MODE, hereafter abbreviated as SUP, MME_{MMI} and
376 MME_{OTS} , respectively. The results of these three MME predictions are weighted
377 ensemble mean forecasts. Thus, they are deterministic forecasts, which we evaluated
378 via MODE taking a threshold of 10mm and a convolution radius of 4 grid points.

379 The object-based scores for both the individual EPSs and the three MME predictions
380 (i.e. MME_{MMI} , MME_{OTS} and SUP) are compared in Fig. 10. As expected, the forecast
381 skill generally decreases with lead time for all predictions. The ECMWF EPS is more
382 skillful than the other EPSs in terms of MMI and OTS and thus it contributes relatively
383 more to the results of MME_{MMI} and MME_{OTS} (Fig. 11). UKMO EPS performs good for
384 lead times of 1–4 days, and NCEP EPS is better for longer lead times. The relative
385 performance of each EPS is shown by their respective weights. The CMA EPS has the
386 lowest scores and thus contributes the least (Fig. 11). The MME predictions weighted
387 by the MMI and OTS metrics perform similarly well, and perform better than both the
388 individual EPSs and the traditional grid point-based MME prediction based on the
389 point-to-point MAE metric for almost all lead times.

390 In order to understand why the MME forecasts based on MMI and OTS are better
391 than the single-model ensemble forecasts and the traditional point-to-point MME, we
392 analyze the four main attribute differences (aspect ratio, orientation angle, zonal and

393 meridional centroid location) between the observed and forecasted fields for all lead
394 times (Fig. 12). For all lead times the forecasted objects are on average more circular
395 than the observed ones (Fig. 12a). The object orientation angles resulting from the
396 traditional super-ensemble forecast are somewhat closer to the observations. The
397 orientation angle of the forecasted objects is on average larger than for the observed
398 objects; thus, the forecasted objects have on average a more meridionally oriented
399 orientation than the observed objects (Fig. 12b). For aspect ratio and orientation angle,
400 the MME_{MMI} and MME_{OTS} forecasts on average are not better than the individual model
401 and traditional point-to-point super-ensemble forecasts, while the centroid locations –
402 both latitude and longitude – are better reproduced by both the MME_{MMI} and MME_{OTS}
403 forecasts and are thus the main reason for their overall better performances given the
404 higher weight the centroid locations get in Eq. (1). The traditional point-to-point super-
405 ensemble forecast is unable to predict the location well in our case, especially for the
406 meridional centroid location. But it still beats some individual EPSs for lead times of 3
407 days and longer. (Figs. 12c,d). The average bias for these four attributes in the MME
408 forecasts is qualitatively similar to the bias of the individual models because all models
409 exhibit similar error characteristics. Accordingly, a MME forecast will suffer from the
410 same errors.

411 We evaluate the equally-weighted MME mean forecasts (EMME) and the two
412 unequally-weighted MME forecasts MME_{OTS} and SUP with the FSS, which is not used
413 for weight determination in the training periods (Fig. 13). The results for MME_{MMI} are

414 similar to those of MME_{OTS} and thus not displayed in Fig. 13. The FSS always increases
415 with scale; accordingly it is easier to predict precipitation probabilities for larger areas.
416 For all spatial scales, the object-based MME forecasting MME_{OTS} is slightly better than
417 the equally-weighted one (EMME) for all lead times. The grid-point-based MME
418 (SUP) provides the best predictions when evaluated with the FSS. There may be two
419 main reasons for this result. First, precipitation objects often have complicated shapes
420 that are not sufficiently represented by the MODE attributes. In this study, only
421 orientation angle and aspect ratio are used to describe the shape of the precipitation
422 object; thus other meaningful precipitation information may be missed. Second, the grid
423 point-based super-ensemble removes the bias of precipitation intensity between the
424 observations and model forecasts, while the object-based MME in our study removes
425 the spatial bias (e.g. centroid location) but not the precipitation intensity bias.

426 **4. Summary and Discussion**

427 Traditional point-to-point verification methods neglect important spatial
428 information, and are usually insensitive to differences in precipitation location and
429 shape errors. Precipitation is regarded as an object by MODE and several object
430 attributes such as number, area, shapes and centroid locations are identified. The
431 differences in object attributes between the model forecasts and the observations could
432 provide important diagnostic information about prediction biases and help forecasters
433 to better use model forecast products.

434 In this study, the ensemble forecasts from five EPSs (ECMWF, NCEP, UKMO, JMA,
435 and CMA EPS) available from the TIGGE datasets are evaluated via object attributes
436 based on MODE. In addition, we investigate a MME technique based on object-based
437 scores and compare it with the equally-weighted multi-model ensemble mean and
438 super-ensemble forecasts based on the point-to-point metric MAE.

439 We first analyze the impact of the convolution radius R and precipitation threshold
440 T on the attributes of the derived precipitation objects. The number of detected objects
441 decreases with increasing convolution radius and precipitation threshold. For all
442 precipitation fields the number of detected objects decreases with increasing object area.

443 In general, the numerical models could capture the distribution of attributes of the
444 observed precipitation objects, and their forecast skill decreases as expected with lead
445 time. The objects aspect ratio varies between 0.3 and 0.9 and the orientation angles are
446 within ± 30 degrees. More objects are found in the eastern/central and southern portion
447 of the domain than in other parts of the domain. In addition, for matched objects -
448 compared to the observed one - the forecasted object centroid positions by all individual
449 model ensembles are more southward and westward. Forecasted objects tend to be more
450 circular and more southwest-northeast orientated compared to the observed ones.
451 Causes for these features of forecast objects are probably related to dynamical errors
452 and model physics.

453 For the five EPSs used in this study, the ECMWF EPS performs best. The MME
454 weighted by the spatial metrics outperforms both all single model EPSs and the

455 traditional point-to-point super-ensemble forecast mainly because of the better
456 forecasted object centroid locations when evaluated using the ensemble mean of the
457 object-based metrics. When all EPSs have similar error characteristics, MMEs will not
458 help much. Thus, the causes for such biases – most probably related to model dynamics
459 and parameterization physics - must be found and the models improved accordingly.

460 When evaluated with the grid point-based (i.e., non-object) metric, FSS, the object-
461 based MME still performs somewhat better than the equally-weighted ensemble mean,
462 but is not as good as the grid point-based MME predictions. This is probably attributed
463 to the use of too few attributes used in our MODE realization and to the inherent bias
464 removal built in the traditional sup-ensemble. MME performance strongly depends on
465 how it is generated, and additional metrics may be used to determine the weights for
466 MME forecasting. Possibly, forecast skill may be further improved by combining
467 different post-processing methods.

468 The rather small differences between object-based and equally-weighted MME
469 forecasts, in terms of MMI and OTS (not shown), are probably due to similar model
470 biases of the five EPSs in our study domain and suggest an extension of such studies to
471 other domains.

472 The precipitation objects are identified in our study from the raw ensemble forecasts
473 without any bias correction. Thus, the object-based scores may be improved by
474 appropriate bias correction. Alternatively, appropriate measures of the objects's
475 precipitation intensity could be developed and added as object attributes both for object

476 pair identification and EPS weight determination and potentially improve the forecast
477 skill of object-based MME above pure grid point-based MME even when evaluated by
478 grid point-based metrics. The object-based MME prediction results may also be further
479 improved by excluding the EPSs performing worst in the training period. In addition,
480 the FSS metric can also be employed to determine the weights of the contributing EPSs.
481 Because precipitation structures become increasingly complex as resolution increases,
482 features such as shape and orientation are hard to define at high resolution, thus the FSS
483 might be an alternative to MODE.

484 **Acknowledgements**

485 This study was supported by the National Natural Science Foundation of China
486 (Grant No. 41575104), the NJCAR key project (Grant No. 2016ZD04), the
487 Postgraduate Research & Practice Innovation Program of Jiangsu Province (Grant No.
488 SJKY19_0934) and the Priority Academic Program Development of Jiangsu Higher
489 Education Institutions (PAPD).

490 **References**

- 491 Bakhshaii, A., and R. Stull, 2009: Deterministic Ensemble Forecasts Using Gene-
492 Expression Programming, *Wea. Forecasting*, **24**, 1431–1451.
- 493 Baldwin, M. E and J. S. Kain, 2006: Sensitivity of several performance measures to
494 displacement error, bias, and event frequency. *Wea. Forecasting*, **21**, 636–648.
- 495 Bauer P., A. Thorpe, G. Brunet, 2015: The quiet revolution of numerical weather
496 prediction. *Nature*, **525**, 47-55.
- 497 Breivik, Ø., O. J. Aarnes, S. Abdalla, J.-R. Bidlot, and P. A. E. M. Janssen, 2014: Wind
498 and wave extremes over the world oceans from very large ensembles, *Geophys.*
499 *Res. Lett.*, **41**, 5122–5131.
- 500 Brown, B. G., R. Bullock, J. H. Gotway, D. Ahijevych, C. Davis, E. Gilleland, and L.
501 Holland., 2007: Application of the MODE object-based verification tool for the
502 evaluation of model precipitation fields. *22nd Conf. on Weather Analysis and*
503 *Forecasting and 18th Conf. on Numerical Weather Prediction*, Park City, Utah.
- 504 Chen, C. H., C. Y. Li, Y. K. Tan, and T. Wang, 2010: Research of the multi-model
505 super-ensemble prediction based on crossvalidation. *J. Meteor. Res.*, **68**, 464–476.
- 506 Davis C. A., B. G. Brown, and R. Bullock, 2006a: Object-based verification of
507 precipitation forecasts. Part I: Methodology and application to mesoscale rain areas.
508 *Mon. Wea. Rev.*, **134**, 1772–1784.

509 Davis, C. A., B. G. Brown, R. Bullock, and J. H. Gotway, 2009: The Method for Object-
510 Based Diagnostic Evaluation (MODE) Applied to Numerical Forecasts from the
511 2005 NSSL/SPC Spring Program. *Wea Forecasting*, **24**, 1252–1267.

512 Dravitzki, S., and J. McGregor, 2011: Predictability of heavy precipitation in the
513 Waikato River basin of New Zealand. *Mon. Wea. Rev.*, **139**, 2184–2197.

514 Fraley, C., A. E. Raftery, and T. Gneiting, 2010: Calibrating multi-model forecast
515 ensembles with exchangeable and missing members using Bayesian model
516 averaging. *Mon. Wea. Rev.*, **138**, 190–202.

517 Gilleland, E., D. Ahijevych, B. G. Brown, B. Casati, and E. E. Ebert, 2009:
518 Intercomparison of Spatial Forecast Verification Methods. *Wea Forecasting*, **24**,
519 1416–1430.

520 Gilleland, E., D. Ahijevych, B. G. Brown, and E. E. Ebert, 2010: Verifying Forecasts
521 Spatially. *Bull. Amer. Meteor. Soc.*, **91**, 1365–1373.

522 Giorgi, F., G. T. Bates, and S. J. Nieman, 1992: Simulation of the arid climate of the
523 southern Great Basin using a regional climate model. *Bull. Amer. Meteor. Soc.*, **73**,
524 1807–1823.

525 Gneiting, T., A. E. Raftery, A. H. Westveld III, and T. Goldman, 2005: Calibrated
526 probabilistic forecasting using ensemble model output statistics and minimum
527 CRPS estimation. *Mon. Wea. Rev.*, **133**, 1098–1118,

528 Golding, B. W., 2000: Quantitative precipitation forecasting in the UK. *J. Hydrol.*, **239**,
529 286–305.

530 He, C. F., X. F. Zhi, Q. L. You, B. Song, and K. Fraedrich, 2015: Multi-model ensemble
531 forecasts of tropical cyclones in 2010 and 2011 based on the Kalman Filter method.
532 *Meteor. Atmos. Phys.* **127**, 467–479.

533 Ji, L. Y., X. F. Zhi, S. P. Zhu and K. Fraedrich, 2019: Probabilistic Precipitation
534 Forecasting over East Asia Using Bayesian Model Averaging. *Wea. Forecasting*,
535 **34**, 377–392.

536 Johnson, A., and X. Wang, 2013: Object-based evaluation of a storm-scale ensemble
537 during the 2009 NOAA Hazardous Weather Testbed Spring Experiment. *Mon.*
538 *Wea. Rev.*, **141**, 1079–1098.

539 Johnson, A., X. Wang, F. Kong, and M. Xue, 2011a: Hierarchical Cluster Analysis of
540 a Convection-Allowing Ensemble during the Hazardous Weather Testbed 2009
541 Spring Experiment. Part I: Development of the Object-Oriented Cluster Analysis
542 Method for Precipitation Fields. *Mon. Wea. Rev.*, **139**, 3673–3693.

543 Johnson, A., X. Wang, F. Kong, and M. Xue, 2011b: Hierarchical Cluster Analysis of
544 a Convection-Allowing Ensemble during the Hazardous Weather Testbed 2009
545 Spring Experiment. Part II: Season-long ensemble clustering and implication
546 for optimal ensemble design. *Mon. Wea. Rev.*, **139**, 694–3710.

547 Johnson, A., X. Wang, F. Kong, and M. Xue, 2013: Object-based evaluation of the
548 impact of horizontal grid points on convection-allowing forecasts. *Mon. Wea. Rev.*,
549 **141**, 3413–3425.

550 Johnson, B., V. Kumar, T. N. Krishnamurti, 2014: Rainfall anomaly prediction using
551 statistical downscaling in a multimodel superensemble over tropical South
552 America. *Climate Dyn.*, **43**, 1731–1752.

553 Jordan, M. R., T. N. Krishnamurti, and C. A. Clayson, 2008: Investigating the utility of
554 using cross-oceanic training sets for superensemble forecasting of eastern Pacific
555 tropical cyclone track and intensity, *Wea. Forecasting*, **23**, 516–522.

556 Ke, Z. J., W. J. Dong, P. Q. Zhang, J. Wang, and T. B. Zhao, 2009: An Analysis of the
557 Difference between the Multiple Linear Regression Approach and the Multimodel
558 Ensemble Mean, *Adv. Atmos. Sci.*, **26**, 1157–1168.

559 Kharin, V. V., and F. W. Zwiers, 2002: Climate predictions with multimodel ensembles.
560 *J. Climate*, **15**, 793–799.

561 Kim, M. K., I. S. Kang, C. K. Park, and Coauthors, 2010: Superensemble prediction of
562 regional precipitation over Korea. *Int. J. Climatol.*, **24**, 777–790.

563 Kim, O. Y., and J. C. L. Chan, 2018: Cyclone-track based seasonal prediction for South
564 Pacific tropical cyclone activity using APCC multi-model ensemble prediction,
565 *Climate Dyn.*, **51**, 3209–3229.

566 Krishnamurti, T. N., C. M. Kishtawal, T. E. Larow, and Coauthors, 1999: Improved
567 Weather and Seasonal Climate Forecasts from Multimodel Superensemble,
568 *Science*, **285**, 1548–1550.

569 Krishnamurti, T. N., C. M. Kishtawal, Z. Zhang, and Coauthors, 2000: Multimodel
570 ensemble forecasts for weather and seasonal climate. *J. Climate*, **13**, 4196–4216.

571 Krishnamurti, T. N., V. Kumar, A. Simon, and Coauthors, 2016: A review of
572 multimodel superensemble forecasting for weather, seasonal climate, and
573 hurricanes. *Rev. Geophys.* **54**, 336–377.

574 Kumar, T. S. V. V., T. N. Krishnamurti, M. Fiorino, and M. Nagata, 2003: Multimodel
575 Superensemble Forecasting of Tropical Cyclones in the Pacific. *Mon. Wea. Rev.*,
576 **131**, 574–583.

577 Li J, K. Hsu, A. AghaKouchak and S. Sorooshian, 2015: An object-based approach for
578 verification of precipitation estimation. *Int. J. Remote Sens.*, **36**, 513–529.

579 Loeser, C. F., M. A. Herrera, and I. Szunyogh, 2017: An assessment of the performance
580 of the operational global ensemble forecast systems in predicting the forecast
581 uncertainty. *Wea. Forecasting*, **32**, 149-164.

582 Mass, C. F., D. Ovens, K. Westrick and B. A. Colle, 2002: Does increasing horizontal
583 resolution produce more skillful forecasts? *Bull. Amer. Meteor. Soc.*, **83**, 407–430.

584 Nohara, D., A. Kitoh, M. Hosaka, et al., 2006: Impact of climate change on river
585 discharge projected by multimodel ensemble. *J. Hydro.*, **7**, 1076–1089.

586 Pan, Y., Y. Shen, J. J. Xu, P. Zhao, 2012: Analysis of the combined gauge-satellite
587 hourly precipitation over China based on the OI technique, *Acta. Meteorologica*
588 *Sinica.*, **70**, 1381-1389. (in Chinese)

589 Parsons, D. B., and Coauthors, 2017: THORPEX Research and the Science of
590 Prediction. *Bull. Amer. Meteor. Soc.*, **98**, 807–830.

591 Peng, P., A. Kumar, and H. Dool, A. G. Barnston, 2002: An analysis of multimodel
592 ensemble predictions for seasonal climate anomalies, *J. Geophys. Res.*, **107(D23)**,
593 ACL-18.

594 Roberts, N., 2008: Assessing the spatial and temporal variation in the skill of
595 precipitation forecasts from an NWP model. *Meteor. Appl.*, **15**, 163–169.

596 Roberts, N. M., and H. W. Lean, 2007: Scale-selective verification of rainfall
597 accumulations from high-resolution forecasts of convective events. *Mon. Wea.*
598 *Rev.*, 136, 78–97.

599 Raftery, A. E., T. Gneiting, F. Balabdaoui, and M. Polakowski, 2005: Using Bayesian
600 model averaging to calibrate forecast ensembles. *Mon. Wea. Rev.*, **133**, 1155–1174.

601 Scheuerer, M., 2014: Probabilistic quantitative precipitation forecasting using
602 ensemble model output statistics. *Quart. J. Roy. Meteor. Soc.*, **140**, 1086–1096.

603 ———, and T. M. Hamill, 2015: Statistical postprocessing of ensemble precipitation
604 forecasts by fitting censored, shifted gamma distributions. *Mon. Wea. Rev.*, **143**,
605 4578–4596.

606 Skamarock, W. C., 2004: Evaluating mesoscale NWP models using kinetic energy
607 spectra. *Mon. Wea. Rev.*, **132**, 3019–3032.

608 Tebaldi, C., L. O. Mearns, D. Nychka, and R. L. Smith, 2004: Regional probabilities of
609 precipitation change: A Bayesian analysis of multimodel simulations. *Geophys.*
610 *Res. Lett.*, **31**, L24213.

611 Vrugt, J. A., M. P. Clark, C. G. H. Diks, Q. Duan, and B. A. Robinson, 2006: Multi-
612 objective calibration of forecast ensembles using Bayesian model averaging,
613 *Geophys. Res. Lett.*, **33**, L19817.

614 Wolff, J. K., M. Harrold, T. Fowler, J. H. Gotway, L. Nance, and B. G. Brown, 2014:
615 Beyond the basics: Evaluating modelbased precipitation forecasts using traditional,
616 spatial, and object-based methods. *Wea Forecasting*, **29**, 1451–1472.

617 Xie, P. P., and A. Y. Xiong, 2011: A conceptual model for constructing high-resolution
618 gauge-satellite merged precipitation analyses. *J. Geophys. Res.*, **116**, D21106.

619 Yuan, H., X. Gao, S. L. Mullen, S. Sorooshian, J. Du, and H. H. Juang, 2007:
620 Calibration of probabilistic quantitative precipitation forecasts with an artificial
621 neural network, *Wea. Forecasting*, **22**, 1287-1303.

622 Yun, W.T., L. Stefanova, A. K. Mitra, T. S. V. Vijaya Kumar, W. Dewar and T. N.
623 Krishnamurti, 2005: A multi-model superensemble algorithm for seasonal climate
624 prediction using DEMETER forecasts, *Tellus*, **57A**, 280–289.

625 Zhang, H. B., X. F. Zhi, J. Chen, and Coauthors, 2015: Study of the modification of
626 multi-model ensemble schemes for tropical cyclone forecasts. *J. Trop. Meteor.*, **21**,
627 389–399.

628 Zhang, L., and X. F. Zhi, 2015: Multi-model consensus forecasting of low temperature
629 and icy weather over central and southern China in early 2008. *J. Trop. Meteor.*,
630 **21**, 67–75.

631 Zhi, X. F., H. X. Qi, Y. Q. Bai, and C. Z. Lin, 2012: A comparison of three kinds of
632 multi-model ensemble forecast techniques based on the TIGGE data. *J. Meteor.*
633 *Res.*, **26**, 41–51.

634 Captions

635 Table 1. Ensemble forecast systems used in this study.

636 Table 2. Weights and confidence values for pair attributes of matched objects used in
637 this study. CD and CDI denotes the centroid distance and centroid distance interest,
638 respectively. AR is the area ratio ($AR = \min(Area_o, Area_f) / \max(Area_o, Area_f)$) and K is the aspect ratio. This table is adopted from Johnson and Wang
639 (2013).
640

641 Fig. 1 Observed objects for 24-h accumulated precipitation on 2 Jun 2013. (a) original
642 precipitation field before smoothing; (b) convoluted precipitation field after
643 smoothing with a 4-gridpoint averaging radius; (c) filtered precipitation field with
644 the precipitation intensity greater or equal to 10mm.

645 Fig. 2 Interest function G_i for (a) area ratio; (b) centroid distance; (c) aspect ratio
646 difference and (d) angle difference. This figure is adopted from Johnson and Wang
647 (2013).

648 Fig. 3. The total number of objects and the average objects area for the observation
649 (solid line) and the ECMWF EPS (dashed line) for different convolution radii
650 (colors) and precipitation thresholds (abscissa).

651 Fig. 4. Number of average precipitation objects vs. average object area for different
652 precipitation thresholds T (line color and type) for the observations (left) and the
653 forecasts of the 51 members of the ECMWF EPS (right) and for increasing
654 convolution radii R (top-to-bottom).

655 Fig. 5. Distribution of objects with specific attribute values as a fraction of the total
656 number of objects for convolution radius $R=4$ grid points and precipitation
657 threshold $T=10\text{mm}$ for observations (black bar) and 24-h lead time predictions
658 from all members of all EPS (white bar). (a) object area, (b) object aspect ratio,
659 (c) object orientation angle, (d) zonal grid point of object centroid, (e) meridional
660 grid point of object centroid.

661 Fig. 6. The objects mean zonal centroid location of the individual members of the five
662 EPS compared to the observation. (a) ECMWF, (b) NCEP, (c) UKMO, (d) JMA,
663 and (e) CMA.

664 Fig. 7. The same as Fig. 6, but for meridional centroid location.

665 Fig. 8. The same as Fig. 6, but for aspect ratio.

666 Fig. 9. The same as Fig. 6, but for orientation angle.

667 Fig. 10. (a) MMI and (b) OTS for five individual EPSs and the three multi-model
668 forecasting with $R=4$ grid points and $T=10\text{mm}$ for the lead time of 1-7 days.

669 Fig. 11. Weights of the five EPSs with lead times of 1-7 days calculated by MMI (right)
670 and OTS (left), respectively.

671 Fig. 12. The average difference between the forecasted (individual EPSs and multi-
672 model ensemble forecasts) and observed object attribute distributions with $R=4$
673 grid points and $T=10\text{ mm}$ as a function of lead time for (a) aspect ratio, (b)
674 orientation angle, (c) zonal grid point of centroid and (d) meridional grid point of
675 centroid.

676 Fig. 13. Fractions skill scores against forecast lead days for spatial scales s of 1 grid
677 point (dots), 2 grid points (asterisks), and 3 grid points (triangles) for the multi-
678 model ensemble predictions based on MAE (SUP), equally-weighted multi-model
679 ensemble mean (EMME), and multi-model ensemble predictions based on object-
680 based scores (MME_{OTS}).

681 **Tables**682 **Table 1.** Ensemble forecast systems used in this study.

Prediction Center	Model spectral resolution	Initial perturbation scheme	Representation of model error and uncertainty	Ensemble members	Max forecast Lead time (day)
ECMWF	T399162/T255L62	Singular vectors and EDA	SKEB/SPPT	51	15
NCEP	T126L28	Ensemble transform and rescaling	STTP	21	15
UKMO	90km	ETKF	SKEB	24	15
JMA	T106	Singular vectors	SPPT	51	11
CMA	T106	Bred vectors	None	15	10

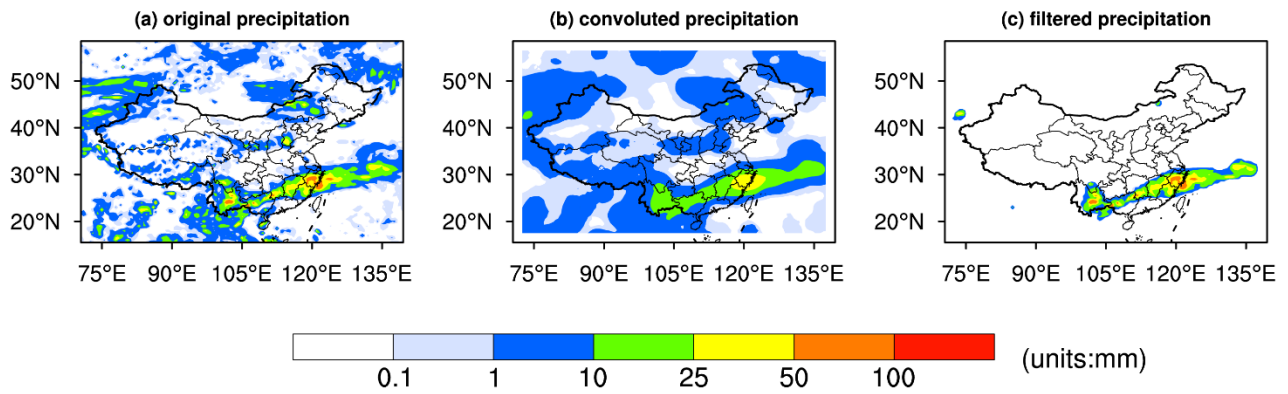
683

684 **Table 2.** Weights and confidence values for pair attributes of matched objects used in
 685 this study. CD and CDI denotes the centroid distance and centroid distance
 686 interest, respectively. AR is the area ratio (AR= min(area(obs),
 687 area(mod))/max(area(obs), area(mod))) and K is the aspect ratio. This table is
 688 adopted from Johnson and Wang (2013).

Pair attributes of matched objects	Weight	Confidence
Centroid distance (CD)	2.0	AR 1.0 if $CD \leq 160\text{km}$
Area ratio (AR)	2.0	$1 - [(CD - 160)/640]$ if $160 < CD < 800\text{km}$ 0.0 if $CD \geq 800\text{km}$
Aspect ratio difference	1.0	$CDI \times AR$
Orientation angle difference	1.0	$CDI \times AR \times \sqrt{a^2 + b^2}$ $a, b = [(K - 1)^2 / (K^2 - 1)]^{0.3}$ for the two matched objects

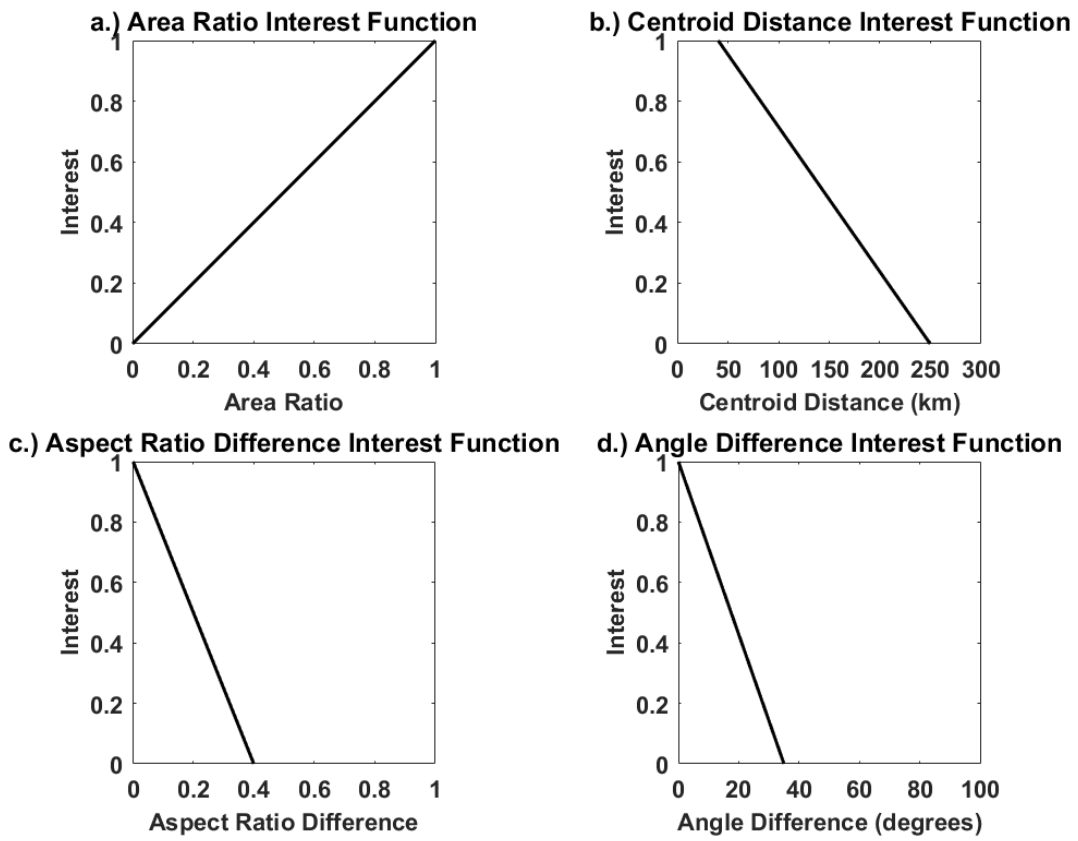
689

690 Figures



691

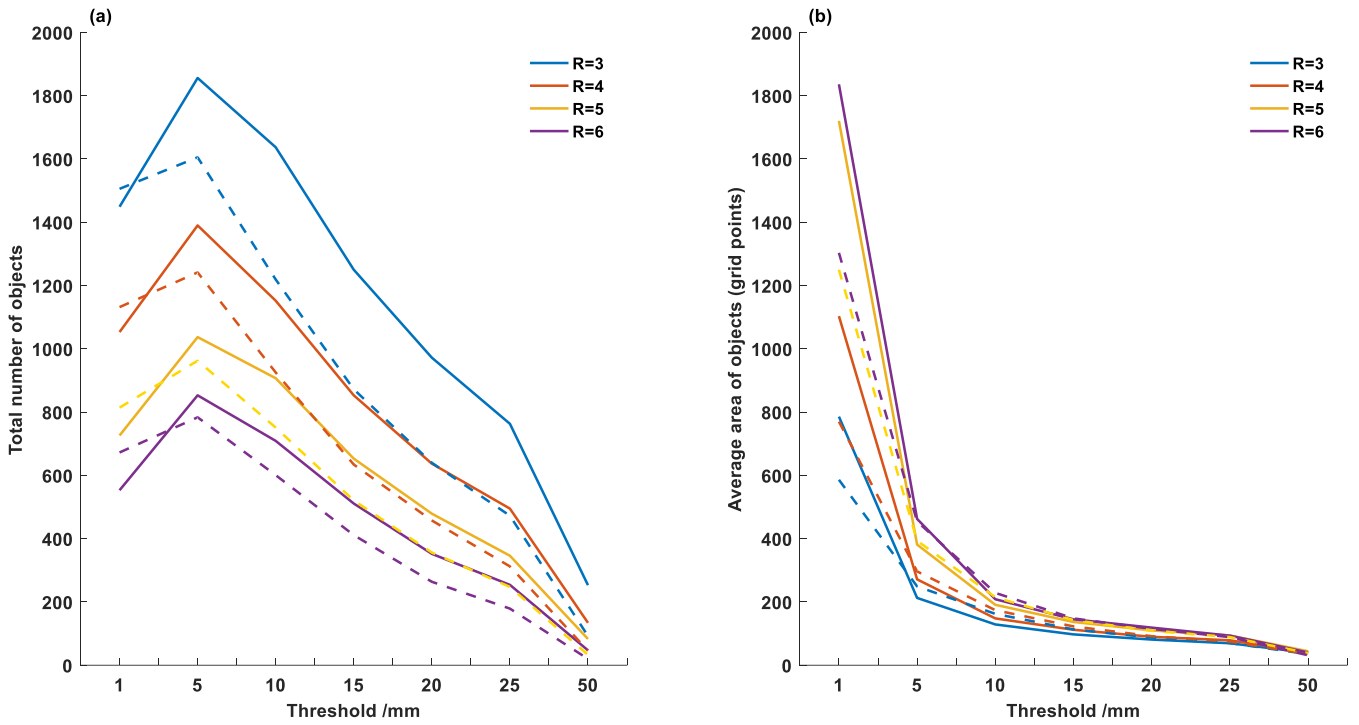
692 **Fig. 1.** Observed objects for 24-h accumulated precipitation on 2 Jun 2013. (a) original precipitation
693 field before smoothing; (b) convoluted precipitation field after smoothing with a 4-gridpoint
694 averaging radius; (c) filtered precipitation field with the precipitation intensity greater or equal to
695 10mm.



696

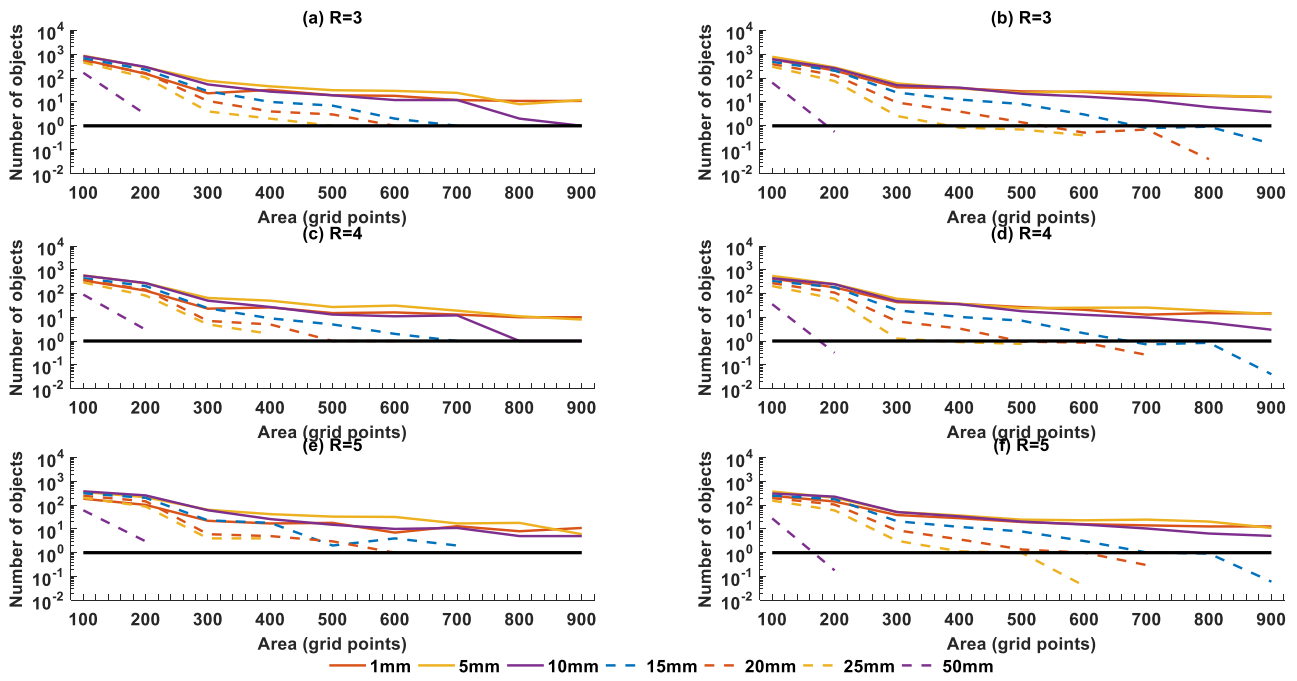
697 **Fig. 2.** Interest function G_i for (a) area ratio; (b) centroid distance; (c) aspect ratio difference and (d)

698 angle difference. This figure is adopted from Johnson and Wang (2013).



699

700 **Fig. 3.** The total number of objects and the average objects area for the observation (solid line) and the
 701 ECMWF EPS (dashed line) for different convolution radii (colors) and precipitation thresholds
 702 (abscissa).



703

704

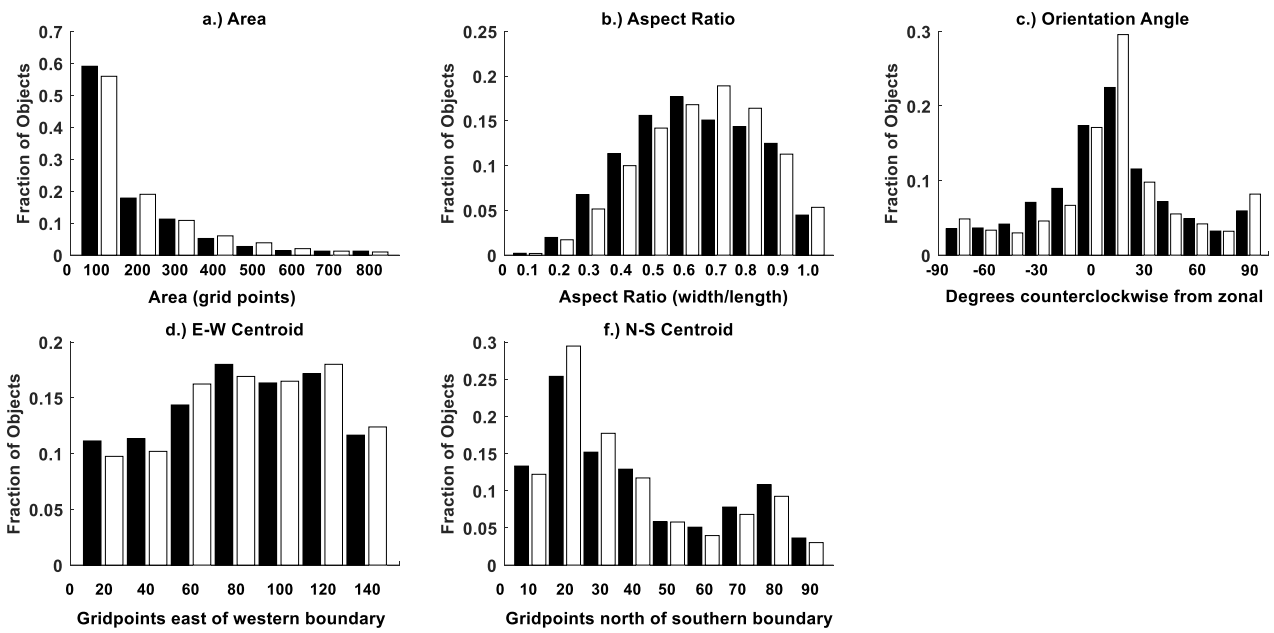
Fig. 4. Number of average precipitation objects vs. average object area for different precipitation

705

thresholds T (line color and type) for the observations (left) and the forecasts of the 51 members

706

of the ECMWF EPS (right) and for increasing convolution radii R (top-to-bottom).



707

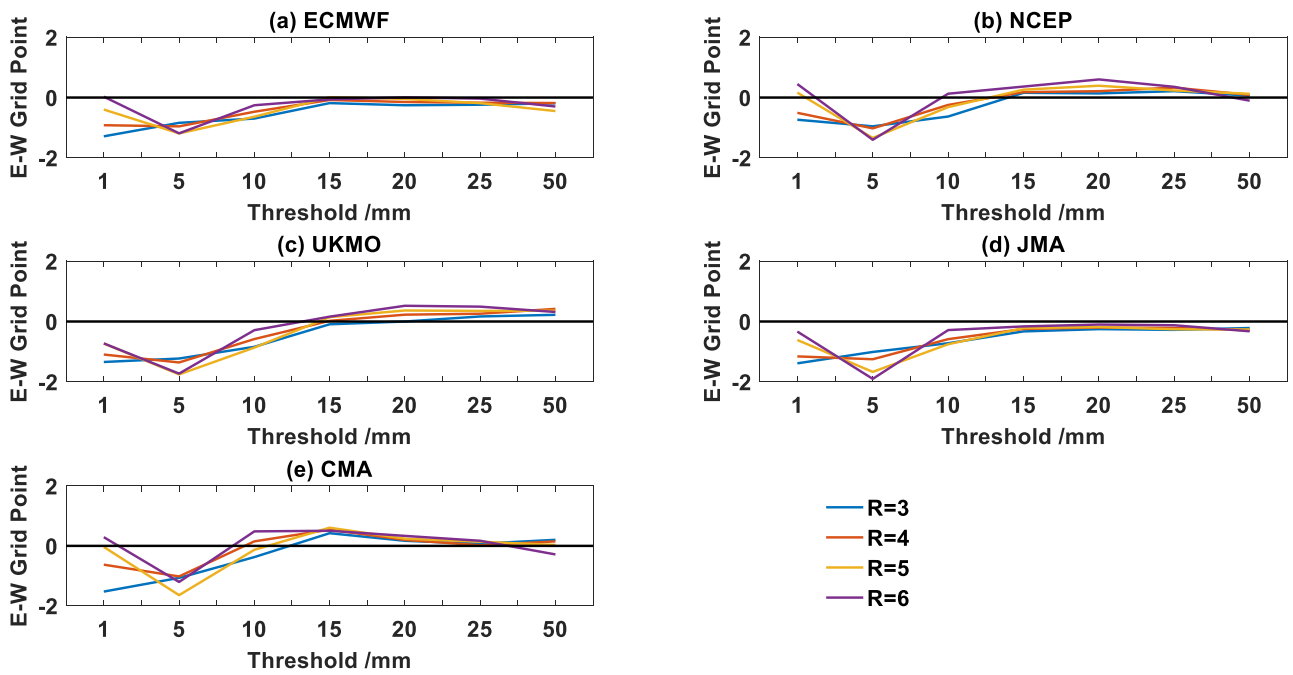
708 **Fig. 5.** Distribution of objects with specific attribute values as a fraction of the total number of objects

709 for convolution radius $R=4$ grid points and precipitation threshold $T=10\text{mm}$ for observations

710 (black bar) and 24-h lead time predictions from all members of all EPS (white bar). (a) object

711 area, (b) object aspect ratio, (c) object orientation angle, (d) zonal grid point of object centroid,

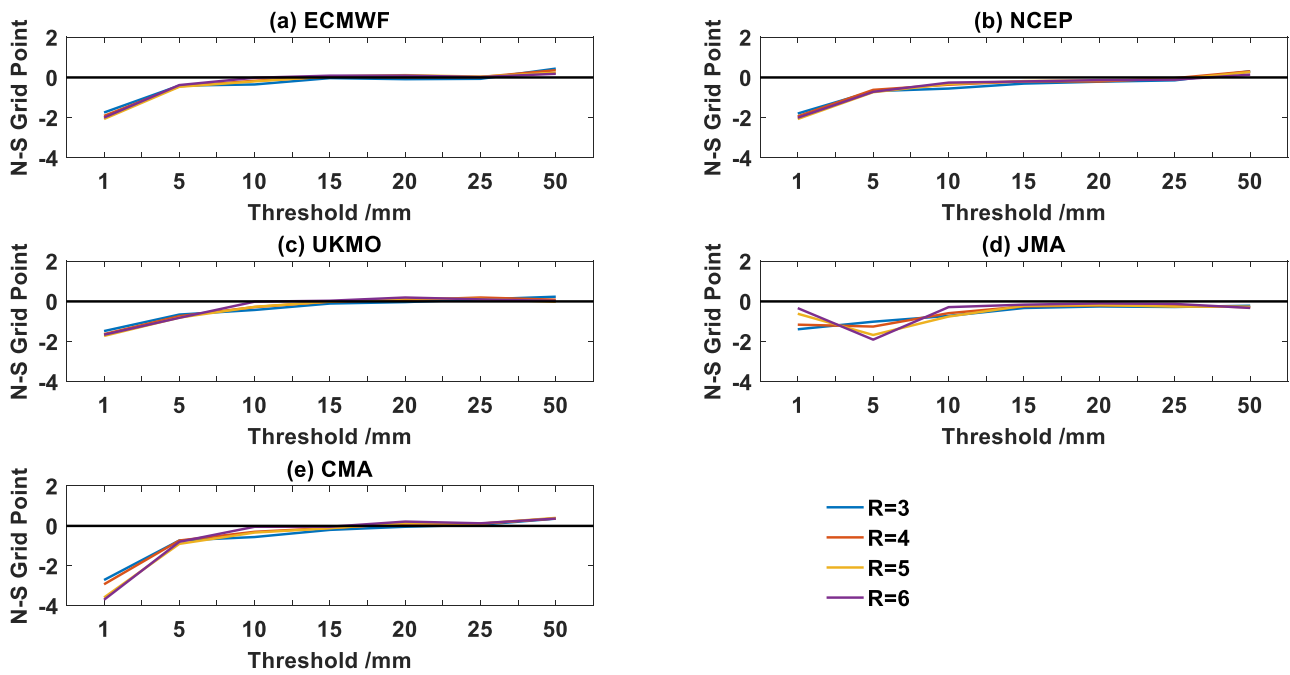
712 (e) meridional grid point of object centroid.



713

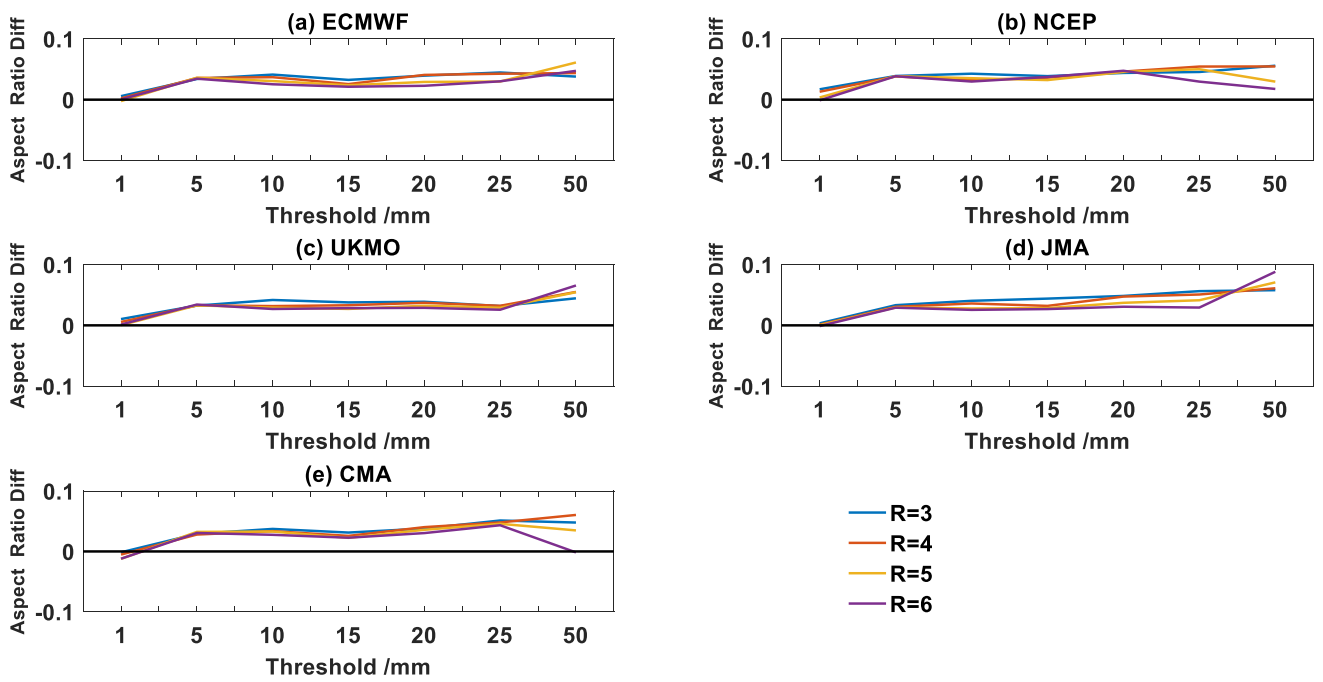
714 **Fig. 6.** The objects mean zonal centroid location of the individual members of the five EPS compared

715 to the observation. (a) ECMWF, (b) NCEP, (c) UKMO, (d) JMA, and (e) CMA.



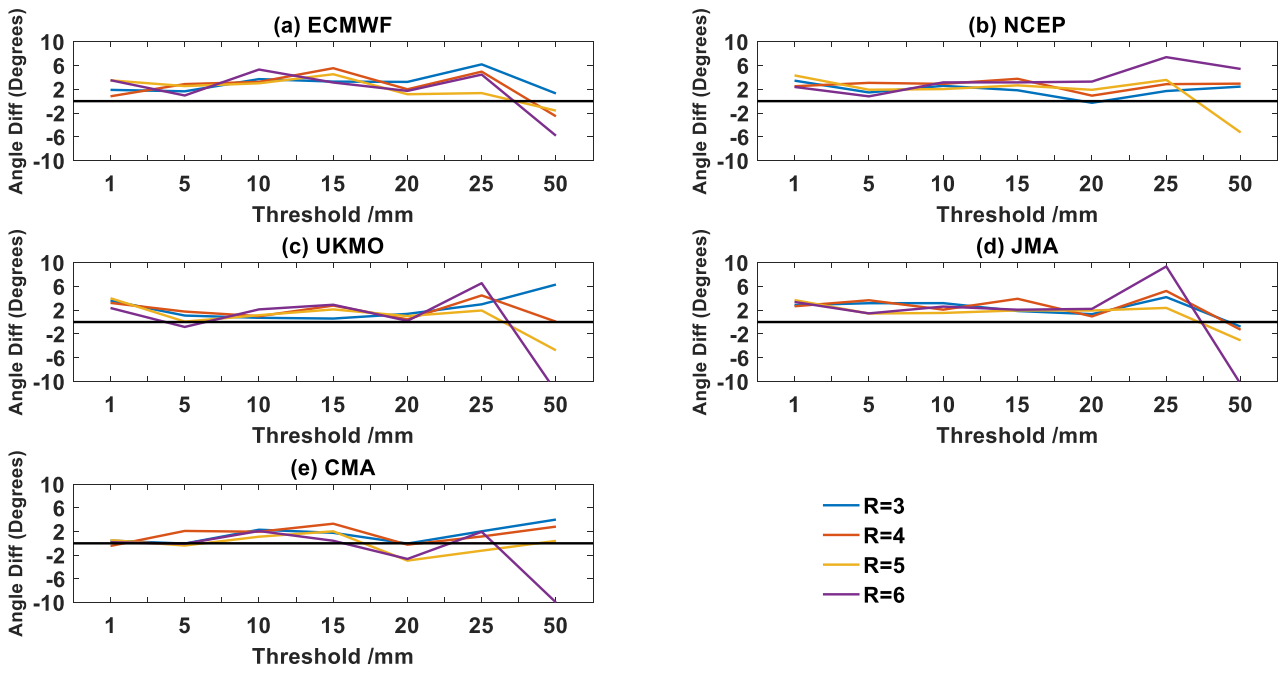
716

717 **Fig. 7.** The same as Fig. 5, but for meridional centroid location.



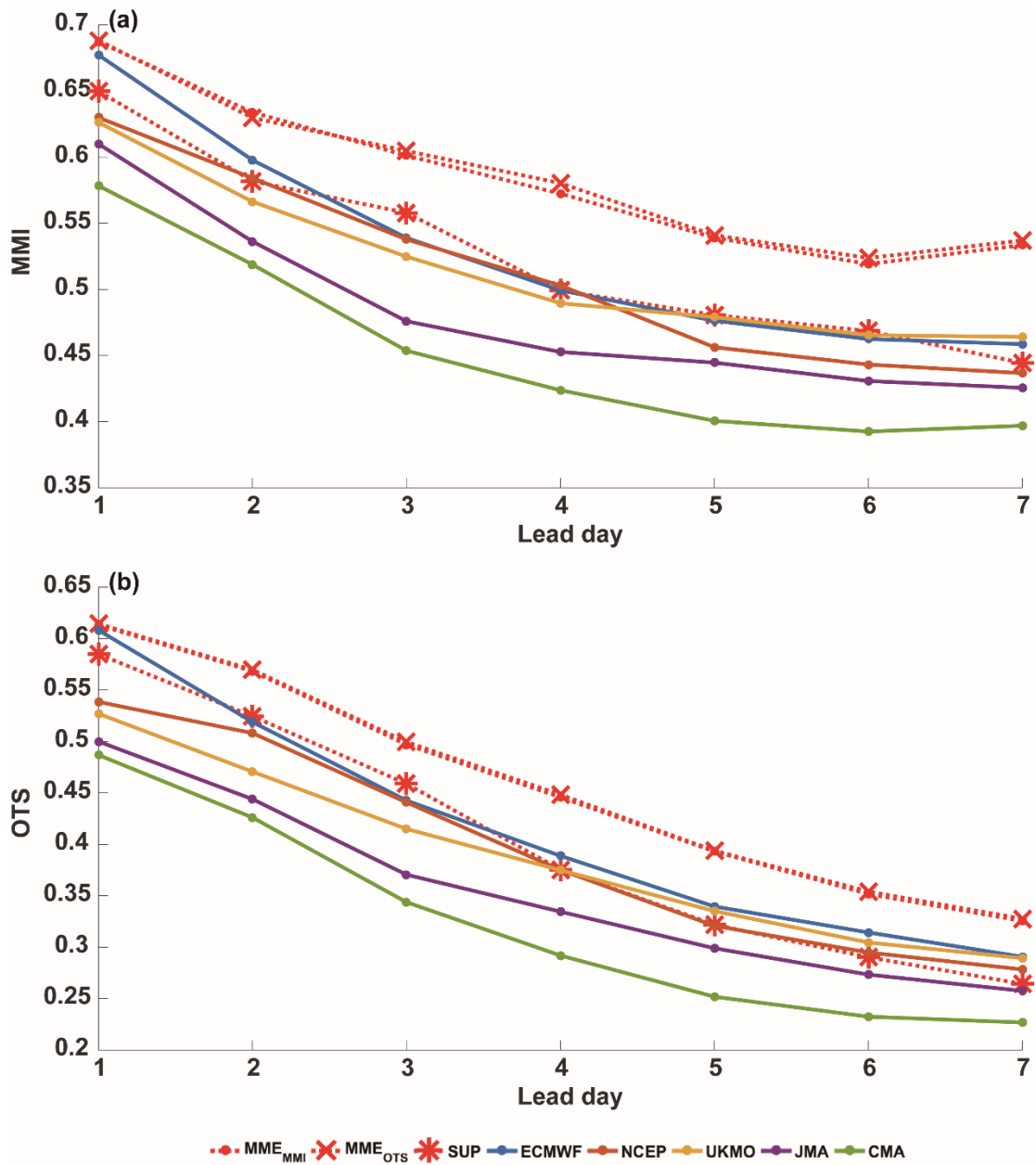
718

719 **Fig. 8.** The same as Fig. 5, but for aspect ratio.



720

721 **Fig. 9.** The same as Fig. 5, but for orientation angle.

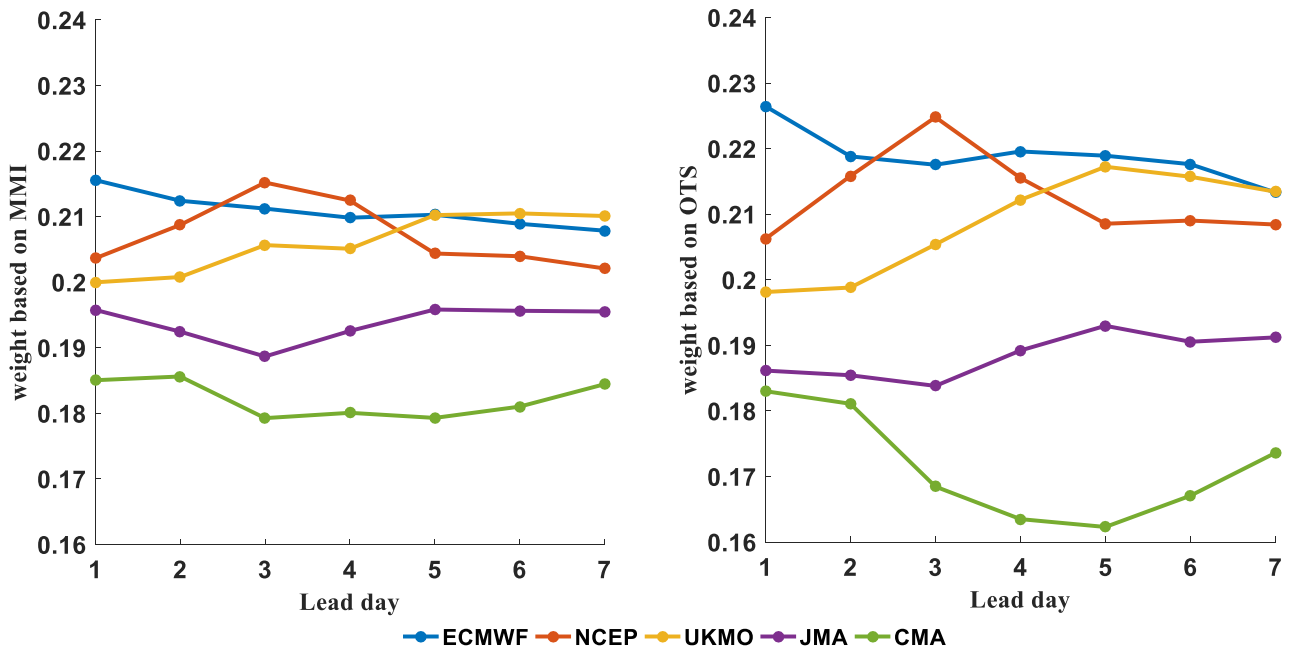


722

723 **Fig. 10.** (a) MMI and (b) OTS for five individual EPSs and the three multi-model forecasting with $R=4$

724

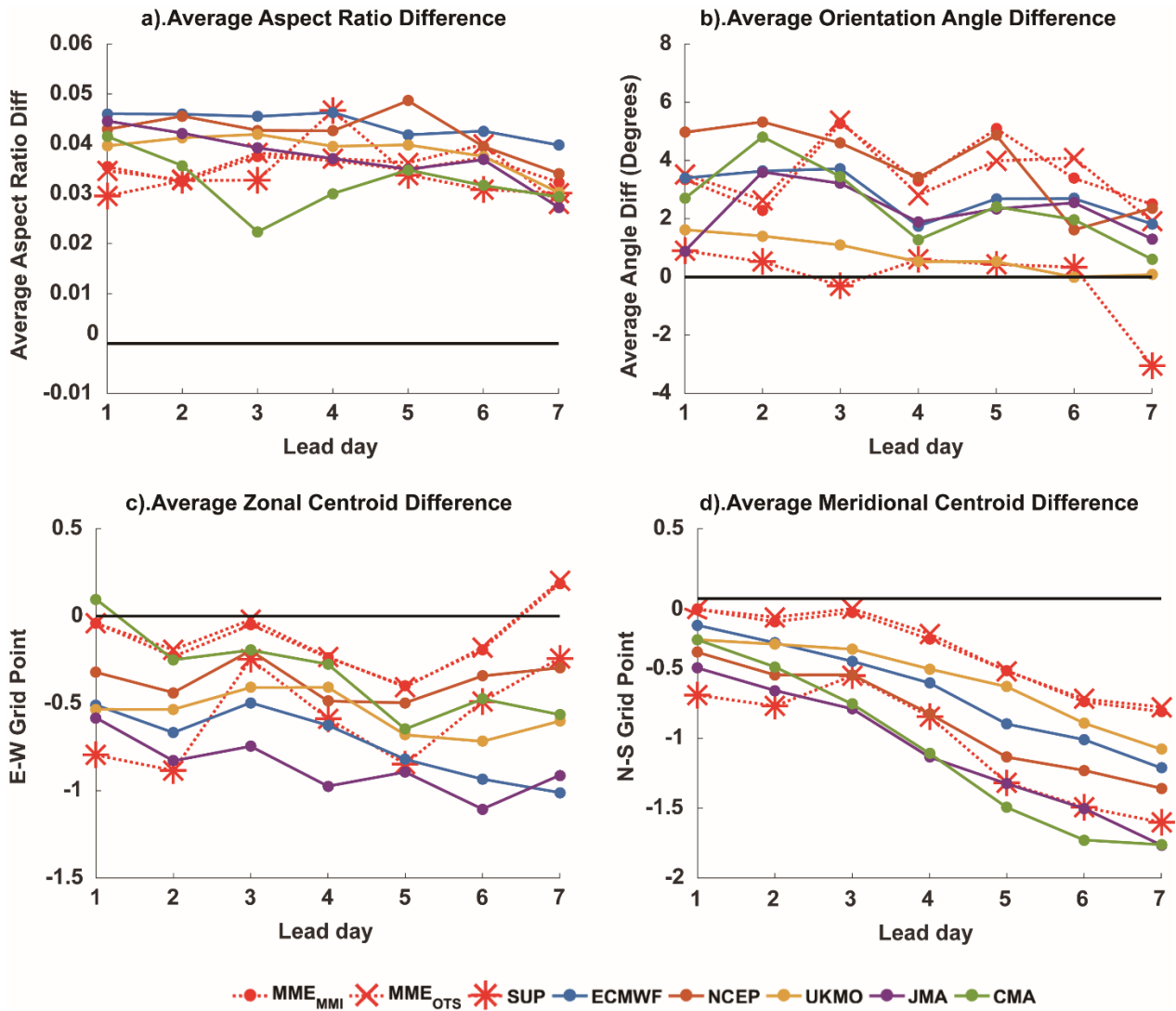
grid points and $T=10\text{mm}$ for the lead time of 1-7 days.



725

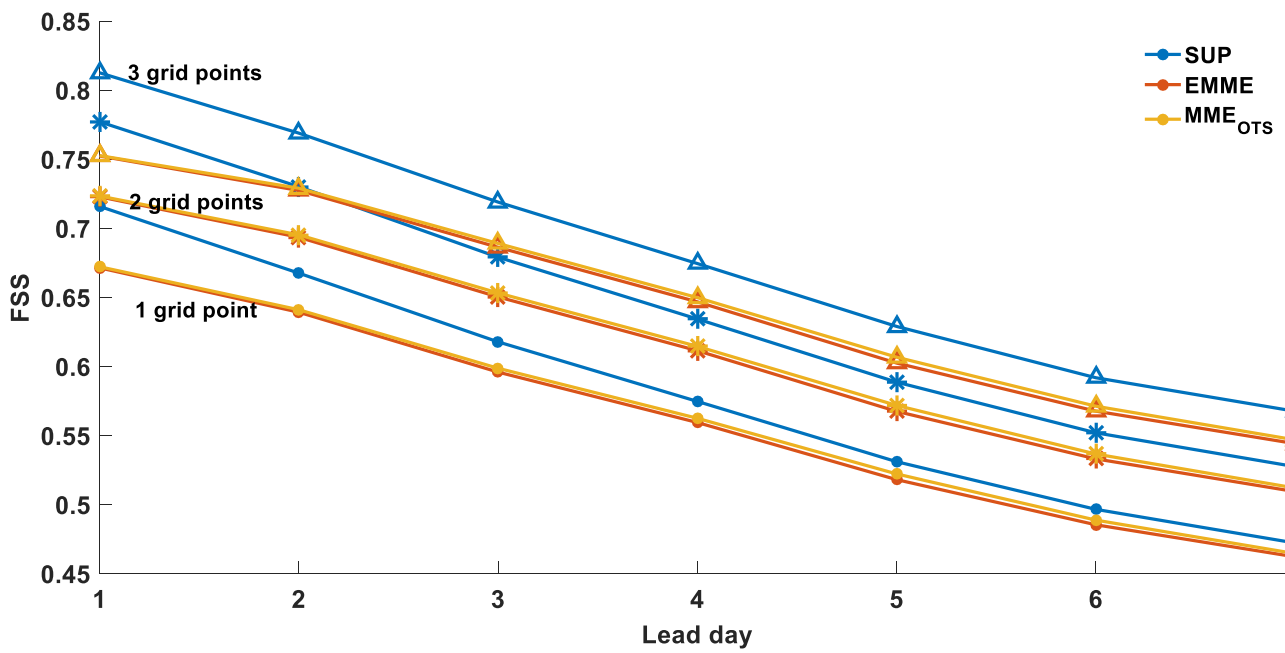
726 **Fig. 11.** Weights of the five EPSs with lead times of 1-7 days calculated by MMI (right) and OTS

727 (left), respectively.



728

729 **Fig. 12.** The average difference between the forecasted (individual EPSs and multi-model ensemble
 730 forecasts) and observed object attribute distributions with $R=4$ grid points and $T=10$ mm as a
 731 function of lead time for (a) aspect ratio, (b) orientation angle, (c) zonal grid point of centroid and
 732 (d) meridional grid point of centroid.



733

734 **Fig. 13.** Fractions skill scores against forecast lead days for spatial scales s of 1 grid point (dots), 2

735 grid points (asterisks), and 3 grid points (triangles) for the multi-model ensemble predictions

736 based on MAE (SUP), equally-weighted multi-model ensemble mean (EMME), and multi-model

737 ensemble predictions based on object-based scores (MME_{OTS}).