

ABSTRACT

7
8 With movement towards kilometer scale ensembles, new techniques are needed for their
9 characterization. We present a new methodology for detailed spatial ensemble character-
10 ization using the Fractions Skill Score (FSS). To evaluate spatial forecast differences the
11 average and standard deviation are taken of the FSS calculated over all ensemble member-
12 member pairs at different scales and lead times. These methods were found to give impor-
13 tant information about the ensemble behavior allowing the identification of useful spatial
14 scales, spin-up times for the model, and upscale growth of errors and forecast differences.
15 The ensemble spread was found to be highly dependent on the spatial scales considered
16 and the threshold applied to the field. High thresholds picked out localized and intense
17 values that gave large temporal variability in ensemble spread: local processes and under
18 sampling dominate for these thresholds. For lower thresholds the ensemble spread increases
19 with time as differences between the ensemble members upscale. Two convective cases were
20 investigated based on the UK Met Office United Model run at 2.2 km resolution. Different
21 ensemble types were considered: ensembles produced using the Met Office Global and Re-
22 gional Ensemble Prediction System (MOGREPS) and an ensemble produced using different
23 model physics configurations. Comparison of the MOGREPS and multiphysics ensembles
24 demonstrated the utility of spatial ensemble evaluation techniques for assessing the impact
25 of different perturbation strategies and the need for assessing spread at different, believable,
26 spatial scales.

1. Introduction

It has been long known that at small spatial scales forecast errors grow more rapidly (Lorenz (1969); Ehrendorfer (1997); Palmer (2000) and references therein) possibly resulting in rapid upscale error growth in high resolution models. In recent years these subjects have again come under discussion as increases in computer power allow models to be run at higher and higher resolutions (Mass et al. (2002) and references therein, Lean et al. (2008)). Hohenegger and Schär (2007a) compared the predictability at large (around 80 km) and convection-permitting (2.2 km) scales and found error doubling times around ten times shorter for the higher resolution simulations. Further work has investigated the links between mesoscale processes and error growth with a focus on moist dynamics (Zhang 2005; Hohenegger et al. 2006) and the separation of equilibrium and triggered convection to distinguish different modes of predictability in convective events (Keil and Craig 2011; Zimmer et al. 2011; Craig et al. 2012; Keil et al. 2013).

Ensemble prediction systems strive to represent the meteorological uncertainty present in a particular forecast and have been widely used to assess error growth in a variety of high-resolution situations (Walser et al. 2004; Walser and Schär 2004; Hohenegger and Schär 2007b; Hanley et al. 2011, 2013). Further investigations have been conducted into different ensemble perturbation strategies for high resolution ensembles including initial condition perturbations (Migliorini et al. 2011; Caron 2013; Kühnlein et al. 2013), physics perturbations (Stensrud et al. 2000; Hacker et al. 2011; Gebhardt et al. 2011; Vié et al. 2012; Baker et al. 2014), perturbation of boundary layer parameters (Martin and Xue 2006; Leoncini et al. 2010; Done et al. 2012) and the use of different physics schemes (Berner et al. 2011; Leoncini et al. 2012).

The aim of this paper is to provide a new methodology for evaluating, thoroughly, the differences between members of a convection permitting ensemble and the dependence of these differences on spatial scale. These methods are based on the Fractions Skill Score (FSS, Roberts and Lean (2008); Roberts (2008)). Various considerations are discussed in-

54 cluding the forecast evolution through different lead times, the effect of considering different
55 threshold values for the fields used to calculate the FSS, and the comparison of different
56 forecast variables. For the demonstrative purposes of this paper two convective cases are
57 considered using ensembles produced as part of the Met Office Global and Regional En-
58 semble Prediction System (MOGREPS, Bowler et al. (2008, 2009)). The spatial spread of
59 the ensemble members is characterized and the realism of the ensemble spread is tested by
60 comparing with the skill against radar derived precipitation accumulations. Radar data is
61 necessary as a verification source because of its high spatial coverage.

62 The technique used to determine spatial differences between members can also be used
63 for the comparison of different model formulations within the ensemble. To demonstrate this,
64 different model physics configurations were considered in addition to the MOGREPS ensem-
65 ble members for the second case study. This specific example is provided to demonstrate the
66 utility of spatial evaluation techniques in the comparison of different ensemble formulations.
67 Note, however, that a complete systematic evaluation comparing different types of physics
68 configuration is outside the scope of this paper. To do this it would be necessary to consider
69 a large number of cases with different convective forcing as detailed by, for example, Stensrud
70 et al. (2000); Keil et al. (2013). The spatial ensemble spread produced by different physics
71 configurations strategies is evaluated and compared to that of the MOGREPS ensemble.
72 In operational frameworks, different physics configurations are often considered in addition
73 to initial and boundary condition perturbations and so the spatial spread produced by an
74 ensemble with different MOGREPS members combined with different physics configurations
75 is also investigated.

76 To evaluate convection permitting ensembles in a sensible way it is necessary to choose
77 a verification approach that considers multiple spatial scales and does not suffer from the
78 double penalty problem where spatial errors are penalized twice: once for being a near miss,
79 and again for being a false positive. Many possible spatial verification approaches have been
80 proposed in recent years; for an overview the reader is referred to the review papers of Ebert

81 (2008), Gilleland et al. (2009), and Johnson and Wang (2013). The spatial approach has
82 also been applied to ensembles (Clark et al. 2011; Johnson et al. 2014; Surcel et al. 2014).
83 Here we have chosen to focus on the Fractions Skill Score (FSS) of Roberts and Lean (2008);
84 Roberts (2008). The FSS is a fuzzy verification measure used to compare two fields within
85 a given square neighborhood.

86 Since its original formulation the FSS has been used for different applications and several
87 further developments have been proposed. Schwartz et al. (2010) consider circular neigh-
88 borhoods to calculate the field of fractions at each grid point and then produce probabilistic
89 guidance using the field of fractions as a neighborhood probability. Duda and Gallus (2013)
90 also use the circular neighborhood approach, verifying the precipitation of mesoscale con-
91 vective systems. In this paper the FSS is considered over a square neighborhood as detailed
92 in Roberts and Lean (2008); Roberts (2008). Duc et al. (2013) extend the FSS to include
93 temporal and ensemble dimensions to give a single FSS value representative of the ensem-
94 ble. A single field of fractions including spatial, temporal and ensemble information is then
95 compared with observations. This is useful for providing an overview of model performance
96 but does not provide information regarding the spread - skill relationship of the ensemble or
97 the spatial differences between individual pairs of ensemble members.

98 Rezacova et al. (2009) use the FSS to calculate the ensemble spread-skill relationship with
99 the ensemble skill calculated from the FSS between ensemble member- radar comparisons and
100 the ensemble spread from the FSS between perturbed ensemble members and the ensemble
101 control. Following on from this Zacharov and Rezacova (2009) determine a relationship
102 between the FSS estimates of ensemble spread and skill and use this to predict the ensemble
103 skill given the spread. Zacharov and Rezacova (2009) consider together FSS results from
104 differently sized neighborhoods. This method was chosen because there is no fixed scale that
105 can give a FSS skill value over different cases. However, as different physical behavior is
106 apparent at different spatial scales (e.g as shown in Roberts (2008)) it is informative also
107 to investigate how the ensemble spread varies with spatial scale which is the subject of

108 this paper. Whereas Rezacova et al. (2009); Zacharov and Rezacova (2009) only consider
109 comparisons between perturbed ensemble members and the control, in this paper the FSS
110 between all independent member-member pairs is considered. Considering all members in
111 this manner is the best representation of total spread as it includes fully the inter-member
112 variability and does not rely on the ensemble mean which is known to lie outside the model
113 manifold (Ansell 2013). Further work by the authors (in preparation, to be submitted shortly
114 to Monthly Weather Review) considers other possible methods of member comparison.

115 Here we present the following: in Section 2 we introduce the two case studies that will
116 provide examples throughout the paper. The model configuration is also discussed along with
117 a justification for our method of using the FSS. Section 3 provides examples of our results
118 for ensembles with different IC and LBC perturbations and results for different physics
119 configurations are discussed in Section 4. Finally, in Section 5 we summarize the conclusions
120 from this work and discuss areas of further investigation.

121 **2. Method**

122 *a. Cases*

123 Two convective cases were chosen for the demonstrative purposes of this paper. In these
124 cases convection occurs in different synoptic situations. The first case, 23 April 2011, was
125 chosen as an example of organized spring convection over England and will be referred to as
126 the ‘organized Spring’ case. This case has a low pressure system centered to the northwest of
127 the UK and a high pressure system centered over Scandinavia. A frontal structure stretches
128 down across the western UK. As the front moves eastward a convergence line forms across
129 eastern England ahead of the front. This convergence line is shown in the UK Met Office
130 analysis at 1800 UTC on the 23 April (Figure 1a). Convective storms developed in the
131 vicinity of this convergence line with precipitation first seen at 1400 UTC on 23 April, and
132 continuing until 0300 UTC on 24 April. At 1800 UTC a band of frontal precipitation enters

133 the model domain from the NW preceding an occluded front which enters the domain at
134 0000 UTC on the 24 April.

135 The second case, 8 July 2011, features a number of convective storms that formed over
136 the UK in an area of instability within the circulation of a decaying low pressure system.
137 At 0600 UTC the low center was situated over Ireland as shown in Figure 1b. Throughout
138 the day the low center then moved towards the northeast reaching the northeast of England
139 by 1800 UTC. By 1400 UTC there were many heavy showers over Scotland as indicated
140 by the Nimrod radar system (not shown). Convective clouds associated with these showers
141 were also seen from visible satellite observations from the Meteosat geostationary satellite.
142 For this case study we focus on one particular storm that formed over the Edinburgh area
143 of eastern Scotland and remained stationary for around four hours producing large rainfall
144 totals (0900 UTC – 2100 UTC radar derived precipitation totals of over 64 mm) and flooding.
145 In future discussion this will be referred to as the ‘flooding’ case. Previous analysis of this
146 case by Leoncini et al. (2011) showed that the Met Office 2.2 km ensemble on this occasion
147 gave a 30% to 40% chance of a flood-producing storm within 25 km of Edinburgh; a level of
148 significant risk.

149 *b. Model Setup*

150 The UK Met Office Unified Model (MetUM) runs with a non-hydrostatic dynamical core
151 with semi-Lagrangian advection (Davies et al. 2005). A comprehensive set of parametriza-
152 tions are used including: surface exchange (Essery et al. 2001), boundary layer mixing (Lock
153 et al. 2000), radiation (Edwards and Slingo 1996) and mixed phase cloud microphysics
154 based on Wilson and Ballard (1999). Version 7.7 of the global ensemble prediction system
155 (MOGREPS-G) was run at a resolution of around 60km in the mid-latitude regions with
156 70 vertical levels. MOGREPS-G provided the initial conditions (ICs) and lateral bound-
157 ary conditions (LBCs) for the North Atlantic and European (NAE) regional model run at
158 18km resolution with 70 vertical levels. Perturbations were generated using an ensemble

159 transform Kalman filter and then added to the Met Office 4D-Var analysis as described by
160 Bowler et al. (2008, 2009). This perturbation strategy includes a stochastic kinetic energy
161 backscatter scheme and localization. Model error is addressed using the “random param-
162 eters” scheme for both ensembles to account for sub-grid processes uncertainty. Both the
163 global and regional ensembles have 23 perturbed members and an unperturbed control.

164 For the case studies described here a high resolution ensemble, run over the Met Office
165 variable resolution UK domain, was one way nested inside the NAE model. This domain has
166 a constant resolution 2.2 km grid over the UK with the grid stretched up to 4 km around the
167 domain edges to reduce the jump in resolution when downscaling from the NAE model. No
168 further data assimilation was included when downscaling from the NAE to UK domain. The
169 global and NAE models were run with a convection scheme based on Gregory and Rowntree
170 (1990) but modified since (Derbyshire et al. 2011). The 2.2 km model has explicit convection
171 only (no convection scheme). The 2.2 km UK domain is shown in Figure 2 in light gray and
172 is approximately 920 km W-E by 1200 km N-S.

173 For the flooding case eleven perturbed members plus a control were run over the 2.2
174 m domain using LBCs and ICs taken from the first eleven members, and control, of the
175 NAE regional ensemble (MOGREPS-R). Twelve simulations were run because this was the
176 ensemble size being considered for an operational 2.2 m ensemble system (MOGREPS-UK,
177 operational since 2013 (Mylne 2013)). To allow the flood producing storm over Edinburgh
178 to be investigated, analysis for this case was also completed over a small 100 km domain
179 surrounding this region. This subdomain is highlighted in Figure 2 in dark gray.

180 For the organized Spring case an ensemble of 8 MOGREPS simulations were run (seven
181 perturbed members plus a control). This reduction in size allowed 5 different physics config-
182 urations to be considered for each MOGREPS simulation (giving a total of 40 simulations).
183 The different model configurations were:

- 184 i. A control ensemble with the standard model settings labeled “standard”.
- 185 ii. An ensemble with a restricted version of the convection scheme (Roberts 2003) as

- 186 would be applied to the Met Office 4 km deterministic model (labeled “conv”).
- 187 iii. An ensemble with the time step increased from 25 s to 50 s labeled “time”. It is
188 interesting to investigate the effects of a longer time step as increasing the time step
189 reduces the computational cost of the simulation but may increase model error.
- 190 iv. An ensemble with increased time step and restricted convection scheme labeled “conv+time”.
- 191 v. An ensemble with modifications to the graupel labeled “grp”. The graupel modification
192 allows the production of graupel through the capture of rain by snow and results in an
193 increased graupel mass. This modification has become a standard option in Met UM
194 versions 8.0 onwards (Wilkinson 2011).

195 It must be emphasized that these model configurations were chosen to demonstrate the
196 methodology presented in this paper, not as possible implementations to the UK Met Office
197 ensemble prediction system. Note also that the model variations are neither stochastic nor
198 designed to represent the model error, although they do, nevertheless, represent plausible
199 alternative formulations. The UK model for the organized Spring case was started at 0600
200 UTC on 23 April 2011, the flooding case at 1800 UTC on 7 July 2011. MOGREPS-G and
201 MOGREPS-R were initiated 6 hrs and 3 hrs respectively before the UK model. For both
202 cases the UK model was run up to lead times of 36 hours.

203 *c. How the FSS is used*

204 The FSS is described in Roberts and Lean (2008) and summarized here for ease of reading.
205 To calculate the FSS a threshold is first selected, say for precipitation, either as a fixed value
206 (e.g 4 mm hr⁻¹) or as a percentile (e.g top 1% of precipitation field). The field is converted
207 to binary form with grid points set to 1 for values above the threshold and 0 otherwise. A
208 neighborhood size is then selected and, for each neighborhood centered upon each grid point,
209 the fraction of grid points with the value ‘1’ within this square is computed. Two fields of

210 fractions (denoted A and B), say from a model and observations, are then compared using
 211 the mean squared error (MSE). For a neighborhood size n and domain size N_x by N_y grid
 212 points this is given by:

$$213 \quad MSE_{(n)} = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [A_{(n)i,j} - B_{(n)i,j}]^2. \quad (1)$$

214 The fractions skill score is computed by comparing $MSE_{(n)}$ with a reference MSE, $MSE_{(n)ref}$.

$$FSS_{(n)} = 1 - \frac{MSE_{(n)}}{MSE_{(n)ref}} \quad (2)$$

215 where $MSE_{(n)ref}$ is the largest possible MSE that can be obtained from fraction fields A
 216 and B :

$$MSE_{(n)ref} = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} [A_{(n)i,j}^2 + B_{(n)i,j}^2]. \quad (3)$$

217 The FSS varies from 0 (complete mismatch between the fields) to one (perfect match between
 218 the fields).

219 Different neighborhood sizes are considered in order to evaluate the FSS at different
 220 spatial scales. Here we define the neighborhood size to be the total width of the square
 221 neighborhood in km. The smallest possible neighborhood is 2.2 km, set by the grid scale. No
 222 bias exists between the binary fields created using percentile thresholds as, by definition, the
 223 same number of points exceed the threshold for both fields. Hence, for percentile thresholds,
 224 the maximum possible spatial disagreement is found for two fields which place the points of
 225 interest at opposite edges of the domain. A perfect match is only obtained between fields
 226 with this maximum disagreement when they are compared over a neighborhood of twice
 227 the smallest dimension of the domain. In other words, the FSS will only equal 1 when the
 228 neighborhood size is equal to twice the smallest dimension of the domain. This sets the
 229 maximum neighborhood size for percentile thresholds. For value thresholds the fields may
 230 be biased and this argument does not hold. For the examples presented here only percentile
 231 thresholds are considered and the maximum neighborhood size is 1848 km for the UK domain
 232 and 200 km for the 100km subdomain.

233 The FSS can be calculated at a particular time between two different forecasts, or between
234 a forecasts and observations, the former giving a measure of spatial spread, the latter of
235 spatial skill. The ensemble spread is characterized by calculating the FSS for all independent
236 member – member pairs ($N_p(N)$, for an ensemble of N members) resulting in

$$N_p(N) = N \times (N - 1)/2 \quad (4)$$

237 comparisons. Here, and for the remainder of this paper, the control is treated as an ad-
238 ditional ensemble member. Hence, for the flooding case we have 12 MOGREPS members
239 (the 11 perturbed members and unperturbed control) and for the organized Spring case we
240 have 8 MOGREPS members for each physics configuration (the 7 perturbed members and
241 unperturbed control). Justification for this method comes from our interest in the total
242 spatial ensemble spread. In this situation the spatial location of a feature in the control
243 forecast is not necessarily at the center of corresponding features in the perturbed members
244 and therefore we do not wish to assign any special status to the control forecast. Figure
245 3 demonstrates the advantages of our method: when considering the control as an addi-
246 tional ensemble member one can distinguish the different spatial spread in Figures 3a and
247 3b, whereas when only comparing against the control the spread in Figures 3a and 3b is
248 indistinguishable.

249 The ensemble skill is assessed by comparing the model hourly precipitation accumulations
250 with those derived from the UK Met Office Nimrod radar system. The Nimrod system
251 includes calibration against rain gauge data and aims to remove common sources of error
252 (Golding 1998; Harrison et al. 2000). For the summer case 1 km Nimrod radar derived hourly
253 precipitation accumulations are interpolated onto the 2.2 km model grid. Nimrod data at
254 1 km resolution was not available for analysis of the organized Spring case so 5 km data
255 was instead used. The area of Nimrod coverage differs slightly from the UK 2.2 km domain
256 over which the model is run and is indicated by the dotted region in Figure 2. All analysis
257 involving radar data, or the comparison of model and radar data, only considers the area
258 with radar coverage. We assume the radar data is representative of the precipitation that

259 occurred and ignore observational errors, which would have to be considered within a routine
260 verification framework. Visual examination of the radar fields found no obvious errors.

261 To assess ensemble skill each model simulation is separately compared with radar obser-
262 vations, whilst to assess ensemble spread we compare all possible pairings of the model runs.
263 Again consider Figure 3, but this time take the filled black circles to represent the location
264 of precipitation in the radar data. As a measure of ensemble skill we are only considering
265 the spatial differences associated with the solid arrows. These measures of ‘spread’ and
266 ‘skill’ consider different numbers of member-member or member-radar pairs, raising ques-
267 tions about a direct comparison of these metrics. However, answering these questions is not
268 the subject of this paper which focuses on the characterization of spatial ensemble spread,
269 with spatial ensemble skill considered only to put the spread into context. Further work by
270 the authors (in preparation, to be submitted shortly to Monthly Weather Review) focuses
271 in more detail on these metrics in the context of the spread-skill relationship.

272 Three different comparison strategies were used for the organized Spring case to char-
273 acterize the differences between spatial spread in the MOGREPS ensemble and that pro-
274 duced through considering different physics configurations. 8 MOGREPS ensemble members
275 ($N = 8$), and 5 different physics configurations ($N = 5$), were considered. Additionally re-
276 sults were produced using a subset of two physics configurations ($N = 2$) to allow spatial
277 differences resulting from individual configurations to be investigated.

278 i. All independent comparisons were made between the MOGREPS members for a given
279 physics configuration, with each physics configuration treated separately. Considering
280 all 5 physics configurations in this manner gives $N_p(8) \times 5 = 140$ comparisons, a
281 strategy denoted as MOGREPS5. Considering 2 physics configurations in this manner
282 gives $N_p(8) \times 2 = 56$ comparisons, denoted as MOGREPS2.

283 ii. All independent comparisons between the different physics configurations for a given
284 MOGREPS member, with each MOGREPS member treated separately. Considering
285 all 5 physics configurations gives $8 \times N_p(5) = 80$ comparisons for this strategy denoted

286 as Physics5. Considering 2 physics configurations gives $8 \times N_p(2) = 16$ comparisons
287 (Physics2).

288 iii. Comparisons between different MOGREPS members which additionally have differ-
289 ent physics configurations. For example, MOGREPS member 2 with the standard
290 physics configuration might be compared with MOGREPS members 1,3,4,...,12 with
291 the physics configurations conv, conv+time, time and grp. Considering all 5 physics
292 configurations with this comparison strategy, referred to as MOGREPS5+Physics5,
293 gives $N_p(8) \times N_p(5) = 280$ comparisons. Considering 2 physics configurations (MO-
294 GREPS2+Physics2), gives $N_p(8) \times N_p(2) = 28$ comparisons.

295 Given the large number of FSS values FSS_i (one calculated for each comparison) it is
296 necessary to consolidate this information to provide an overview of spatial ensemble behavior.
297 In this paper the mean is taken over the relevant set of FSS_i . When calculated over member-
298 member pairs this is referred to as dFSSmean where ‘d’ indicates that this is a measure
299 of ensemble dispersion. When calculated over member-radar pairs this is referred to as
300 eFSSmean where ‘e’ indicates that this is a measure of ensemble error. dFSSmean gives an
301 indication of the average spatial agreement within the ensemble for a given neighborhood
302 size. In other words, we can select a level of spatial agreement for the ensemble, represented
303 by the value of dFSSmean, and ask at what neighborhood size this agreement is obtained.

304 As the ensemble members do not necessarily have an even spatial distribution, a range of
305 FSS_i will be obtained from the different ensemble member-member pairs. For example, if the
306 majority of ensemble members place rain at the same spatial location but a small number of
307 members place the rain far away this may produce a similar value of dFSSmean as a situation
308 in which all ensemble members place the rain at slightly different spatial locations. Hence it
309 is also important to investigate the range of FSS values surrounding dFSSmean. To do this
310 the standard deviation of FSS values, dFSSstdev, is used. dFSSstdev is closely linked to the
311 standard error in dFSSmean, $\frac{dFSSstdev}{\sqrt{N_{FSS}}}$ where N_{FSS} is the number of FSS_i samples used to
312 calculate dFSSmean. As the purpose of this paper is to focus on the spatial distribution of

313 ensemble members, we consider dFSSstdev and avoid the $\frac{1}{\sqrt{N_{FSS}}}$ dependence on ensemble
314 size. This allows the spatial distribution of differently sized ensembles to be compared.

315 In order to make a spatial comparison between different ensembles it is necessary to find
316 scales which are believable and have a reasonable level of spatial agreement. For the purposes
317 of this paper, ‘believable’ scales for the intercomparison of ensemble members are derived
318 in an equivalent manner to those scales that would be considered skillful if the comparison
319 was instead against observations (assuming that the ensemble is well spread). This scale is
320 quantified using the methodology of Roberts and Lean (2008) where a neighborhood size is
321 considered believable (‘skillful’) if a FSS value of

$$FSS \geq 0.5 + \frac{f_0}{2} \quad (5)$$

322 is obtained for that neighborhood. f_0 is equal to fraction of the field considered in the FSS
323 calculation (for example, considering the top 99th percentile threshold would give $f_0 = 0.01$)
324 and Equation 5 simplifies to an equality when the neighborhood is twice the spatial difference
325 between two binary fields (Roberts and Lean 2008; Roberts 2008). As f_0 is small Equation
326 5 can be approximated as $FSS \geq 0.5$.

327 *d. Thresholding*

328 The FSS can be calculated using either fixed value or percentile thresholds. Following on
329 from the work of Roberts (2008); Mittermaier and Roberts (2010) this paper focuses on the
330 use of percentile thresholds to allow the spatial distribution of phenomena to be investigated.
331 Higher percentile thresholds are associated with smaller, more extreme forecast features, and
332 lower percentile thresholds are associated with larger-scale smoother features (Roberts 2008).
333 Note that here, and in all future discussion, the percentile threshold is applied over the whole
334 domain, including areas both with and without precipitation.

335 To understand the effect of applying percentile thresholds it is informative to investigate
336 the values corresponding to each threshold. Examples for hourly precipitation values corre-

337 sponding to the 90th and 99th percentile thresholds are given in Figure 4. These percentile
338 thresholds are used as examples throughout this paper. All ensemble members (gray solid
339 lines) and radar (black lines) are shown for the organized Spring case (top) and Summer
340 flooding case (bottom). From both cases and thresholds it can be seen that the radar per-
341 centile thresholds generally correspond to lower precipitation values than the model. This
342 bias in the model compared to radar is an important consideration for model evaluation.
343 However, it is also important to investigate the spatial distribution of precipitation: using
344 percentile thresholds allows us to focus on this despite the model bias.

345 For the Spring case at the 90th percentile threshold (Figure 4a) the radar values drop to
346 zero after 16 hours. After this time radar derived precipitation covers less than 10% of the
347 domain. This demonstrates that the 90th percentile, and other percentile thresholds below
348 the 90th, are not a suitable threshold for radar precipitation accumulations for this case. For
349 all cases (apart from the unlikely event of 100% coverage) there will be a limited area covered
350 by precipitation in both the model and observations, and a corresponding minimum suitable
351 percentile threshold. In an operational situation this minimum threshold could easily be
352 calculated from the fraction of precipitation coverage. All FSS results presented in this
353 paper have been calculated using percentile thresholds above this minimum value.

354 For the Spring case the 8 MOGREPS members from the standard physics configuration
355 are shown in dark gray in Figure 4a and Figure 4b and, although differing by up to 2.5
356 mm in accumulation values (for the 99th percentile threshold), follow the same overall trend
357 throughout the day. This suggests that the ensemble members produce precipitation fea-
358 tures, such as that associated with frontal passage, at similar times. The simulations for all
359 MOGREPS members and the other 4 physics configurations are shown in light gray with the
360 different physics configurations clustering around the corresponding MOGREPS member. In
361 these experiments the different physics configurations have little effect on the precipitation
362 value corresponding to a given percentile threshold. Interestingly, Figure 4a and Figure 4b
363 show peaks in precipitation values at different times: Figure 4a (90th percentile) at a lead

364 time of 20 hrs and Figure 4b (99th percentile) at a lead time of 12 hrs. The higher threshold
365 considers only the areas of convective precipitation, giving a corresponding value that peaks
366 when these storms are strongest whereas the lower threshold includes frontal precipitation
367 and peaks where this is heaviest.

368 The 12 members for the Summer flooding case are shown for thresholds calculated over
369 the full UK domain (dark gray) and limited area domain (light gray). Beyond a lead time
370 of 15 hours, when convection occurred over Edinburgh, values for the limited domain are up
371 to 5 times larger than those over the UK domain. Considering this area separately using
372 percentile thresholds allows the flood producing storm to be investigated. It should be noted
373 that using high value thresholds over the UK domain would also select the Edinburgh area.
374 However, for this highly variable case some ensemble members missed the convection over
375 Edinburgh, and do not produce sufficiently high precipitation values. It is not possible to
376 choose a value threshold that is high enough to select only the area of convection, and yet
377 low enough to include all the ensemble members. Again, this demonstrates the utility of
378 percentile thresholds.

379 **3. Results for LBC and IC perturbations**

380 *a. dFSSmean and eFSSmean*

381 First we consider the realism of the spatial ensemble spread by comparing dFSSmean
382 and eFSSmean for both cases. Both dFSSmean and eFSSmean were calculated over the
383 section of the 2.2 km UK domain with radar coverage (highlighted by the dotted region in
384 Figure 2). Figure 5 shows dFSSmean (left) and eFSSmean (right) for the organized Spring
385 case (top) and flooding case (bottom) calculated for the 99th percentile threshold over the
386 whole UK domain. These results were computed for the 12 members of the flooding cases
387 and 8 MOGREPS members with standard physics for the organized Spring case. To check
388 the validity of comparing these differently sized ensembles, results were also produced for

389 the flooding case when only considering the first 8 ensemble members (not shown). These
390 8 member results differed only in small details from those calculated from 12 members, and
391 lead to the same conclusions, so it was decided to show the results from the full 12 member
392 comparisons.

393 Comparison of the dispersion measures (dFSSmean) for the two cases (Figures 5a and
394 5c) shows that, although these cases are synoptically different, with different convective
395 forcing, the overall behavior is broadly similar. At small scales ensemble members are very
396 different resulting in low values of FSS. FSS values increase as the members become more
397 similar when considered at larger scales. The temporal variability present in the ensemble
398 spread, as measured by dFSSmean, is also clear at this threshold with the scale at which
399 $FSS = 0.5$ varying between 50-500 km for the organized Spring case and 100-250 km for the
400 flooding case. These scales are large because in both cases there is considerable uncertainty
401 in the locations of the showers and showery areas. The temporal variability can be related
402 to the evolution of physical processes. For example, in Figure 5a the area of larger ensemble
403 spread (decrease in dFSSmean) at lead times 13-20 hrs can be linked to greater convective
404 activity and the highest rainfall instances (compare with Figure 4b) and the increase in
405 dFSSmean (lower spread) from 20-25 hrs can be related to a area of spatially predictable
406 frontal precipitation moving into the domain.

407 Overall there is less temporal variability in the FSS for the flooding case. This can
408 again be related to the meteorology of the cases: precipitation in the flooding case was
409 the result of one mechanism, instability associated with a decaying low pressure system,
410 whereas precipitation in the Spring case was associated with both convective showers and
411 frontal passage. Coincidentally, for both cases, the spatial ensemble spread increases with
412 forecast lead time after 20 hours. This up-scaling of forecast spatial differences should be
413 expected from a statistical evaluation of a large number of cases, but not necessarily from
414 individual case studies where the physical processes of the day dominate. Using dFSSmean
415 for individual case studies allows these processes, and their effect on the spatial ensemble

416 spread and upscale growth of forecast differences, to be examined.

417 The error measures (eFSSmean, Figures 5b and 5d) show a similar structure to the
418 dispersion measures with a similar magnitude for ensemble spread and skill. There are
419 times, such as for the Spring case at a lead time of 20 hrs (Figure 5b), or the flooding case
420 at lead times 0-5 hrs (Figure 5d) when the ensemble is clearly under-spread. For the Spring
421 case a timing error results from a front passing into the domain in all members earlier than
422 seen in the radar; for the flooding case convective showers present in the radar have yet to
423 spin up in the model. In both cases there is little evidence that the ensemble is over-spread.

424 For the flooding case dFSSmean and eFSSmean have also been calculated over the 100
425 km limited area domain containing the flooding event. Selecting a subdomain in this man-
426 ner allows us to focus on the spatial predictability of a specific event which can be very
427 different from the UK domain averaged results. Differences between the domains can also
428 be seen in the values corresponding to each percentile threshold as discussed in Section 2d.
429 dFSSmean and eFSSmean, calculated over the 100km domain are shown in Figures 6a and
430 6b respectively at forecast lead times 17 hrs - 26 hrs when convection was seen over Ed-
431 inburgh. Comparison of Figure 6a and Figure 6b suggests that the ensemble spread and
432 skill are similar and that, over this area, the ensemble is capturing the spatial variability of
433 the precipitation well. This gives confidence in the ensemble for a spatially unpredictable
434 flooding event. There are some differences between dFSSmean and eFSSmean, in particular
435 that eFSSmean is more variable with time. This may be partly due to both the smaller num-
436 ber of comparisons in the error calculation, and also reflects differences between the model
437 and observations in the temporal evolution of the storm. Note that, as the 99th percentile
438 threshold corresponds to different precipitation values over the UK and Edinburgh domains,
439 we cannot do a direct comparison between Figures 5 and 6. This also suggests that we are
440 indeed looking at different processes or phenomena with the different domains and confirms
441 the need to use a suitable domain size to examine the spatial variability of particular fea-
442 tures. The domain must be large enough to give representative results, but small enough to

443 focus on the phenomena of interest. Of course, the same remarks will be true of any spatial
444 measure.

445 *b. dFSSstdev in addition to dFSSmean*

446 In this section we discuss the benefits of considering dFSSstdev in addition to dFSSmean.
447 Figure 7 shows dFSSmean and dFSSstdev calculated for the organized Springcase (top) and
448 flooding case (bottom) when considering the 99th percentile threshold for hourly precipitation
449 accumulations. The FSS was calculated over the whole UK domain. dFSSstdev is shown
450 in Figure 7c and Figure 7d and presents results consistent with those from dFSSmean. For
451 example, the largest values of dFSSstdev occur in areas where low dFSSmean values extend
452 to large scales. The greater spatial spread associated with low values of dFSSmean results
453 in a wider range of possible values for FSS_i and larger dFSSstdev.

454 However, there is also some further information given by the standard deviation. In
455 particular, for the flooding case (Figure 7d) there is an area of higher standard deviation
456 seen in the first two hours of the forecast at neighborhood sizes up to 500km which is
457 associated with the spin-up of the model. This effect is even more apparent in results for
458 the 99.9th percentile threshold (not shown) and is the result of the convection permitting
459 model having to spin up showers during the first few hours of the forecast. Because the
460 ensemble members spin-up showers at different locations, lower values of dFSSmean and a
461 large range of values of FSS_i (resulting in a large dFSSstdev) are obtained. For the spring
462 case (Figures 7a and 7b) convective showers are not present at the forecast start time and
463 do not need to be spun-up from initial conditions. Hence spin-up effects are not seen in
464 the precipitation diagnostics. It is useful to examine how the standard deviation behaves at
465 different scales. The smallest values are found at both the grid scale, where differences are
466 so large that similarly low values of the FSS are expected for all member pairs, and also at
467 the largest scales, where all members are effectively the same.

468 *c. Other fields and thresholds*

469 The use of different percentile thresholds allows more information to be gained about
470 the ensemble spread for different ranges of forecast values, for example a higher threshold
471 will select more extreme values compared to a lower threshold which will select values that
472 are more widespread. An example is given in Figure 8 for the organized Spring case where
473 results for the top 99th (LHS) and 85th (RHS) percentiles are compared. This time we show
474 a different diagnostic field, the 10 m horizontal wind speed. Like the hourly precipitation
475 accumulations this field was selected as a suitable candidate for calculation of the FSS
476 because of its high spatial variability. 10 m wind speeds are also used by the Met Office for
477 routine forecast verification.

478 The 99th percentile threshold selects only the highest wind speeds in the domain. At
479 lead times 0-10 hrs these are found in to areas to the north of the UK near the low pressure
480 center. The exact placement of the highest winds varied considerably between the ensemble
481 members, with some placing them to the northwest and others to the northeast of the UK.
482 Hence there were large spatial differences between the members resulting in low dFSSmean
483 values extending to large neighborhoods at a lead time of 10 hrs as shown in Figure 8a.
484 At lead times greater than 10 hours there is high spatial agreement amongst the ensemble
485 members resulting in high values of dFSSmean. All members place the highest winds to the
486 northwest of the UK associated with the frontal feature that enters the domain at this time.

487 Comparing Figure 8a and Figure 8b we see the unusual result that for a lead time of 12
488 hours, and after 28 hours, there is more agreement (larger FSS values) for the 99th than for
489 the 85th percentile for a given neighborhood size. This behavior suggests that care must be
490 taken in the interpretation the 99th percentile threshold for the wind speed field. For the
491 wind speed, local variability is superimposed upon a background gradient from the large scale
492 situation. The 99th percentile is likely to include both local variability from points where the
493 background field is moderate and also larger scale variability where the background field is
494 high. Consequently, unlike for precipitation, we cannot cleanly examine local features in the

495 wind speed field simply by selecting a high threshold value. It is necessary to also consider a
496 lower threshold that includes features of the larger scale flow such as, for this case, the 85th
497 percentile threshold. Figure 8b shows that, at lead times 12-20 hrs, the FSS values for the
498 85th percentile are particularly high. These areas of small spatial spread can be related to
499 the synoptic situation: at a lead time of 12 hrs a highly predictable frontal feature entered
500 the domain from the NW and the top 15 % of wind speeds in the domain were closely
501 associated with the flow in the vicinity of this front. Hence, there was very high spatial
502 agreement between the members at these times. Before the front entered the domain the
503 highest winds were associated with a less predictable decaying cold front. Moreover, after
504 the front had progressed further into the domain greater differences between the members
505 emerged at larger scales for the winds to the south of the occluded front.

506 The effect of different thresholds on the FSS for hourly precipitation accumulations can
507 be seen by comparing Figures 5a and 5c with Figures 9a and 9b respectively. The latter show
508 dFSSmean calculated for the 90th percentile threshold. In particular, it can be seen that the
509 large temporal variability seen in Figures 5a and 5c for the 99th threshold has been replaced
510 in the 90th percentile results by a trend for ensemble spread to increase systematically with
511 time. This trend is expected climatologically as forecast differences grow from small to larger
512 scales with increasing forecast lead time. The rate of increase is different for the two cases.
513 For the flooding case (Figure 9b) scales at which dFSSmean=0.5 increase gradually from 5
514 km to 100 km over 36 hours as forecast differences grow from small to larger scales. For the
515 Spring case, dFSSmean values greater than 0.5 are seen even at the grid scale for lead times
516 up to 25 hrs. After this time the scale at which dFSSmean=0.5 increases rapidly to 225km.
517 This pattern is in agreement with the behavior seen for the 99th threshold and has the
518 same interpretation: after 25 hrs an area of precipitation moves out of the domain but with
519 timing differences between the members. Overall, there is better spatial agreement between
520 the ensemble members at the 90th percentile threshold than at the 99th: the broader-scale
521 features selected by the lower threshold are more predictable. When considering a range of

522 different thresholds from the 99th to 80th percentile (not shown) the transition from large
523 temporal variability to a trend of upscale growth of forecast differences with increasing lead
524 time was found to be smooth: there is no sudden transition. It is likely that the range of
525 thresholds over which such a transition occurs will be highly case dependent as the relative
526 importance of local and large scale features changes. The FSS allows such behavior to be
527 investigated.

528 4. Results assessing different physics configurations

529 In this section we present an application of dFSSmean to the comparison of the mul-
530 tiphysics and MOGREPS ensembles for the organized Spring case. Thus we compare the
531 spatial ensemble spread associated with LBC and IC perturbations to that generated through
532 different physics configurations as described in Section 2c. The examples presented are for
533 the 99th percentile threshold of precipitation accumulation: lower thresholds showed smaller
534 spatial differences (larger dFSSmean values) but lead to the same general conclusions. Note
535 that the purpose is not to evaluate the merits of particular physics configurations but to
536 show a method that can be used to examine the behavior of stochastic processes or physics
537 changes in ensembles.

538 Figure 10b shows dFSSmean comparing the configuration with restricted convection
539 scheme and increased time step (conv+time) to that with the modified treatment of graupel
540 (grp) using the Physics2 comparison strategy (comparison strategy ii in Section 2c). This
541 comparison strategy is shown because it gives larger spatial differences than those found
542 when comparing any other physics configuration pairs, or considering all physics configura-
543 tions (the Physics5 comparison strategy). In Figure 10b FSS values of 0.5 are reached by
544 a neighborhood size of 5 km, and no spatial differences are seen for neighborhoods greater
545 than 100 km (where FSS=1). The lowest values of dFSSmean occur between lead times
546 of 12 hrs and 16 hrs when the heaviest convective showers were present: it is during these

547 events that modifications to the treatment of graupel are most noticeable.

548 Results from comparing only the MOGREPS members from conv+time and grp (com-
549 parison strategy MOGREPS2, i in Section 2c) are shown in Figure 10a. These differ only
550 in minor details from those shown in Figure 7a (dFSSmean calculated for the MOGREPS
551 ensemble with the standard physics configuration). The MOGREPS2 results show that FSS
552 values of 0.5 are reached on scales greater than 60 km, scales at which the Physics2 mem-
553 bers are almost identical. In other words, the spatial variation introduced through different
554 physics configurations is only seen close to the grid scale. If we consider FSS values lower
555 than $FSS = 0.5$ to represent fields so different that the forecast is no longer useful, then
556 the different physics configurations applied here, for this particular case, are simply mov-
557 ing around features that are known to be unpredictable from the MOGREPS ensemble.
558 Of course, this is not to say that physics changes in general are unimportant for improving
559 model performance, or that using different physics configurations is not sometimes a valuable
560 component of an ensemble system, or that adding small scale perturbations is undesirable
561 or that, for another case or for other physics perturbations the effects might be very differ-
562 ent. Our purpose is simply to demonstrate a methodology that allows the spatial effects of
563 different ensemble configurations to be thoroughly investigated and set into the context of
564 other aspects of forecast uncertainty.

565 It is possible that, although the evaluation of Physics2 only showed forecast differences
566 at small spatial scales, combining the different physics configurations with those from the
567 MOGREPS2 ensemble may lead to large changes in the growth of spatial differences. To
568 assess this, the comparison strategy MOGREPS2+Physics2 (comparison strategy iii in Sec-
569 tion 2c) is employed. Again, examples are shown for the physics configurations conv+time
570 and grp which show the largest spatial differences. The results of MOGREPS2+Physics2
571 are shown in Figure 10c. Differences between Figure 10c and Figure 10a are very small and
572 hence, to aid interpretation, Figure 10d shows the difference between the MOGREPS2 and
573 the MOGREPS2+Physics2 results. The differences are over an order of magnitude smaller

574 than the dFSSmean values in Figures 10a and 10c. It is interesting that both positive and
575 negative differences are seen: modifying the different physics configuration both adds and
576 removes spatial spread. From Figure 10d it can also be seen that differences between MO-
577 GREPS2 and the MOGREPS2+Physics2 extend, with similar magnitude, across all spatial
578 scales. However, in terms of the fractional difference relative to dFSSmean the differences at
579 small neighborhoods have more importance. At a lead time of 15 hrs the fractional differ-
580 ence in dFSSmean varies from 7% at 50 km to 3% at 250 km. It should be noted that these
581 differences are still very small, especially at the larger more predictable scales (as indicated
582 by the point where $FSS \geq 0.5$ in the MOGREPS ensemble).

583 Analysis of the combined MOGREPS+Physics comparisons supports the conclusions
584 drawn previously that the introduction of these differences in the physics only influences
585 scales much smaller than the predictable scales of the system (in this particular experiment).
586 In practical terms, the variability of those scales could be addressed with spatial post pro-
587 cessing and without the need for additional ensemble members. On the other hand, if the
588 scales of the physics changes were to upscale to scales greater than the system's predictable
589 scales then the performance of the ensemble might benefit from more perturbed-physics
590 members. Systematic application of the methods shown here would provide a sound basis
591 for making these decisions.

592 **5. Discussion and conclusions**

593 In this paper we have presented, with examples, a new methodology for the detailed anal-
594 ysis of ensemble spread for high resolution forecasts focusing on spatial variability. In par-
595 ticular we focused on two different measures of ensemble spread: dFSSmean and dFSSstdev,
596 the mean and standard deviation of the FSS calculated over all ensemble member-member
597 pairs. dFSSmean gives a measure of the FSS value for the whole ensemble indicating the
598 average spatial agreement within the ensemble over a particular size of neighborhood i.e at a

599 given spatial scale. dFSSstdev provides some further useful information about the range of
600 FSS values used in the calculation of dFSSmean. A large range of FSS values, corresponding
601 to a large value of dFSSstdev, indicates that the ensemble members are unevenly distributed.

602 To demonstrate the utility of these measures results were presented from two case studies.
603 It was shown that dFSSmean and dFSSstdev allowed investigation of, for example, the
604 temporal evolution of ensemble spread, model spin up, and saturation of forecast differences.
605 Considering different percentile thresholds allowed information to be gained about the spatial
606 spread of the ensemble for different physical regimes. In particular it was found that, for
607 hourly precipitation accumulations, the dFSSmean for the 99th percentile threshold had high
608 temporal variability. This contrasted with the dFSSmean for the 90th percentile threshold
609 for which spatial differences between the ensemble members increased with time.

610 The realism of the ensemble spatial distribution was also tested by comparison with
611 another metric, the mean FSS calculated over all member-radar pairs, denoted eFSSmean.
612 This error measure can be compared with dFSSmean to investigate the spread-skill rela-
613 tionship of the ensemble at different times and spatial scales. For the two cases considered
614 here these measures suggested that ensemble spread was reasonable. On occasion the en-
615 semble was under-spread and this was linked to timing errors between the simulations and
616 the observations and to the need for spin up of showers in a convection permitting model.

617 For one case study, results were presented for a comparison of spread between differently
618 generated ensembles, including multiple physics configurations. This application illustrates a
619 methodology for identifying the spatial scales that are influenced by modifications to physical
620 processes. Examining the FSS for different spatial scales and over a range of times allowed
621 a quantification of the effects of using different physics configurations compared to LBC and
622 IC perturbations. For the case described here it was concluded that modifying the physics
623 for this case did not influence the ensemble evolution at scales where the forecast has skill.
624 These results are not to be interpreted as general: well chosen physics modifications can and
625 do improve forecasts as demonstrated by, for example by Stensrud et al. (2000); Keil et al.

626 (2013). The key point is that evaluation techniques presented here allow clear statements
627 about the impacts of physics modifications to be made since different ensemble configurations
628 can be thoroughly investigated and the spatial impact of the changes quantified.

629 The work presented here provides a framework through which spatial ensemble spread
630 can be analyzed. There are some limitations to this study: in particular the consideration of
631 two cases only and the limited consideration of physics perturbations. It is left to future work
632 to apply these methods to a larger sample of cases, and different, more realistic, multiphysics
633 ensembles or other model error inclusion schemes. Another limiting factor is the methodology
634 of calculating a single value of the FSS that is representative of a comparison across a whole
635 domain. As discussed above this can mean that different meteorological phenomena, such as
636 convective and frontal precipitation, are compared together, when each individually may have
637 an inherently different predictability and ensemble spread. It is possible to select a smaller
638 domain to consider events of interest, as highlighted with respect to Figure 6, although this
639 is only useful in hindsight once the event has occurred. Hence, future work is intended to
640 develop a spatially varying and scale dependent measure of ensemble spread that does not
641 suffer from this drawback.

642 Despite these limitations there are some important conclusions from this work. In partic-
643 ular, we have stressed how the ensemble spread is highly dependent on the scales considered
644 for evaluation. Consequently, to investigate the ensemble behavior thoroughly it is neces-
645 sary to consider multiple scales, and be mindful of the different expectations for skill at these
646 scales. Forecasts should be verified, and the benefits of forecast model changes assessed, at
647 scales that are believable. This paper has provided a methodology for determining such
648 believable scales and their temporal evolution. With future movement to higher and higher
649 resolution models the distinction between the grid scale and the believable scales is becoming
650 increasingly important.

651 *Acknowledgments.*

652 The authors would like to thank the three anonymous reviewers for their detailed com-
653 ments that have helped improve the quality and clarity of this article. S.Dey acknowledges
654 support from a NERC PhD studentship with CASE support from the UK Met Office. Initial
655 work contributing to this paper was completed by the same author during a Summer Place-
656 ment at MetOffice@Reading. S.Migliorini acknowledges support from the NERC National
657 Center for Earth Observation.

REFERENCES

- 660 Ancell, B. C., 2013: Nonlinear characteristics of ensemble perturbation evolution and their
661 application to forecasting high-impact events. *Wea. Forecasting*, **28** (6).
- 662 Baker, L., A. Rudd, S. Migliorini, and R. Bannister, 2014: Representation of model error in
663 a convective-scale ensemble prediction system. *Nonlinear Proc. Geophys.*, **21** (1), 19–39.
- 664 Berner, J., S.-Y. Ha, J. Hacker, A. Fournier, and C. Snyder, 2011: Model uncertainty in
665 a mesoscale ensemble prediction system: Stochastic versus multiphysics representations.
666 *Mon. Wea. Rev.*, **139** (6), 1972–1995.
- 667 Bowler, N. E., A. Arribas, S. E. Beare, K. R. Mylne, and G. J. Shutts, 2009: The local
668 ETKF and SKEB: Upgrades to the MOGREPS short-range ensemble prediction system.
669 *Quart. J. Roy. Meteor. Soc.*, **135** (640), 767–776.
- 670 Bowler, N. E., A. Arribas, K. R. Mylne, K. B. Robertson, and S. E. Beare, 2008: The MO-
671 GREPS short-range ensemble prediction system. *Quart. J. Roy. Meteor. Soc.*, **134** (632),
672 703–722.
- 673 Caron, J.-F., 2013: Mismatching perturbations at the lateral boundaries in limited-area
674 ensemble forecasting: A case study. *Mon. Wea. Rev.*, **141** (1), 356–374.
- 675 Clark, A. J., et al., 2011: Probabilistic precipitation forecast skill as a function of ensemble
676 size and spatial scale in a convection-allowing ensemble. *Monthly Weather Review*, **139** (5),
677 1410–1418.
- 678 Craig, G. C., C. Keil, and D. Leuenberger, 2012: Constraints on the impact of radar rain-
679 fall data assimilation on forecasts of cumulus convection. *Quart. J. Roy. Meteor. Soc.*,
680 **138** (663), 340–352.

681 Davies, T., M. J. P. Cullen, A. J. Malcolm, M. H. Mawson, A. Staniforth, A. A. White, and
682 N. Wood, 2005: A new dynamical core for the met office’s global and regional modelling
683 of the atmosphere. *Quart. J. Roy. Meteor. Soc.*, **131 (608)**, 1759–1782.

684 Derbyshire, S., A. Maidens, S. Milton, R. Stratton, and M. Willett, 2011: Adaptive detrain-
685 ment in a convective parametrization. *Quart. J. Roy. Meteor. Soc.*, **137 (660)**, 1856–1871.

686 Done, J. M., G. C. Craig, S. L. Gray, and P. A. Clark, 2012: Case-to-case variability
687 of predictability of deep convection in a mesoscale model. *Quart. J. Roy. Meteor. Soc.*,
688 **138 (664)**, 638–648.

689 Duc, L., K. Saito, and H. Seko, 2013: Spatial-temporal fractions verification for high-
690 resolution ensemble forecasts. *Tellus A*, **65**.

691 Duda, J. D. and W. A. Gallus, Jr, 2013: The impact of large-scale forcing on skill of simulated
692 convective initiation and upscale evolution with convection-allowing grid spacings in the
693 WRF. *Wea. Forecasting*, **28 (4)**.

694 Ebert, E. E., 2008: Fuzzy verification of high-resolution gridded forecasts: a review and
695 proposed framework. *Meteorol. Appl.*, **15 (1)**, 51–64.

696 Edwards, J. M. and A. Slingo, 1996: Studies with a flexible new radiation code. I: Choosing
697 a configuration for a large-scale model. *Quart. J. Roy. Meteor. Soc.*, **122 (531)**, 689–719.

698 Ehrendorfer, M., 1997: Predicting the uncertainty of numerical weather forecasts: a review.
699 *Meteorol. Z.*, **6**, 147–183.

700 Essery, R., M. Best, and P. Cox, 2001: MOSES 2.2 technical documentation. Tech. rep.,
701 Hadley Centre Technical Note.

702 Gebhardt, C., S. Theis, M. Paulat, and Z. Ben Bouallègue, 2011: Uncertainties in COSMO-
703 DE precipitation forecasts introduced by model perturbations and variation of lateral
704 boundaries. *Atmos. Res.*, **100 (2)**, 168–177.

- 705 Gilleland, E., D. Ahijevych, B. G. Brown, B. Casati, and E. E. Ebert, 2009: Intercomparison
706 of spatial forecast verification methods. *Wea. Forecasting*, **24** (5), 1416–1430.
- 707 Golding, B. W., 1998: Nimrod: a system for generating automated very short range forecasts.
708 *Meteorol. Appl.*, **5** (1), 1–16.
- 709 Gregory, D. and P. R. Rowntree, 1990: A mass flux convection scheme with representation of
710 cloud ensemble characteristics and stability-dependent closure. *Mon. Wea. Rev.*, **118** (7),
711 1483–1506.
- 712 Hacker, J., C. Snyder, S.-Y. HA, and M. Pocerlich, 2011: Linear and non-linear response to
713 parameter variations in a mesoscale model. *Tellus A*, **63** (3), 429–444.
- 714 Hanley, K., D. Kirshbaum, N. Roberts, and G. Leoncini, 2013: Sensitivities of a squall line
715 over central europe in a convective-scale ensemble. *Mon. Wea. Rev.*, **141** (1), 112–133.
- 716 Hanley, K. E., D. J. Kirshbaum, S. E. Belcher, N. M. Roberts, and G. Leoncini, 2011:
717 Ensemble predictability of an isolated mountain thunderstorm in a high-resolution model.
718 *Quart. J. Roy. Meteor. Soc.*, **137** (661), 2124–2137.
- 719 Harrison, D. L., S. J. Driscoll, and M. Kitchen, 2000: Improving precipitation estimates from
720 weather radar using quality control and correction techniques. *Meteorol. Appl.*, **7** (2), 135–
721 144.
- 722 Hohenegger, C., D. Lüthi, and C. Schär, 2006: Predictability mysteries in cloud-resolving
723 models. *Mon. Wea. Rev.*, **134** (8), 2095–2107.
- 724 Hohenegger, C. and C. Schär, 2007a: Atmospheric predictability at synoptic versus cloud-
725 resolving scales. *Bull. Amer. Meteor. Soc.*, **88** (7), 1783–1793.
- 726 Hohenegger, C. and C. Schär, 2007b: Predictability and error growth dynamics in cloud-
727 resolving models. *J. Atmos. Sci.*, **64** (12), 4467–4478.

728 Johnson, A. and X. Wang, 2013: Object-based evaluation of a storm-scale ensemble during
729 the 2009 NOAA hazardous weather testbed spring experiment. *Mon. Wea. Rev.*, **141 (3)**,
730 1079–1098.

731 Johnson, A., et al., 2014: Multiscale characteristics and evolution of perturbations for warm
732 season convection-allowing precipitation forecasts: Dependence on background flow and
733 method of perturbation. *Mon. Wea. Rev.*, **142 (3)**, 1053–1073.

734 Keil, C. and G. C. Craig, 2011: Regime-dependent forecast uncertainty of convective pre-
735 cipitation. *Meteorol. Z.*, **20 (2)**, 145–151.

736 Keil, C., F. Heinlein, and G. C. Craig, 2013: The convective adjustment time-scale as
737 indicator of predictability of convective precipitation. *Quart. J. Roy. Meteor. Soc.*, doi:
738 10.1002/qj.2143.

739 Kühnlein, C., C. Keil, G. C. Craig, and C. Gebhardt, 2013: The impact of downscaled initial
740 condition perturbations on convective-scale ensemble forecasts of precipitation. *Quart. J.*
741 *Roy. Meteor. Soc.*, doi:10.1002/qj.2238.

742 Lean, H. W., P. A. Clark, M. Dixon, N. M. Roberts, A. Fitch, R. Forbes, and C. Halli-
743 well, 2008: Characteristics of high-resolution versions of the Met Office Unified Model for
744 forecasting convection over the United Kingdom. *Mon. Wea. Rev.*, **136 (9)**, 3408 – 3424.

745 Leoncini, G., R. S. Plant, S. L. Gray, and P. A. Clark, 2010: Perturbation growth at the
746 convective scale for CSIP IOP18. *Quart. J. Roy. Meteor. Soc.*, **136 (648)**, 653–670.

747 Leoncini, G., R. S. Plant, S. L. Gray, and P. A. Clark, 2012: Ensemble forecasts of a flood-
748 producing storm: comparison of the influence of model-state perturbations and parameter
749 modifications. *Quart. J. Roy. Meteor. Soc.*, **139 (670)**, 198–211.

750 Leoncini, G., N. Roberts, and B. Golding, 2011: 8th July 2011 Floods in Scotland. report
751 for Scottish Environment Protection Agency, Met Office, 18 pp.

- 752 Lock, A., A. Brown, M. Bush, G. Martin, and R. Smith, 2000: A new boundary layer
753 mixing scheme. Part I: Scheme description and single-column model tests. *Mon. Wea.*
754 *Rev.*, **128** (9), 3187–3199.
- 755 Lorenz, E. N., 1969: The predictability of a flow which possesses many scales of motion.
756 *Tellus*, **21** (3), 289–307.
- 757 Martin, W. J. and M. Xue, 2006: Sensitivity analysis of convection of the 24 May 2002 IHOP
758 case using very large ensembles. *Mon. Wea. Rev.*, **134** (1), 192–207.
- 759 Mass, C. F., D. Ovens, K. Westrick, and B. A. Colle, 2002: Does increasing horizontal
760 resolution produce more skillful forecasts? *Bull. Amer. Meteor. Soc.*, **83** (3), 407–430.
- 761 Migliorini, S., M. Dixon, R. Bannister, and S. Ballard, 2011: Ensemble prediction for now-
762 casting with a convection-permitting model-I : description of the system and the impact
763 of radar-derived surface precipitation rates. *Tellus A*, **63** (3), 468–496.
- 764 Mittermaier, M. and N. Roberts, 2010: Intercomparison of spatial forecast verification meth-
765 ods: Identifying skillful spatial scales using the fractions skill score. *Wea. Forecasting*,
766 **25** (1), 343–354.
- 767 Mylne, K., 2013: Scientific framework for the ensemble prediction system for the UKV.
768 MOSAC PAPER 18.6, UK Meteorological Office. URL [http://www.metoffice.gov.uk/
769 media/pdf/q/0/MOSAC_18.6_Mylne.pdf](http://www.metoffice.gov.uk/media/pdf/q/0/MOSAC_18.6_Mylne.pdf).
- 770 Palmer, T. N., 2000: Predicting uncertainty in forecasts of weather and climate. *Reports on*
771 *Progress in Physics*, **63** (2), 71–116.
- 772 Rezacova, D., P. Zacharov, and Z. Sokol, 2009: Uncertainty in the area-related QPF for
773 heavy convective precipitation. *Atmos. Res.*, **93** (1), 238–246.
- 774 Roberts, N., 2008: Assessing the spatial and temporal variation in the skill of precipitation
775 forecasts from an NWP model. *Meteorol. Appl.*, **15** (1), 163–169.

776 Roberts, N. M., 2003: The impact of a change to the use of the convection scheme
777 to high resolution simulations of convective events. Forecasting Research Techni-
778 cal Report 407, Met Office, 30 pp. URL [http://www.metoffice.gov.uk/archive/
779 forecasting-research-technical-report-407](http://www.metoffice.gov.uk/archive/forecasting-research-technical-report-407), available from the National Meteorolo-
780 gical Libray and Archive.

781 Roberts, N. M. and H. W. Lean, 2008: Scale-selective verification of rainfall accumulations
782 from high-resolution forecasts of convective events. *Mon. Wea. Rev.*, **136** (1), 78–97.

783 Schwartz, C. S., et al., 2010: Toward improved convection-allowing ensembles: Model physics
784 sensitivities and optimizing probabilistic guidance with small ensemble membership. *Wea.*
785 *Forecasting*, **25** (1), 263–280.

786 Stensrud, D. J., J.-W. Bao, and T. T. Warner, 2000: Using initial condition and model
787 physics perturbations in short-range ensemble simulations of mesoscale convective systems.
788 *Mon. Wea. Rev.*, **128** (7), 2077–2107.

789 Surcel, M., I. Zawadzki, and M. K. Yau, 2014: On the filtering properties of ensemble
790 averaging for storm-scale precipitation forecasts. *Mon. Wea. Rev.*, **142** (3), 1093–1105.

791 Vié, B., G. Molinié, O. Nuissier, B. Vincendon, V. Ducrocq, F. Bouttier, and E. Richard,
792 2012: Hydro-meteorological evaluation of a convection-permitting ensemble prediction
793 system for mediterranean heavy precipitating events. *Nat. Hazards Earth Sys.*, **12** (8),
794 2631–2645.

795 Walser, A., D. Lüthi, and C. Schär, 2004: Predictability of precipitation in a cloud-resolving
796 model. *Mon. Wea. Rev.*, **132** (2), 560–577.

797 Walser, A. and C. Schär, 2004: Convection-resolving precipitation forecasting and its pre-
798 dictability in alpine river catchments. *J. Hydrol. (Amsterdam, Neth.)*, **288** (12), 57 –
799 73.

- 800 Wilkinson, J., 2011: The large-scale precipitation parametrization scheme. Unified Model
801 Documentation Paper 26, Unified Model Version: 8.0, Met Office, 53 pp.
- 802 Wilson, D. R. and S. P. Ballard, 1999: A microphysically based precipitation scheme for
803 the UK Meteorological Office Unified Model. *Quart. J. Roy. Meteor. Soc.*, **125** (557),
804 1607–1636.
- 805 Zacharov, P. and D. Rezacova, 2009: Using the fractions skill score to assess the relationship
806 between an ensemble QPF spread and skill. *Atmos. Res.*, **94** (4), 684–693.
- 807 Zhang, F., 2005: Dynamics and structure of mesoscale error covariance of a winter cyclone
808 estimated through short-range ensemble forecasts. *Mon. Wea. Rev.*, **133** (10), 2876–2893.
- 809 Zimmer, M., G. C. Craig, C. Keil, and H. Wernli, 2011: Classification of precipitation events
810 with a convective response timescale and their forecasting characteristics. *Geophys. Res.*
811 *Lett.*, **38**, doi:10.1029/2010GL046199.

List of Figures

- 812
- 813 1 UK Met Office surface analysis valid at (a) 18 UTC on 23 April 2011 and (b)
814 06 UTC on 8 July 2011. Courtesy of the Met Office. Crown copyright. 36
- 815 2 Domains of the UK 2.2 km model (light gray), 100 km subdomain for the
816 Summer flooding case (dark gray) and areas of radar coverage (dotted). 37
- 817 3 Two different idealized spatial distributions of precipitation. Individual en-
818 semble members (shown in white) position the precipitation in different spatial
819 locations. The control simulation (shown in filled black) may produce precip-
820 itation in the center of that produced by individual ensemble members as
821 shown in (a) or at the edge of the ensemble as shown in (b). Considering only
822 the spatial separation of member-member pairs (solid arrows) indicates that
823 (a) and (b) have the same spatial ensemble spread. Including both member-
824 control and member-member pairs allows the differences in spread between
825 (a) and (b) to be detected. 38
- 826 4 Hourly precipitation accumulation values corresponding to the 90th (a,c) and
827 99th (b,d) percentile thresholds. For the organized Spring case (top) results
828 from all the simulations are shown. To highlight the grouping of members
829 those with the standard physics configuration are shown in dark gray and those
830 from other physics configurations in light gray. For the flooding case (bottom)
831 percentile thresholds calculated using data for the full the UK domain are
832 shown in dark gray, and those for the limited area domain are shown in light
833 gray. Radar data is shown from the area of the UK domain with radar coverage
834 (black with circles) and, in (c,d) over the limited area domain (black with
835 crosses). 39

836 5 dFSSmean (a,c) and eFSSmean (b,d) for the organized Spring case (top) and
837 the Summer flooding case (bottom). The standard physics configuration and
838 the 99th percentile threshold are considered. The white dashed line at 0.5
839 represents the believable scale. Results were calculated over the area of the
840 UK domain with radar coverage. 40

841 6 FSS calculations over the Edinburgh subdomain: (a) dFSSmean and (b) eF-
842 SSmean. The 99th percentile threshold is considered. The white dashed line
843 at 0.5 represents the believable scale. 41

844 7 dFSSmean (a,c) and dFSSstdev (b,d) for the organized Spring case (top) and
845 the flooding case (bottom). The white dashed line in (a,c) at 0.5 represents
846 the believable scale. To guide the eye, in (b,d) the white dashed line at 0.05
847 represents the neighborhood at which dFSSstdev is an order of magnitude
848 smaller than the believable scale. The 99th percentile threshold is considered
849 and results are calculated over the whole UK domain. 42

850 8 Comparison of dFSSmean calculated for the (a) 99th and (b) 85th percentile
851 thresholds for the 10 m horizontal wind speed field and the organized Spring
852 case. Results are calculated over the whole of the UK domain and only the
853 standard physics configuration is considered. The white dashed line at 0.5
854 represents the believable scale. 43

855 9 dFSSmean calculated using the 90th percentile threshold of hourly precipita-
856 tion accumulations for (a) the organized Spring and (b) the Summer flooding
857 case. Results are calculated over the whole of the UK domain and only the
858 standard physics configuration is considered. The white dashed line at 0.5
859 represents the believable scale. 44

860 10 dFSSmean comparisons of the restricted convection with increased time step
861 and graupel physics configurations for the 99th percentile threshold of hourly
862 precipitation accumulations. Results from different comparison strategies are
863 shown: (a) MOGREPS2, (b)Physics2 and (c) MOGREPS2+Physics2. (d)
864 shows the difference between sub-figures (c) and (a). Results are calculated
865 over the whole of the UK domain. The white dashed line at 0.5 represents
866 the believable scale.

45

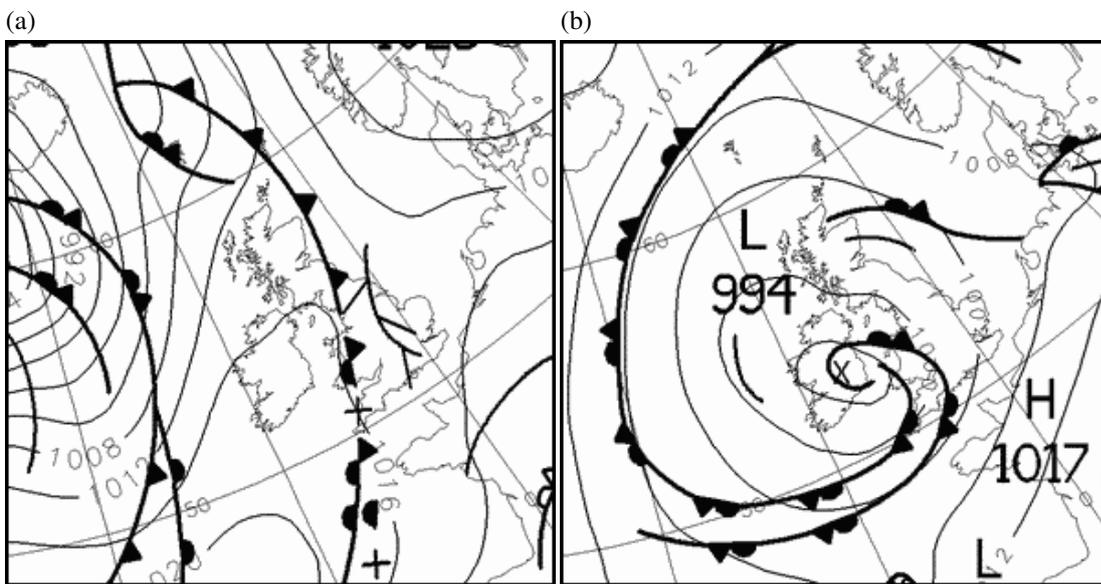


FIG. 1. UK Met Office surface analysis valid at (a) 18 UTC on 23 April 2011 and (b) 06 UTC on 8 July 2011. Courtesy of the Met Office. Crown copyright.

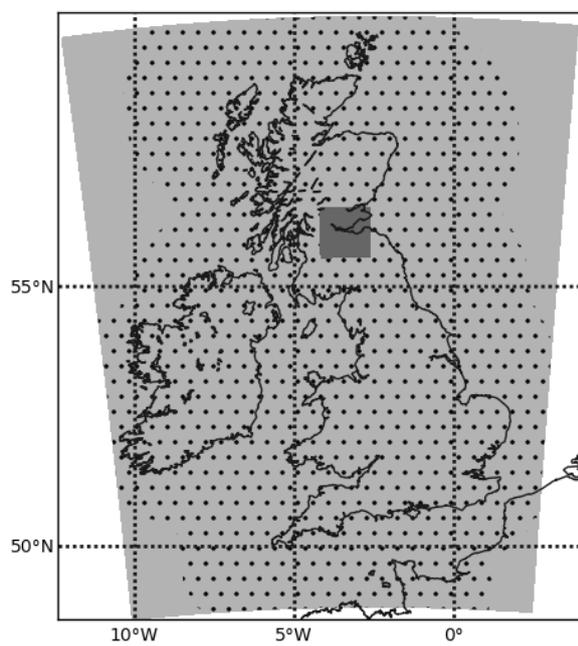


FIG. 2. Domains of the UK 2.2 km model (light gray), 100 km subdomain for the Summer flooding case (dark gray) and areas of radar coverage (dotted).

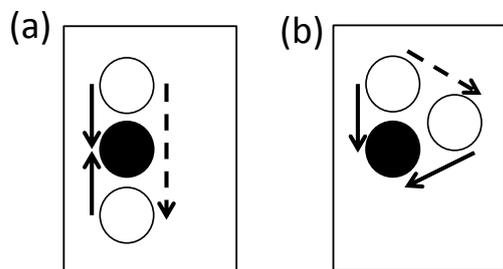


FIG. 3. Two different idealized spatial distributions of precipitation. Individual ensemble members (shown in white) position the precipitation in different spatial locations. The control simulation (shown in filled black) may produce precipitation in the center of that produced by individual ensemble members as shown in (a) or at the edge of the ensemble as shown in (b). Considering only the spatial separation of member-member pairs (solid arrows) indicates that (a) and (b) have the same spatial ensemble spread. Including both member-control and member-member pairs allows the differences in spread between (a) and (b) to be detected.

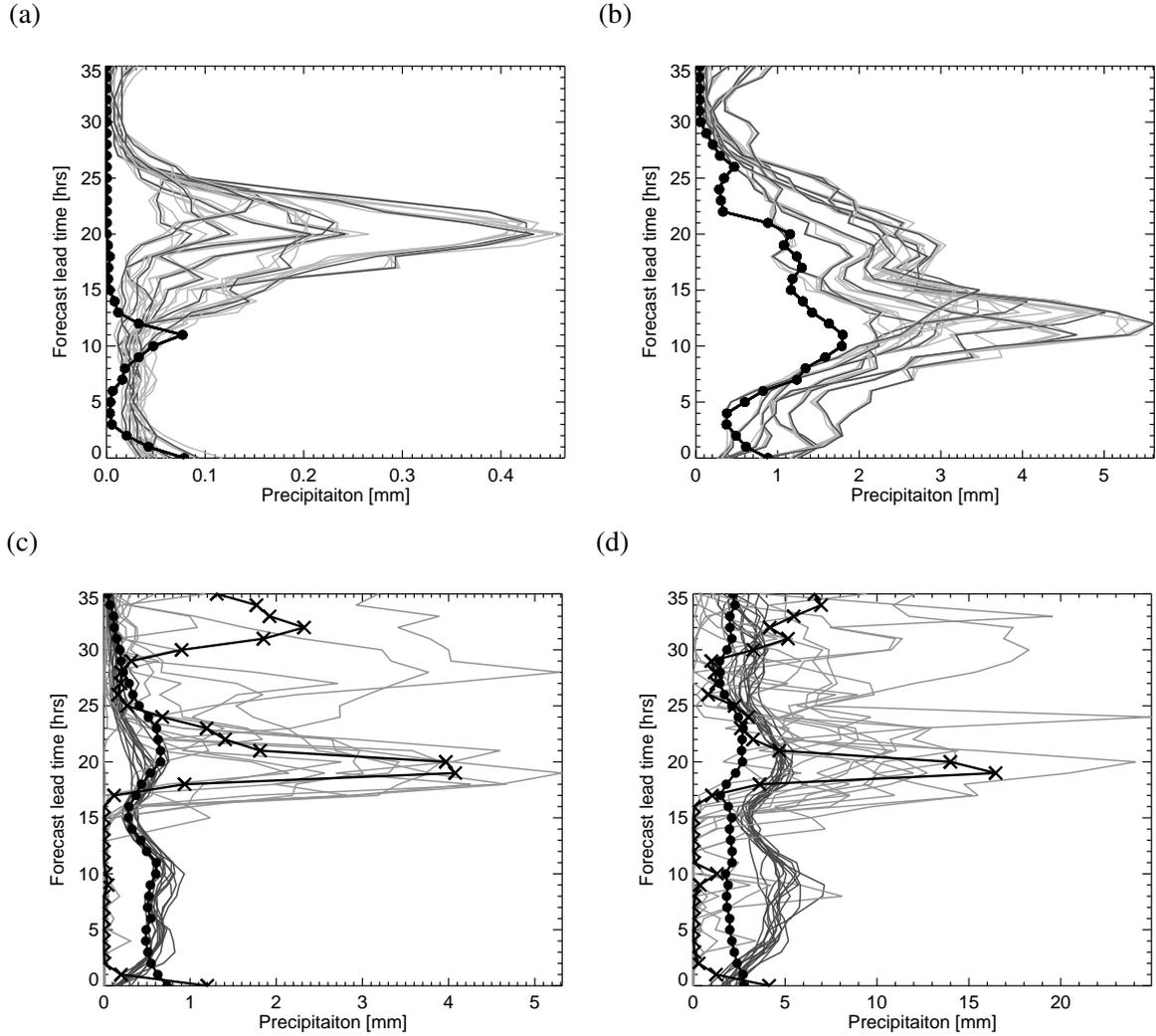


FIG. 4. Hourly precipitation accumulation values corresponding to the 90th (a,c) and 99th (b,d) percentile thresholds. For the organized Spring case (top) results from all the simulations are shown. To highlight the grouping of members those with the standard physics configuration are shown in dark gray and those from other physics configurations in light gray. For the flooding case (bottom) percentile thresholds calculated using data for the full the UK domain are shown in dark gray, and those for the limited area domain are shown in light gray. Radar data is shown from the area of the UK domain with radar coverage (black with circles) and, in (c,d) over the limited area domain (black with crosses).

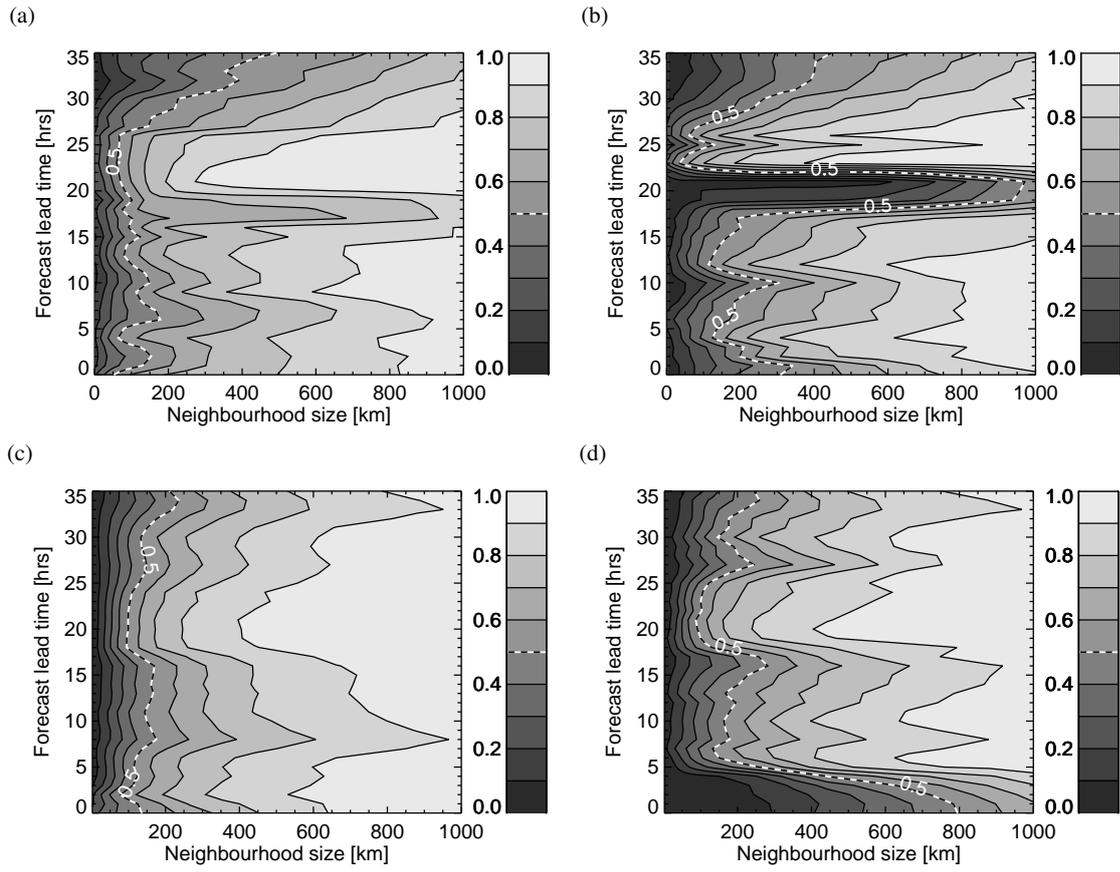


FIG. 5. dFSSmean (a,c) and eFSSmean (b,d) for the organized Spring case (top) and the Summer flooding case (bottom). The standard physics configuration and the 99th percentile threshold are considered. The white dashed line at 0.5 represents the believable scale. Results were calculated over the area of the UK domain with radar coverage.

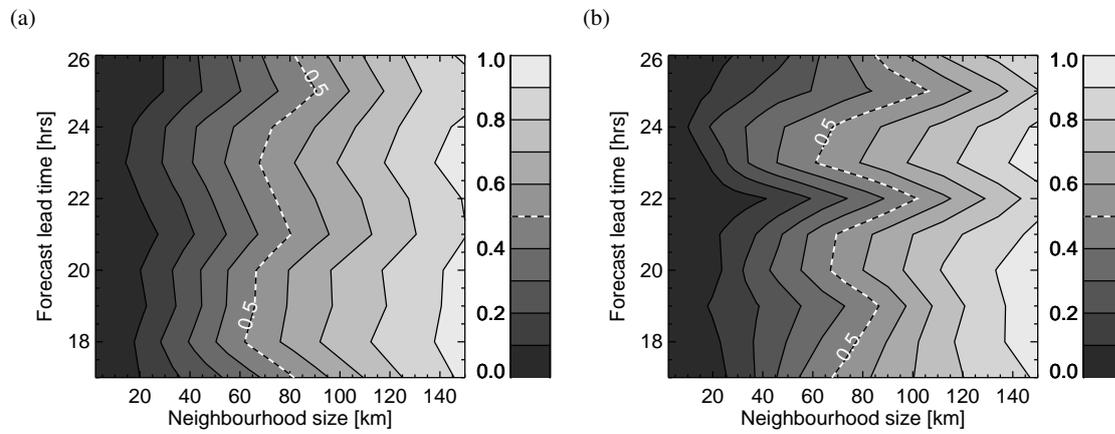


FIG. 6. FSS calculations over the Edinburgh subdomain: (a) dFSSmean and (b) eFSSmean. The 99th percentile threshold is considered. The white dashed line at 0.5 represents the believable scale.

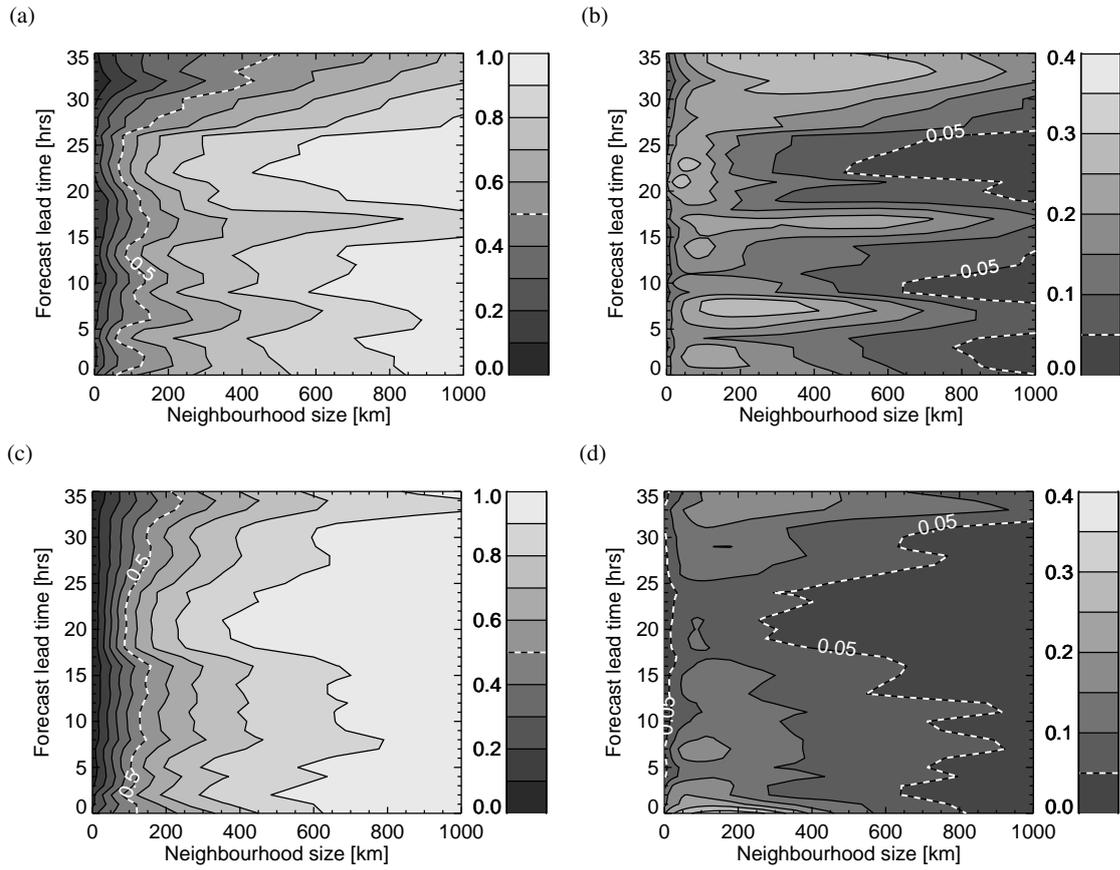


FIG. 7. dFSSmean (a,c) and dFSSstdev (b,d) for the organized Spring case (top) and the flooding case (bottom). The white dashed line in (a,c) at 0.5 represents the believable scale. To guide the eye, in (b,d) the white dashed line at 0.05 represents the neighborhood at which dFSSstdev is an order of magnitude smaller than the believable scale. The 99th percentile threshold is considered and results are calculated over the whole UK domain.

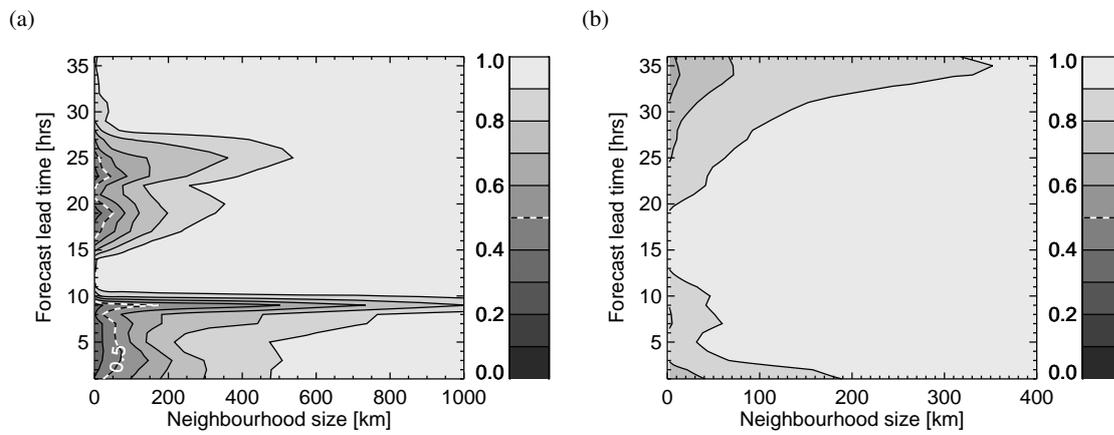


FIG. 8. Comparison of dFSSmean calculated for the (a) 99th and (b) 85th percentile thresholds for the 10 m horizontal wind speed field and the organized Spring case. Results are calculated over the whole of the UK domain and only the standard physics configuration is considered. The white dashed line at 0.5 represents the believable scale.

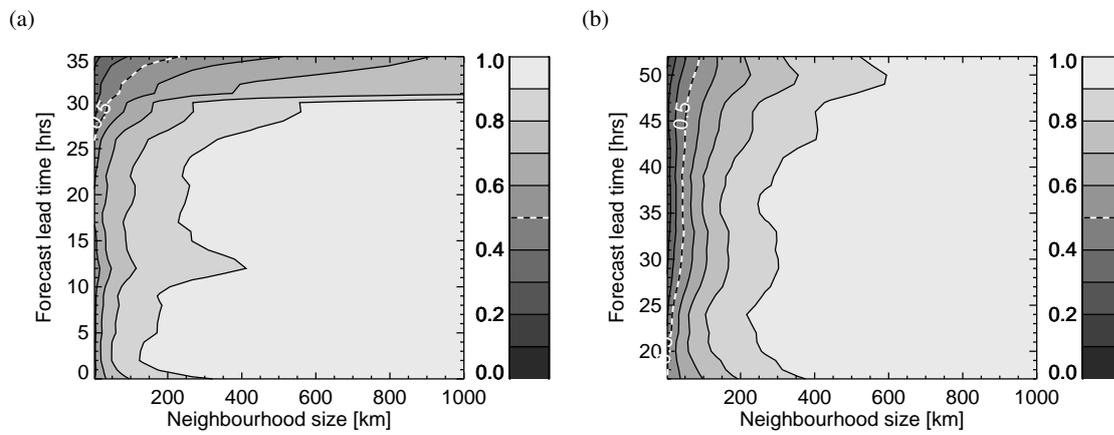


FIG. 9. dFSSmean calculated using the 90th percentile threshold of hourly precipitation accumulations for (a) the organized Spring and (b) the Summer flooding case. Results are calculated over the whole of the UK domain and only the standard physics configuration is considered. The white dashed line at 0.5 represents the believable scale.

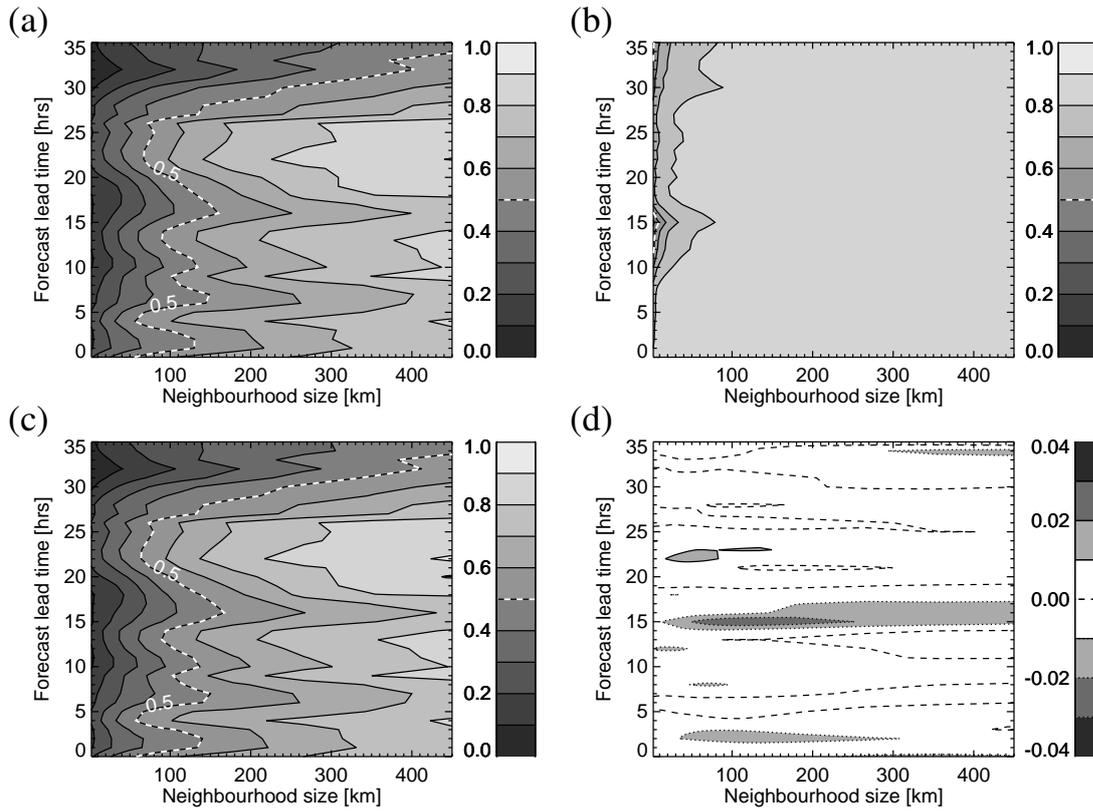


FIG. 10. dFSSmean comparisons of the restricted convection with increased time step and graupel physics configurations for the 99th percentile threshold of hourly precipitation accumulations. Results from different comparison strategies are shown: (a) MORGREPS2, (b) Physics2 and (c) MORGREPS2+Physics2. (d) shows the difference between sub-figures (c) and (a). Results are calculated over the whole of the UK domain. The white dashed line at 0.5 represents the believable scale.