



**A Methodology for Sensitivity Analysis of Spatial Features in Forecasts:
The Stochastic Kinetic Energy Backscatter Scheme**

Journal:	<i>Meteorological Applications</i>
Manuscript ID	MET-16-0165.R2
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	25-May-2018
Complete List of Authors:	Marzban, Caren; University of Washington, Applied Physics Lab; Dept of Statistics, Tardif, Robert; University of Washington, Atmospheric Sciences Hryniw, Natalia; University of Washington, Atmospheric Sciences Sandgathe, Scott; University of Washington, Applied Physics Laboratory
Keywords:	Sensitivity Analysis < Modelling, Statistical Models < Modelling, NWP < Modelling
Manuscript keywords:	

SCHOLARONE™
Manuscripts

A Methodology for Sensitivity Analysis of Spatial Features in Forecasts: The Stochastic Kinetic Energy Backscatter Scheme

Caren Marzban^{1,2}, Robert Tardif³, Natalia Hryniw³, Scott Sandgathe¹

¹ Applied Physics Laboratory,

² Department of Statistics,

³ Department of Atmospheric Sciences,

Univ. of Washington, Seattle, WA 98195 USA

Abstract

1 Stochastic Kinetic Energy Backscatter Schemes (SKEBS) are introduced in numerical
2 weather forecast models to represent uncertainties related to unresolved subgrid-scale
3 processes. These schemes are formulated using a set of parameters that must be
4 determined using physical knowledge and/or to obtain a desired outcome. Here, a
5 methodology is developed for assessing the effect of four factors on spatial features of
6 forecasts simulated by the SKEBS-enabled Weather Research and Forecasting (WRF)
7 model. The four factors include two physically motivated SKEBS parameters
8 (determining amplitude of perturbations applied to streamfunction and potential
9 temperature tendencies), a purely stochastic element (a seed used in generating random
10 perturbations), and a factor reflecting daily variability. A simple threshold-based
11 approach for identifying coherent objects within forecast fields is employed, and the
12 effect of the four factors on object features (e.g., number, size, and intensity) is assessed.
13 Four object types are examined: upper-air jet streaks, low-level jets, precipitation areas,
14 and frontal boundaries. The proposed method consists of a set of standard techniques in
15 experimental design, based on the analysis of variance, tailored to sensitivity analysis.
16 More specifically, a Latin Square Design is employed to reduce the number of model
17 simulations necessary for performing the sensitivity analysis. Fixed effects and random
18 effects models are employed to assess the main effects and the percentage of the total
19 variability explained by the four factors. It is found that the two SKEBS parameters do
20 not have an appreciable and/or statistically significant effect on any of the examined
21 object features.

22 Keywords: Sensitivity analysis, statistical models, parametrization, NWP, analysis of
23 variance.

24

25 1. Introduction

26

27 Stochastic Kinetic Energy Backscatter Schemes (SKEBS) are introduced in numerical
28 weather forecast models to enhance their skill in the production of probabilistic forecasts.
29 First introduced in Large Eddy Simulation models (e.g., Leith 1990; Mason and Thomson
30 1994), SKEBS are used to represent energetic contributions to flows from unresolved
31 physical processes through stochastic perturbations. For atmospheric flows, these
32 perturbations are added to model tendencies resulting in better calibrated forecast

33 ensembles (i.e., better match between mean errors and forecast uncertainties as
34 represented by the variance in the ensemble). Such schemes are formulated, as with any
35 other parameterization scheme, using a number of parameters that must be determined
36 based on physical knowledge and intuition, or tuned to obtain a desired outcome such as
37 increasing the variance in ensemble forecast members by a given amount. As such, it is of
38 interest to understand what effect the SKEBS parameter values have on the evolution of
39 simulated atmospheric states, especially if a specific effect is desired. For example, if
40 increased ensemble variance is the end goal, then it is useful to know which parameters to
41 vary to that end. Or, if one is performing object-oriented forecast verification (e.g.,
42 Gilleland *et al.*, 2009; Marzban *et al.*, 2009), then it is important to know how features of
43 the objects are affected by model parameters. All of these issues can be examined under
44 the umbrella of sensitivity analysis.

45
46 It is important to distinguish two distinct categories of Sensitivity Analysis (SA). In one
47 category SA is done primarily for the purpose of model tuning and/or data assimilation,
48 e.g., Ancell, Hakim (2007), Järvinen *et al.*, (2012), Laine *et al.*, (2012), and Ollinaho *et*
49 *al.*, (2014). In this category the SA is only a component of a complex optimization
50 problem where one seeks specific values of parameters (or initial conditions, etc.) that
51 optimize some quantity gauging the agreement between forecasts and observations.
52 Another way in which observations play a central role in this category of works is through
53 data assimilation. By contrast, in the second category, SA does not involve any
54 optimization or data assimilation (Alpert, 1993; Aires *et al.*, 2013; Marzban, 2013;
55 Marzban *et al.*, 2014; Marzban *et al.* 2018a; Marzban *et al.* 2018b; Yang *et al.*, 2014;
56 Dasari, Salgado, 2015; Smith, *et al.*, 2015); there the main purpose of SA is to assess the
57 effect of the parameters on the forecasts. The main goal is not to optimize forecasts but
58 rather gain knowledge on the relationship between model parameters and forecasts. This
59 knowledge may, in turn, be used for improving forecasts, or it may shed light on the
60 underlying physics of the phenomenon under study. There are (at least) two reasons that
61 render this latter approach to SA nontrivial: 1) The effect, on forecasts, of a given
62 parameter cannot be assessed independently of other parameters because the underlying
63 physics is inherently multivariate, and 2) natural variability must be taken into account in
64 order to establish the statistical significance of the results. Properly attending to these
65 issues is a complicated task that has led to a large body of literature on this flavor of SA
66 (Alpert 1993; Sobol', 1993; Oakley, O'Hagan, 2004; Fasso, 2006; Saltelli *et al.*, 2010;
67 Zhao, Tiede, 2011; Aires *et al.*, 2013; Marzban, 2013; Marzban *et al.*, 2014; Marzban *et*
68 *al.* 2018a; Marzban *et al.* 2018b). The present work falls into the latter category.

69
70 The approach adopted here consists of assessing the sensitivity of object features of
71 meteorological interest. Four object types are considered: upper-air jet streaks, low level
72 jets, precipitation areas, and frontal boundaries (i.e., baroclinic zones). Because the
73 SKEBS parameters affect the amount of energy that is injected into the flow, one expects
74 that large-scale features that rely on energetic growth (such as growing baroclinic modes)
75 would be affected by different parameter values.

76

77 Here the SKEBS in the Weather Research and Forecasting (WRF) model (Skamarock,
78 Klemp, 2008) is used. This SKEBS implementation introduces stochastic perturbations to
79 the simulated tendencies of potential temperature and non-divergent wind, which are
80 controlled through several user-specified parameters. Some of the parameters are
81 deterministic in nature, such as those used to control the amplitude of the perturbations,
82 which represent the total amount of backscattered energy in potential temperature and
83 non-divergent wind. However, since the perturbations are generated using an
84 Autoregressive process, there is also an element of pure randomness, hereafter referred to
85 as the purely stochastic component of SKEBS. This component is controlled by a seed
86 parameter that affects the random number generation in SKEBS. The reader is referred to
87 Berner *et al.* (2011) for more details on this SKEBS implementation. Here, the effect of
88 both types of parameters are evaluated and contrasted using WRF forecasts generated up
89 to 120 hours. There exist many more parameters in SKEBS whose impact on features of
90 objects is worthy of consideration. Here, the analysis is restricted to only two model
91 parameters in order to simplify the demonstration of the methodology. (In a work to
92 presented separately, as many as eight model parameters are being examined by the
93 authors).

94

95 Serving as the central piece in this evaluation are the four aforementioned object types
96 identified within gridded forecast fields. Section 4 describes a simple threshold-based
97 method for identifying the objects. In addition to the number of identified objects,
98 various quantities characterizing each object are recorded. For this study, these quantities
99 serve as the response variable in linear models, and methods of experimental design
100 provide a setting wherein the effect of several factors on these responses can be
101 quantified.

102

103 Two of the factors are key SKEBS parameters (the amplitude of perturbations to
104 rotational wind and potential temperature), and a third factor is the replication of SKEBS
105 itself (i.e., the seed used to generate sequences of random perturbations in SKEBS); this
106 factor represents the purely stochastic component of SKEBS. The fourth factor
107 represents the effect of daily variability. The third factor can be viewed as generating an
108 ensemble, and the fourth factor is motivated by the expectation that forecasts are sensitive
109 to initial conditions. The effect of these four factors is estimated for forecast hours 0-120.

110

111 As further explained in Section 2, the design of the experiment involves nine days, 41
112 forecasts at 3hr intervals (between 0 and 120 hours), nine values of each of the two
113 SKEBS parameters, and six SKEBS replications, which in a full factorial design leads to
114 a large number of experiments (or “ensemble members”); to reduce the number of
115 experiments, a special type of a fractional factorial design (called a Latin Square Design)
116 is used.

117

118 Experiments of this type are often called *computer experiments* because the resulting data
119 are not the result of a real experiment in any sense of the word (Sacks *et al.*, 1989; Welch
120 *et al.*, 1992; Santner, *et al.*, 2003; Fang, *et al.*, 2006). The defining characteristic of
121 computer experiments is that the experimental error is zero, because re-running the

122 computer model (here WRF/SKEBS) leads to the same set of outcomes. Without an
 123 estimate of experimental error it is impossible to perform any of the statistical tests
 124 designed to assess statistical significance (Santner *et al.*, 2003; Fang *et al.*, 2006).
 125 However, as long as one is interested in main effects only (i.e., no higher-order
 126 interactions), then standard methods of experimental design can be used for assessing
 127 statistical significance, because all of the contributions to variance from higher-order
 128 interactions can act as a proxy for experimental error (Montgomery, 2009).

129

130 2. Experimental Design: A Brief Introduction

131

132 This study aims to determine how certain spatial features of forecasts are affected by four
 133 factors, including two model parameters: the amplitude of perturbations to 1) rotational
 134 wind, and 2) potential temperature, denoted Par1 and Par2, respectively. Additionally,
 135 another factor is also examined – one that measures how the effects vary across (here, 9)
 136 days; it is denoted Day. One important question is: How does the effect of the
 137 deterministic parameters (Par1 and Par2) compare with the effect of the purely stochastic
 138 component of SKEBS? Therefore, in addition to the three factors Day, Par1, and Par2, a
 139 fourth factor - denoted Rep - is introduced to measure the effect of replicating the
 140 experiment. Finally, it is useful to examine how all these effects vary with forecast (valid)
 141 time, denoted Fhour (here, varying from 0 to 120 hours).

142

143 In the field of experimental design (Montgomery, 2009), linear models are often
 144 employed to estimate the effect of various factors on the response. One simple model is

145

$$146 \quad y_{ijkl} = \mu + Day_i + Par1_j + Par2_k + Rep_l + \varepsilon_{ijkl}, \quad (1)$$

147

148 where the response y_{ijkl} denotes a measurement of some quantity of interest (e.g., the
 149 number of jet streaks) on the i^{th} Day, for the j^{th} and k^{th} values of Par1 and Par2,
 150 respectively, and for the l^{th} replication of the experiment. The factor Fhour is not
 151 included in the model, because the model is developed at each value of Fhour. The terms
 152 appearing on the right side of Eq. (1) are all parameters (not to be confused with SKEBS
 153 parameters) to be estimated from data on the response y and the factors. The ε term is a
 154 random variable whose variance σ_ε^2 is another quantity that must be estimated from data,
 155 not only for assessing goodness-of-fit, but also for performing statistical tests. It can be
 156 shown (Montgomery, 2009) that the least-squares estimates of these parameters generally
 157 involve sample means of the response, or the difference between two sample means. For
 158 example, the least-squares estimate of the μ parameter is the sample mean y_{\dots} , also called
 159 the *grand mean*. The parameter Day₁ is estimated by the difference ($y_{1\dots} - y_{\dots}$). In all of
 160 these expressions a “dot” refers to a sample mean over the corresponding index. The
 161 other components - the Day factor, and the other factors in the model - are all estimated
 162 through similar difference between sample means. Given that the estimates of the factors
 163 are differences from the grand mean, these estimates are also called *main effects*. The
 164 machinery of experimental design aims to perform statistical/hypothesis tests of whether
 165 the true/population main effects are zero; (see next paragraph for another measure of a
 166 factor's effect.) The model in Eq. (1) is strictly linear, but it is possible to introduce

167 nonlinear terms. Such terms generally appear as terms with multiple indices, and they are
 168 called *interaction effects*. For example, a term like X_{ij} (called a 2-way interaction)
 169 measures how the effect of Par1 on response varies across days.

170
 171 Although tests of main effects are performed for the problem at hand, there exists an
 172 alternative approach which is also appropriate. Strictly speaking, the main effects
 173 discussed above are estimates of fixed, population parameters, and for this reason they are
 174 called *fixed effects*. Any conclusions based on a fixed effects model are specific only to
 175 the particular values assigned to the various factors. However, one may choose to view
 176 these particular values as a random sample taken from a larger space of parameter values,
 177 in which case it makes no sense to speak of the main effect of a factor, because any
 178 notion of an effect is itself a random variable. Effects of this type are called *random*
 179 *effects*, and any conclusions based on a random effects model pertain to the population of
 180 **all** possible values that the factors may take, not the specific values appearing in the
 181 sample only. In such models, the main aim is not to test whether or not an effect is zero,
 182 but rather to test whether or not any portion of the variability in the response can be
 183 explained by each of the factors in the model. Specifically, for random effects models
 184 one writes

$$185 \sigma_{Response}^2 = \sigma_{Day}^2 + \sigma_{Par1}^2 + \sigma_{Par2}^2 + \sigma_{Rep}^2 + \sigma_{\epsilon}^2, \quad (2)$$

186 and the goal is to estimate and then test whether any of the *variance components* on the
 187 right hand side of Eq. (2) are zero.

188
 189
 190 To clarify the difference between a fixed effects model and a random effects model,
 191 suppose the Day factor takes d values (i.e., the number of days in the study). Treating the
 192 Day factor as a fixed factor would allow one to test whether there is a difference between
 193 the sample means of the response across the d days. A significant result would then
 194 suggest that the mean response varies across the specific d days, i.e., the Day factor has an
 195 effect on the response. However, one may choose to consider the d days in the study as a
 196 random sample taken from the population of all days, in which case it is more appropriate
 197 to treat the Day factor as a random factor. Then, one can test the null hypothesis $\sigma_{Day}^2 = 0$
 198 which constitutes a test of whether a nonzero portion of the total variability in the
 199 response $\sigma_{Response}^2$ can be accounted for by daily variability. A significant result would
 200 suggest that the mean response varies across **all** days (not just the d days appearing in the
 201 data). Similarly, one can treat Rep, Par1, and Par2 as fixed or random factors. Although
 202 fixed effects models provide intuitive measures of effects, random effects models have
 203 the advantage that the final conclusions are not specific to the values of the factors chosen
 204 for the study. As such, both model types are useful.

205
 206 Therefore, here, both types of models are developed. First, the factors are treated as fixed
 207 parameters. The estimate of each factor represents the sensitivity of the response with
 208 respect to that factor, i.e., the main effect of that factor. Then, random effects models are
 209 developed wherein the sensitivity of the response with respect to a given factor is
 210 measured by the variance component of that factor. It is more useful to report the variance

211 component as the fraction of the total variance. For example, the sensitivity for the Day
 212 factor is best reported as the so-called *intraclass correlation*

$$\rho_{Day} = 100 \frac{\sigma_{DAY}^2}{\sigma_{Response}^2} ; \quad (3)$$

213
 214

215 similarly for the other variance components. Another advantage of examining the
 216 intraclass correlation is that analytic formulas exist for its confidence intervals
 217 (Montgomery, 2009). Such confidence intervals are critical for assessing the statistical
 218 significance of the sensitivity results.

219

220 In a full factorial design involving the four factors Day, Par1, Par2, and Rep, the number
 221 of model runs would be equal to the product of the number of values of each factor. That
 222 number of runs is often impractically large, and so, there exist a number of experimental
 223 designs whose goal is to reduce the number of runs. The Latin Square Design (LSD) is
 224 one such design, and it is briefly explained in the Appendix. In order to illustrate the basic
 225 idea, consider a problem involving three factors (and a response), with each factor taking
 226 three possible values. Ideally, one must observe the response at all 3^3 possible values of
 227 the factors, because then one can estimate the effect of the three factors as well as all of
 228 the interactions between them. However, it can be shown (Montgomery, 2009) that 3^2
 229 runs are sufficient for estimating the main effects of the factors, if the values of the 3
 230 factors for the 3^2 runs are selected according to a special prescription best displayed as a
 231 square table. An example of such a square is shown in Table 1, where the factors are
 232 denoted A, B, and C, and the subscripts denote the value of each factor. For example, the
 233 bottom/right element in that square corresponds to a run where the factors A and B are set
 234 to their third value, and the factor C is set to its second value. If the three factors have p
 235 levels, then the square table is $p \times p$, and so, the necessary number of runs is only p^2 . This
 236 example involves three factors, but it can be shown that the number of necessary runs is
 237 p^2 regardless of the number of factors (See Appendix). Such tables are called Latin
 238 Squares, and by virtue of being square tables, designs that follow such tables dramatically
 239 reduce the number of necessary runs, although at the cost of making all interactions
 240 between the factors inestimable (Montgomery, 2009). The inability of the LSD to
 241 estimate interaction effects is not a major concern because the main effects are generally
 242 much larger than interaction effects. The expectation that higher-order interactions are
 243 weaker than main effects is generally borne out due to several principles: the principle of
 244 hierarchical ordering, the principle of effect sparsity, and the principle of effect hierarchy;
 245 see pages 192, 230, 272, 314, 329 in (Montgomery 2009), and pages 33-34 in (Li,
 246 Sudarsanam, and Frey 2006). In the case of precipitation, Marzban et al. (2014) also find
 247 the interactions to be much smaller than main effects.

248

249 3. Data

250

251 Version 3.7.0 of the WRF-ARW model was used for this work, with lateral boundary
 252 conditions specified every 6 hours from output of the Global Forecast System (GFS). All
 253 of the standard WRF parameters were the default “out of box” parameters, with a 25-Km

254 grid-spacing for a domain 200 (east-west) by 140 (north-south), over the Continental US.
255 Nine days are selected between December 2014 and March 2015. Each initial forecast
256 hour is 10 days apart in this time period, ensuring minimal temporal association between
257 days. The specific dates are as follows: Dec. 01, 11, 21, 31, Jan. 10, 20, 30, and Feb. 9,
258 and 19. Winter months were chosen for the high degree of variability with regards to jet
259 stream activity and mid-latitude cyclone activity.

260

261 For this study three factors Day, Par1, and Par2 were sampled according to the LSD,
262 thereby reducing the necessary number of runs from 9^3 to 9^2 . As a result, it is assumed
263 that the interactions between these three factors are much smaller than the main effects.
264 Because of the LSD, Par1 and Par2 take nine values as well. The range of the nine values
265 are chosen to be centered on the recommended SKEBS values, but in order to examine
266 the full range of possible effects, they span one order of magnitude smaller and one order
267 of magnitude larger than the default values. The nine specific values are (0.1, 0.325,
268 0.550, 0.775, 1.000, 3.250, 5.500, 7.750, 10.000) $\times 10^{-5}$ for Par1, and (0.1, 0.325, 0.550,
269 0.775, 1.000, 3.250, 5.500, 7.750, 10.000) $\times 10^{-6}$ for Par2. As mentioned previously, in
270 the random effects model inference of the sensitivities pertains to all possible values of
271 the parameters, not just to the specific nine values; for this reason, the specific nine
272 values selected here do not play an important role in the final analysis. Indeed, in an
273 earlier version of the analysis, the following Par1 values produced very similar results:
274 (0.5, 2.875, 5.25, 7.625, 1.0, 12.375, 14.75, 17.125, 19.5) $\times 10^{-5}$.

275

276 One of the main goals here is to assess the effect of Rep (i.e., the purely stochastic
277 component of SKEBS) and how it compares with the effect of the other factors.
278 Therefore, more computational effort is dedicated to that end. Specifically, the Rep factor
279 and the combination of the other three factors (Day, Par1, and Par2) follow a full factorial
280 design. In other words, all 9^2 LSD runs involving Day, Par1, and Par2, are replicated at
281 every value of the Rep factor. Here the number of replications is six.

282

283 All of these factors are treated first as fixed factors, and then as random factors.
284 However, the factor Fhour is treated as a fixed factor, because it varies across 41 fixed
285 values, from 0 to 120 hours, in increments of 3 hours. As mentioned above, Fhour is not
286 included in the model, because the model in Eq. (1) is developed at each of the 41 values
287 of Fhour. Consequently, all of the results found here take the form of “time series” of the
288 main effects, variance components, or intraclass correlations as a function of Fhour.

289

290 Given the above design, the total number of runs is $9^2 \times 6 \times 41 = 19,926$. Although this
291 is a large number of runs, it is significantly smaller than what would be necessary in a full
292 factorial design: $9^3 \times 6 \times 41 = 179,334$.

293

294 **4. The Response**

295

296 As previously mentioned, sensitivity to SKEBS parameters is assessed with respect to
297 various features of the following meteorologically significant object types: 1) upper-air jet
298 streaks, 2) low-level jets, 3) precipitation events, and 4) frontal boundaries. Jet streaks

299 are defined at 250 *hPa* as regions with contiguous model grid points having wind speeds
300 in excess of 50 m s^{-1} (approximately 100 *knots*). Similarly, low level jets are defined as
301 regions at 850 *hPa* characterized by winds stronger than 20 m s^{-1} . Precipitation events are
302 contiguous regions where the total precipitation accumulation at the surface is above 1*mm*
303 in 3-hour forecast intervals. Frontal boundaries (i.e., baroclinic zones) are identified using
304 the horizontal gradients of the 1000-700 *hPa* geopotential thickness field (McCann,
305 Whistler, 2001). A threshold value of 0.3 m Km^{-1} for the magnitude of the geopotential
306 height gradient is used to identify significant baroclinic zones. Although it is possible to
307 develop more sophisticated means of identifying such objects, the focus of this study is
308 on the development of an object-based SA method, regardless of how the objects are
309 identified.

310

311 By definition, all of these objects are characterized by relatively well-defined spatial
312 extent. For every available forecast hour, objects meeting the above criteria are
313 identified. Figure 1 shows an example of jet streak objects identified in the 250 *hPa*
314 WRF wind field. Three jet streak objects are identified in this particular WRF forecast.
315 The smallest and weakest is located over states in the northwestern United States, a
316 second is located over eastern Canada, while the largest and strongest extends from the
317 southwestern Four Corners states to the mid-Atlantic states.

318

319 Here we point out that the five grid points nearest the model domain's lateral boundaries
320 are omitted from the analysis in order to prevent any direct influence from the imposed
321 lateral boundary conditions taken from deterministic GFS forecasts. This way, only grid
322 points in the interior of the domain where solutions are fully influenced by SKEBS
323 perturbations are considered. Also, identified objects are restricted to those composed of
324 at least 50 grid cells, corresponding to areas larger than about $31,000 \text{ Km}^2$ in order to
325 minimize any “noise” in the resulting object datasets that could be associated with
326 spurious appearance/disappearance of small areas with wind speeds changing to values
327 just above/below the threshold. Despite the application of such conditions, spurious
328 changes in object characteristics may occur as objects merge or separate solely due to
329 subtle changes in the underlying continuous field; for example, it is possible for two
330 nearby jet streak objects at a particular forecast hour to merge at the following forecast
331 hour due to an increase in wind speed above the threshold in the region separating the two
332 jets. Associated changes to the response variable (e.g., number of objects or their size,
333 intensity, and location) can be described as “measurement error” because the variability
334 introduced by these changes is not due to any of the factors included in the model (Eq. 1).

335

336

337 The features examined here are 1) the number of objects, 2) their size, 3) intensity, and 4)
338 location. The size of each object is computed as the number of grid points included in that
339 object. The intensity is measured as the mean intensity of the field across the object, and
340 their location is recorded at the latitude and longitude of the center-of-mass of the object.
341 Panels a-c in Figure 2 show the histograms of number, size, and intensity for precipitation
342 objects across all four factors (i.e., days, replications, Par1 and Par2); other object types
343 have similar histograms. The histograms of latitude and longitude are not shown because

344 that figure shows no useful information. It can be seen that the number of objects can
345 vary between 1 and 13, with the most common value around 3 or 4. By contrast, the size
346 of objects has an exponential-looking histogram, and so the data examined consists of
347 mostly small objects (i.e., consisting of 50 grid points). Mean intensity values (panel c)
348 vary between about 1 and 17 *m/s*, with the most common value around 2.3 *m/s*.

349

350 Given the similarity in the shape of the histograms of number and size of objects, one
351 may wonder if these two features are correlated. In fact, given that the size of the forecast
352 domain is fixed, one may suspect a negative correlation. Panel d in Figure 2 shows the
353 scatterplot of these two features. Although for the extreme case where there are as many
354 as 13 objects, their size is restricted to be around 500 grid points, for cases with four
355 objects, their size can vary from the smallest possible value (50) to 3500 grid points. As
356 such, it can be seen that there is no linear association between the two features.

357

358 The histograms discussed above are constructed from the object features that arise in the
359 data across all values of the four factors. But even for given values of the four factors,
360 there exists a distribution of features. Here, that distribution is summarized by two
361 quantities - the minimum and maximum; (the 25th and 75th percentiles of the histograms
362 were also examined, but the results were statistically equivalent to those based on the
363 minimum and maximum) In short, the aim is to study the effect of the aforementioned
364 four factors on the following response/feature variables: 1) Number of objects (e.g., jet
365 streaks) across the forecast domain, 2) Minimum and 3) Maximum size of (i.e., smallest
366 and largest) objects across the domain, and 4) Minimum and 5) Maximum intensity (i.e.,
367 weakest and strongest) of objects across the domain. As for the location feature, the
368 minimum, median, and maximum of latitude and longitude are also examined; it can be
369 argued that the two SKEBS parameters considered here may have an effect on the
370 location of the objects because they control propagation and development rates. (The
371 authors acknowledge an anonymous Reviewer for this suggestion).

372

373 Although all four object types (upper-air jet streaks, low-level jets, precipitation events,
374 and frontal boundaries) have been analyzed, only sensitivity results pertaining to upper-
375 air jet streaks are shown in the next section. Results with respect to the other object types
376 were found to be similar especially in terms of the relative magnitude of the effect of the
377 four factors. Of the various features considered here, specific results pertaining to
378 latitude and longitude are not shown, because they are similar to those pertaining to the
379 intensity feature.

380

381 **5. Results**

382

383 Before developing the aforementioned models, it is useful to examine the simulated data,
384 first. Figure 3 shows the values of the five responses/features as a function of forecast
385 time (Fhour), on one day, with model parameters set to default values, and for the six
386 replications (in colours). The thick/black line corresponds to a run wherein all of SKEBS
387 has been turned off. It can be seen that the coloured curves (i.e., different replications of
388 SKEBS with default parameters) generally fluctuate about the curve of this control run.

389 Moreover, evidently, all five response variables have significant variability across
390 forecast times. Part of this variability is “real” in the sense that objects can appear and
391 disappear in a forecast field across three hours. The remainder of the variability is due to
392 the aforementioned measurement error; for example, although the actual size of an object
393 may not change in a 3-hour interval, the thresholding procedure adopted here for
394 identifying objects may give a slightly different value for the size. This measurement
395 error is not a stumbling block for the analysis; its only effect is to magnify the variance of
396 the ϵ term in Eq. (1), and thereby reduce statistical power. Also, as mentioned at the end
397 of Section 1, these empirical errors are necessary for performing statistical tests of
398 significance in computer experiments.

399

400 The variability of the response variable plays an important role in both fixed effects and
401 random effects models. Figure 4 shows the variability of the five response variables at
402 each forecast time. The slow modulations of all of these curves correspond to the natural
403 evolution of weather patterns in the nine days examined here. To obtain a sense of the
404 variability of these results, 95% confidence intervals are also shown (as vertical bars). It
405 is evident that all five response variables have nonzero variance at all forecast hours.
406 Recall that the goal of random effects models is to determine how these variances are
407 apportioned across the various factors in the model.

408

409 The linear model in Eq. (1) is developed at each forecast time. Treating the factors as
410 fixed factors allows one to perform F (or t) tests on the main effects. The resulting p-
411 values are summarized in Figure 5. The variability in the boxplots is across the 120
412 forecast hours. Here, a significance level (e.g., 0.05 or 0.01) is not selected to assess
413 statistical significance. Instead, the boxplot of the p-values is examined to provide a
414 visual assessment of the “strength” of the statistical significance. A tight boxplot, near 0,
415 suggests that the corresponding effect is statistically significant. By contrast, if the
416 boxplot of p-values is near 1 or extends across the full range from 0 to 1, then the
417 corresponding factor is deemed non-significant, i.e., that there is insufficient evidence
418 from data to conclude that the factor has an effect. This practice is consistent with a
419 fundamental theorem in statistics stating that the distribution of p-values is given by a
420 uniform distribution between 0 and 1, if the null hypothesis (of no-effect) is true.

421

422 Here (Figure 5) it can be seen that the factor Day has a significant effect on all five
423 responses. This is not surprising, because it is known that the responses vary across the
424 nine days in the study. By contrast, the near-1 location of the boxplots for Par1 and Par2
425 in all five panels suggests that there is no evidence from data to suggest that these two
426 parameters have any effect on any of the response variables. The Rep factor plays a more
427 complex role; although the p-values do extend to relatively large values, the bulk of their
428 histogram is skewed toward smaller values, in all five panels. In other words, the Rep
429 factor does appear to have an effect on all five response variables, but not at all forecast
430 hours.

431

432 Although it is possible to examine the p-values in the fixed-effects model at each forecast
433 hour, it is more useful to examine the forecast-hour-dependence of the results in the

434 random effects model. Treating the factors as random variables leads to consideration of
435 the variance components, and in turn, intraclass correlations, ρ in Eq. (3), and their
436 confidence intervals. Figure 6 shows the 95% confidence interval for ρ at different
437 forecast times. The “Day, Number” panel shows the effect of the Day factor on the
438 number of objects in the domain. It can be seen that the effect of the Day factor
439 diminishes very quickly, and falls to near-zero values for Fhour beyond nine. The effect
440 of the Day factor on the size of objects is shown in the panel marked “Day, Size”;
441 although the smallest (black) objects are mostly unaffected by the Day factor, the effect
442 on the largest (red) objects is less trivial. On forecast times scales from 0 to 120 hours,
443 for very short forecast times (3 to 9 hours) the Day factor can explain 60% to 90% of the
444 variability; even for longer forecasts, the effect is non-zero, leveling-off at values in the
445 5% to 10% range. In other words, even for very long forecast times, daily variability
446 contributes a significant portion of the total variability in the size of objects. The effect of
447 the Day factor on the (mean) intensity of objects has a similar behavior (panel “Day,
448 Intensity”), although for longer forecast times, the effect is generally weaker than the
449 effect on object size. Said differently, for short forecast hours the variability in object
450 intensity can be explained by daily variability, but for very long forecast times that
451 variability is not due to daily changes in weather.

452

453 The effect of the Rep factor can be seen in the second row of panels in Figure 6. For all
454 five response variables (number, min. size, max. size, min. intensity, and max. intensity),
455 Rep can explain only about 0.1% to 0.5% of the variability in the data. The large
456 confidence intervals make it difficult to interpret the results; the lower end of the intervals
457 are generally above zero, suggesting that the corresponding effects are nonzero, consistent
458 with the small p-values observed in Figure 5. Although the top end of the intervals is
459 erratic, it is important to note the scale on the y-axis of these panels - 0 to 1% - and so,
460 the effect of Rep is generally quite small.

461

462 The effect of the parameters (Par1, Par2) on all response variables is even weaker than
463 that of the Rep factor (third and fourth rows in Figure 6). The ρ values are generally
464 below 0.1%. In other words, even when the effect of the parameters is statistically
465 significant (i.e., nonzero at 95% confidence level), the magnitude of the effects is
466 extremely small. The fact that the effect of the parameters is weaker than that of Rep is
467 important, and is further discussed in the next section.

468

469 The last row of panels in Figure 6 shows ρ_ϵ , i.e., the percentage of the variability in the
470 data that cannot be explained by the four factors Day, Rep, Par1, and Par2. As such, it is
471 useful for assessing the combined effect of the four factors. Evidently, for forecast hours
472 longer than three hours nearly 100% of the variability in the number of objects cannot be
473 explained by any of the four factors. This is expected from the panels in the first column
474 of Figure 6, because none of the four factors appear to have an effect on the number of
475 objects for long forecast times.

476

477 When the response is object size (bottom row, middle panel), or object intensity (bottom
478 row, right panel), the variability that cannot be explained by the four factors generally

479 increases with forecast time. For the smallest of objects (black curve) the increase is quite
480 abrupt – from 0 to 100% as one goes from 0hr to 3hr forecasts and beyond. For the largest
481 objects (red curve), although the increase is more gradual, the percentage of unexplained
482 variance approaches 100% by forecast hour 100. The undulations in the curves, caused by
483 the natural variability in the data across the 120 hours, make it difficult to pinpoint a
484 specific forecast time beyond which the four factors become useless.

485 In summary, examining all of the panels in Figure 6, it appears that when the factors do
486 contribute to the variability in the response, most of that variability is due to the Day
487 factor. The next important factor is Rep; and Par1 and Par2 have nearly no effect. It is
488 also clear that Par1 and Par2 have a much smaller effect than Rep, at every forecast hour.
489 This suggests that the two tunable SKEBS parameters examined here may not produce
490 the expected variability in the specific objects under consideration, since the purely
491 stochastic component (which is not as controllable as the tunable parameters)
492 overwhelms the variability in the forecasts.

493

494 **6. Conclusion and Discussion**

495

496 SKEBS has been designed to introduce variability into the forecasts in a manner
497 consistent with the physics that are unresolved by the model. One would then expect that
498 SKEBS parameters (Par1, Par2) would have some effect on the forecasts, and that the
499 effect of these parameters would be more prominent than that of the purely stochastic
500 component of SKEBS (Rep). Here, forecasts of jet streaks, low-level jets, precipitation,
501 and baroclinic zones are considered, although only the analysis on jet streaks is presented.
502 A simple method is employed to identify these objects within continuous forecast fields; a
503 suite of methods from experimental design are then woven together to assess the effect of
504 four factors (Day, Par1, Par2, Rep) on five features of these objects (number, minimum
505 and maximum size, and minimum and maximum intensity. The impact of the four factors
506 on the location (latitude and longitude) of the objects is also examined; but it is not
507 presented because the results are similar to that of intensity. It is shown that the number
508 of objects in these fields does not appear to be affected by any of the factors. It is also
509 shown that for forecast times when the factors do have a nonzero effect on the size and
510 intensity of objects, apart from the effect of the Day factor, the effects of the other three
511 factors are quite small, explaining only a few percentage points of the variability observed
512 in the data. More importantly, it is found that the effect of Par1 and Par2 is much less
513 than that of Rep.

514

515 This suggests that the variability produced by varying the two SKEBS parameters does
516 not appear to have a significant effect on the specific object types and their features
517 examined here; the purely stochastic part is the main driver of any SKEBS-induced
518 variability. It is important to emphasise that this conclusion pertains only to the specific
519 object types and features examined here. It does not reflect on the connection between
520 SKEBS and the physical processes it seeks to represent, and whether the physically-
521 motivated model behind SKEBS has a consistent effect on model forecast evolution at
522 large. In practice, then, if one is interested in the specific objects and features examined

523 here, it is best if the resources for tuning or calibrating the model parameters are directed
524 away from the physical SKEBS parameters. However, see next paragraph.

525

526 Armed with the methodology developed here, the above analysis can be generalized in a
527 number of ways. For instance, the criteria for identifying objects can be revised; the
528 number of parameters, and their range and values can be extended, and/or other response
529 variables can be examined. Although the two SKEBS parameters under consideration do
530 not appear to have an effect on the four object types examined here, it will be useful to
531 find other meteorologically relevant objects that are affected by these SKEBS
532 parameters. As pointed out by an anonymous reviewer, it is known that the SKEBS
533 parameters examined here do affect the reliability/skill of large-scale ensemble forecasts.
534 As such, the null effect of the model parameters may seem contradictory, but then it is
535 important to recall that the sizes of the objects considered here fall on the smaller end of
536 the resolved scales in the model simulations.

537

538 One may also consider more/other SKEBS parameters, in which case Graeco-Latin
539 Square Designs (GLSD; see Appendix) can be used to reduce the number of runs
540 necessary for estimating main effects. A desirable feature of GLSD is that the necessary
541 number of runs for estimating main effects is the square of the number of values each
542 factor takes, independent of the number of factors in the study. In fact, fixed-effects and
543 random-effects models with as many as eight parameters are currently under
544 investigation, and preliminary results suggest that even when some of the SKEBS
545 parameters do affect the spatial structure of the forecasts, their effect is still overwhelmed
546 by daily variability and variability due to replication. That finding also raises the
547 possibility of examining the effect of the factors on the spatial structure of the forecasts,
548 independently of the existence of any objects in the forecast field. Generalizations can
549 also be made to the statistical modeling effort. For example, the fixed-effects and
550 random-effects models employed here are linear models commonly employed in
551 experimental design (Montgomery 2009). These can be generalized to include higher
552 order interactions. Alternatively, it is possible to replace these models entirely with fully
553 nonlinear models - often called metamodels (Santner *et al.*, 2003; Aires, 2013). Many of
554 these questions are currently under consideration.

555

556 A comparison of the current work and that reported in Marzban *et al.* 2018b is in order.
557 First, and foremost, whereas the objects here are identified by a simple thresholding
558 method, those in the latter work are identified via two different clustering algorithms.
559 Second, the (11) model parameters in the latter study are continuous parameters which
560 necessitates a different (than LSD) method for sampling the parameter space. The reason
561 the model parameters are different between the two studies is that the underlying model
562 in the latter work is COAMPS® (Coupled Ocean/Atmosphere Mesoscale Prediction
563 System). The impact of the 11 parameters in COAMPS on the spatial structure of
564 forecasts (i.e., without reference to any objects) has also been examined (Marzban *et al.*
565 2018a).

566

567 **7. Appendix: Latin Square Designs**

568

569 Consider an experiment involving three factors, A , B , and C , each taking three values
570 denoted $A1, A2, A3, B1, B2, B3$, and $C1, C2, C3$. (In statistics, the values a discrete
571 variable can take are referred to as levels. Here we avoid the term level in order to
572 minimize confusion with the use of that term in meteorology.) A full factorial design
573 refers to 3^3 runs necessary to consider all possible combinations of the values each factor
574 can take. It can be shown (Montgomery, 2009) that in a full factorial design one can
575 estimate all main effects, all interactions, and the variance of the errors, σ^2_ϵ . If, however,
576 interactions are not of interest, then only the specific runs shown in Table 1 are sufficient.
577 In other words, only the nine runs $(A1, B1, C1), (A1, B2, C2), (A1, B3, C3), \dots, (A3, B3,$
578 $C3)$ are sufficient for estimating the main effects (and the error variance). An experiment
579 involving only such specific runs is said to follow a Latin Square Design (LSD).
580 Interactions, however, cannot be estimated. Technically, in LSD, main effects and
581 interactions effects are said to be *aliased*, meaning that the effects one can estimate are a
582 combination of main effects and interaction effects, and one cannot disentangle the two.
583 As such, when one computes main effects in an LSD, the assumption is that the
584 interaction effects are negligible. Latin squares as in Table 1 are constructed by assigning
585 the columns to the values of one factor, the rows to the values of another factor, and then
586 cyclicly permuting the values of the last factor within the body of the square. This assures
587 that every combination of the three values appears precisely one time - a unique and
588 defining characteristic of the LSD. The factors may take more than three values, in which
589 case the Latin square will simply be larger.

590

591 The Graeco-Latin Square Design (GLSD) is the generalization of the LSD to four or more
592 factors, with each factor taking any number of values; the only constraint is that all
593 factors must have the same number of values. So, in the present study, if nine days are
594 selected for the analysis, then each of the two parameters (Par1 and Par2) must take nine
595 values. More examples of LSDs and GLSDs can be found in Montgomery (2009).

596

597 It is worth mentioning that in an LSD involving the three factors Day, Par1 and Par2, on
598 no single day are the two parameters varied across all their values. Consequently, one
599 cannot assess the sensitivity of the two parameters for each day. This may appear to be a
600 limitation; however, it is important to point out that knowledge of sensitivities for any
601 given day is useless; only the sensitivities across all days have practical utility. And the
602 LSD allows one to estimate those sensitivities with only 9^2 runs (instead of 9^3).

603

604 It is important to distinguish LSDs (or GLSDs) with another sampling design with a
605 similar name, namely Latin Hypercube Sampling (LHS). Although frequently used in SA
606 (Hacker *et al.*, 2011; Marzban, 2013; Marzban *et al.*, 2014), the LHS is a completely
607 different sampling scheme, and is most suitable for situations where the covariates
608 (independent variables) are continuous quantities, not discrete factors; there, one specifies
609 the desired sample size, first. Then, each of the factors is subdivided into that many bins,
610 and a sample is drawn such that any combination of the bins appears precisely one time.
611 The utility of the LHS derives from the fact that LHS estimates of model parameters are
612 more precise (at least, no less precise) than estimates based on Simple Random Sampling

613 (McKay *et al.*, 1979). Note that by contrast to the LSD (or GLSD) where the sample size
614 is simply the square of the number of values in a factor, the sample size in LHS is not
615 determined by the number of values of a factor, or the number of factors; instead, it is
616 specified by the user.

617

618 **8. Acknowledgments**

619

620 This work has received support from Office of Naval Research (N00014-12-G-0078 task
621 29) and National Science Foundation (AGS-1402895).

622

623 **9. References**

624

625 Aires F, Gentine P, Findell K, Lintner B, Kerr C. 2013. Neural Network-Based
626 Sensitivity Analysis of Summertime Convection over the Continental United States. *J.*
627 *Climate*, **27**, 1958-1979.

628

629 Ancell B, and Hakim G. 2007. Comparing Adjoint- and Ensemble-Sensitivity Analysis
630 with Applications to Observation Targeting. *Mon. Weather. Rev.*, **135**, 4117-4134.

631

632 Berner J, Ha SY, Hacker JP, Fournier A, Snyder C. 2011. Model uncertainty in a
633 mesoscale ensemble prediction system: Stochastic versus multiphysics representations.
634 *Mon. Weather. Rev.*, **139**, 1972-1995.

635

636 Dasari HP, Salgado R. 2015. Numerical modeling of heavy rainfall event over Madeira
637 Island in Portugal: sensitivity to different micro physical processes. *Meteorol. Appl.*, **22**,
638 113-127.

639

640 Gilleland E, Ahijevych D, Brown BG, Casati B, Ebert E. 2009. Inter comparison of
641 spatial forecast verification methods. *Wes. Forecasting*, **24**, 1416-1430.

642

643 Fasso A. 2006. Sensitivity Analysis for Environmental Models and Monitoring Networks.
644 In: Voinov A, Jakeman AJ, Rizzoli, AE (eds). Proceedings of the iEMSs Third Biennial
645 Meeting: Summit on Environmental Modeling and Software. International Environmental
646 Modeling and Software Society, Burlington, USA, July 2006.
647 Internet: <http://www.iemss.org/iemss2006/sessions/all.html>

648

649 Fang K-T, Li R, Sudjianto A. 2006. *Design and Modeling for Computer Experiments*,
650 Chapman & Hall/CRC, 290 pp.

651

652 Hacker JP, Snyder C, Ha S-Y, Pocerlich M. 2011. Linear and non-linear response to
653 parameter variations in a mesoscale model. *Tellus A*, **63**, 429-444.

654

655 Järvinen H, Laine M, Solonen A, Haario H. 2012. Ensemble prediction and parameter
656 estimation system: the concept. *Q. J. R. Meteorol. Soc.*, **138**, 281-288.

657

- 658 Laine M, Solonen A, Haario H, Järvinen H. 2012. Ensemble prediction and parameter
659 estimation system: the method. *Q. J. R. Meteorol. Soc.*, **138**, 289-297.
660
- 661 Leith CE. 1990. Stochastic backscatter in a subgrid-scale model: Plane shear mixing
662 layer. *Physics of Fluids A: Fluid Dynamics*, **2.3**, 297-299.
663
- 664 Ollinaho P, Järvinen H., Bauer P, Laine M, Bechtold P, Susiluoto J, Haario H. 2014.
665 Optimization of NWP model closure parameters using total energy norm of forecast error
666 as a target. *Geoscientific Model Development*, **7**, 1889-1900.
667
- 668 Marzban C, Sandgathe S, Lyons H, Lederer N. 2009. Three Spatial Verification
669 Techniques: Cluster Analysis, Variogram, and Optical Flow. *Wea. Forecasting*, **24**, 1457-
670 1471.
671
- 672 Marzban C. 2013. Variance-based Sensitivity analysis: An illustration on the Lorenz '63
673 model. *Mon. Weather. Rev.*, **141**, 4069-4079.
674
- 675 Marzban C, Sandgathe S, Doyle JD, Lederer NC. 2014. Variance-based sensitivity
676 analysis: Preliminary results in COAMPS. *Mon. Weather. Rev.*, **142**, 2028-2042.
677
- 678 Marzban, C, Du X, Sandgathe S, Doyle JD, Jin Y, Lederer NC. 2018a: Sensitivity
679 analysis of the spatial structure of forecasts in mesoscale models: Continuous model
680 parameters. *Mon. Wea. Rev.* **146**, 967-983.
681
- 682 Marzban, C, Jones C, Li N, Sandgathe S. 2018b: On the effect of model parameters on
683 forecast objects. *Geoscientific Model Developmet*, **11**, 1-14.
684
- 685 Mason PJ, Thomson DJ. 1994. Stochastic backscatter in large-eddy simulations of
686 boundary layers. *Journal of Fluid Mechanics*, **242**, 51-78.
687
- 688 McCann DW, Whistler JP. 2001. Problems and solutions for drawing fronts objectively.
689 *Meteorol. Appl.*, **8**, 195-203.
690
- 691 McKay MD, Beckman RJ, Conover WJ. 1979. A Comparison of Three Methods for
692 Selecting Values of Input Variables in the Analysis of Output from a Computer Code.
693 *Technometrics*, **21**, 239-245 .
694
- 695 Montgomery DC. 2009. *Design and Analysis of Experiments*, 7th Edition, John Wiley &
696 Sons, 656 pp.
697
- 698 Oakley JE, O'Hagan A. 2004. Probabilistic sensitivity analysis of complex models: a
699 Bayesian approach. *J. R. Statist. Soc., B*, **66**, 751-769.
700

- 701 Saltelli A, Annoni P, Azzini I, Campolongo F, Ratto M, Tarantola S, 2010. Variance
702 based sensitivity analysis of model output: Design and estimator for the total sensitivity
703 index. *Computer Physics Communications*, **181**, 259–270.
704
- 705 Sacks J, Welch WJ, Mitchell TJ, Wynn HP. 1989. Design and Analysis of Computer
706 Experiments. *Statistical Science*, **4**, 409-423.
707
- 708 Santner TJ, Williams BJ, Notz WI. 2003. *The Design and Analysis of Computer*
709 *Experiments*. Springer, 299pp.
710
- 711 Smith SA, Vosper SB, Field PR. 2015. Sensitivity of orographic precipitation
712 enhancement to horizontal resolution in the operational Met Office Weather
713 forecasts. *Meteorol. Appl.*, **22**, 14-24.
714
- 715 Skamarock WC, Klemp JB. 2008. A time-split nonhydrostatic atmospheric model for
716 weather research and forecasting applications. *J. Comp. Phys.*, **227**, 3465-3485.
717
- 718 Sobol' IM. 1993. Sensitivity estimates for nonlinear mathematical models. *Mathematical*
719 *Modeling and Computational Experiments*, **1**, 407-414.
720
- 721 Stein U. Alpert P. 1993. Factor separation in numerical simulations. *J. Atmos. Sci.*, **50**,
722 2107-2115.
723
- 724 Welch WJ, Buck RJ, Sacks J, Wynn HP, Mitchell TJ, Morris MD. 1992. Screening,
725 Predicting, and Computer Experiments. *Technometrics*, **34**, 15-25.
726
- 727 Yang Y, Uddstrom M, Revell M, Moore S. 2014. Soil moisture simulation by JULES in
728 New Zealand: verification and sensitivity tests. *Meteorol. Appl.*, **21**, 888-897.
729
- 730 Zhao J, Tiede C. 2011. Using a variance-based sensitivity analysis for analyzing the
731 relation between measurements and unknown parameters of a physical model. *Nonlin.*
732 *Processes Geophys.*, **18**, 269276.

733 Figure Captions

734

735 Figure 1. The histogram of a) the number of precipitation objects, and their b) size and c)
736 intensity. Panel d shows the scatterplot of size versus the number of objects.

737

738 Figure 2. Jet streak objects identified within the WRF wind field at 250 hPa for a 42-hour
739 forecast initialized at 00 UTC on February 9 2015. Jet streaks are identified by white
740 contour lines, and the location of the maximum wind speed within each object is
741 identified by the white-contoured black dot.

742

743 Figure 3. The “time series” of the five response variables: The Number of objects (a), the
744 size of the smallest (b) and largest (c) objects, and the intensity of the weakest (d) and
745 strongest (e) objects. The colors correspond to the six replications, and the thick/black
746 line corresponds to a control run wherein SKEBS has been turned off. Par1 and Par2 are
747 set to their default SKEBS values (10^{-5} and 10^{-6} , respectively). Size refers to the number
748 of grid points in an object, and intensity is measured in m/s.

749

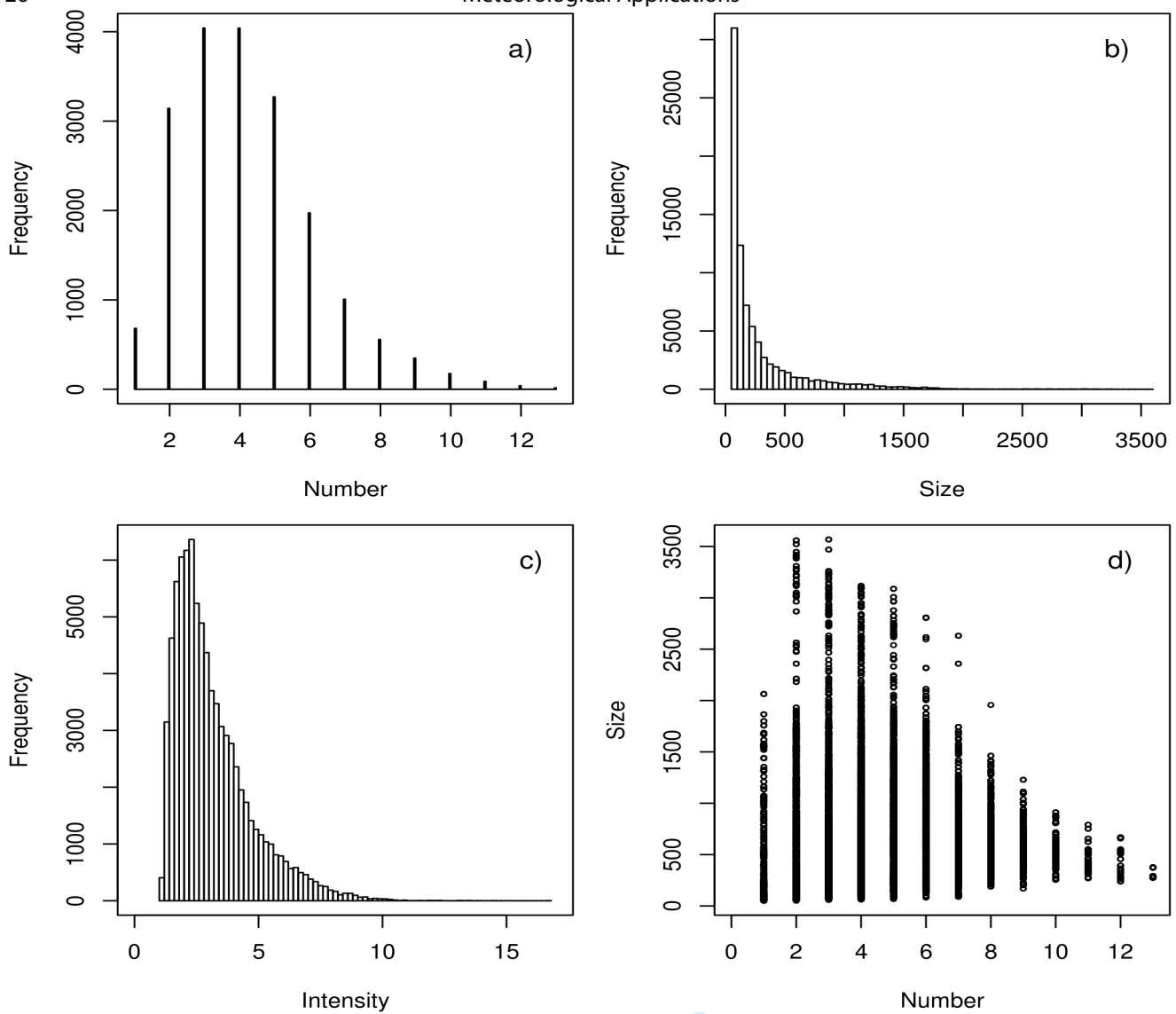
750 Figure 4. The variance (across all factors - Day, Rep, Par1, and Par2) of the five response
751 variables - Number of objects (top), minimum size (black) and maximum size (red)
752 (middle panel), and minimum intensity (black) and maximum intensity (red) (bottom
753 panel). The vertical lines are 95% confidence intervals, displaying the uncertainty in these
754 variance estimates. Size refers to the number of grid points in an object, and intensity is
755 measured in m/s.

756

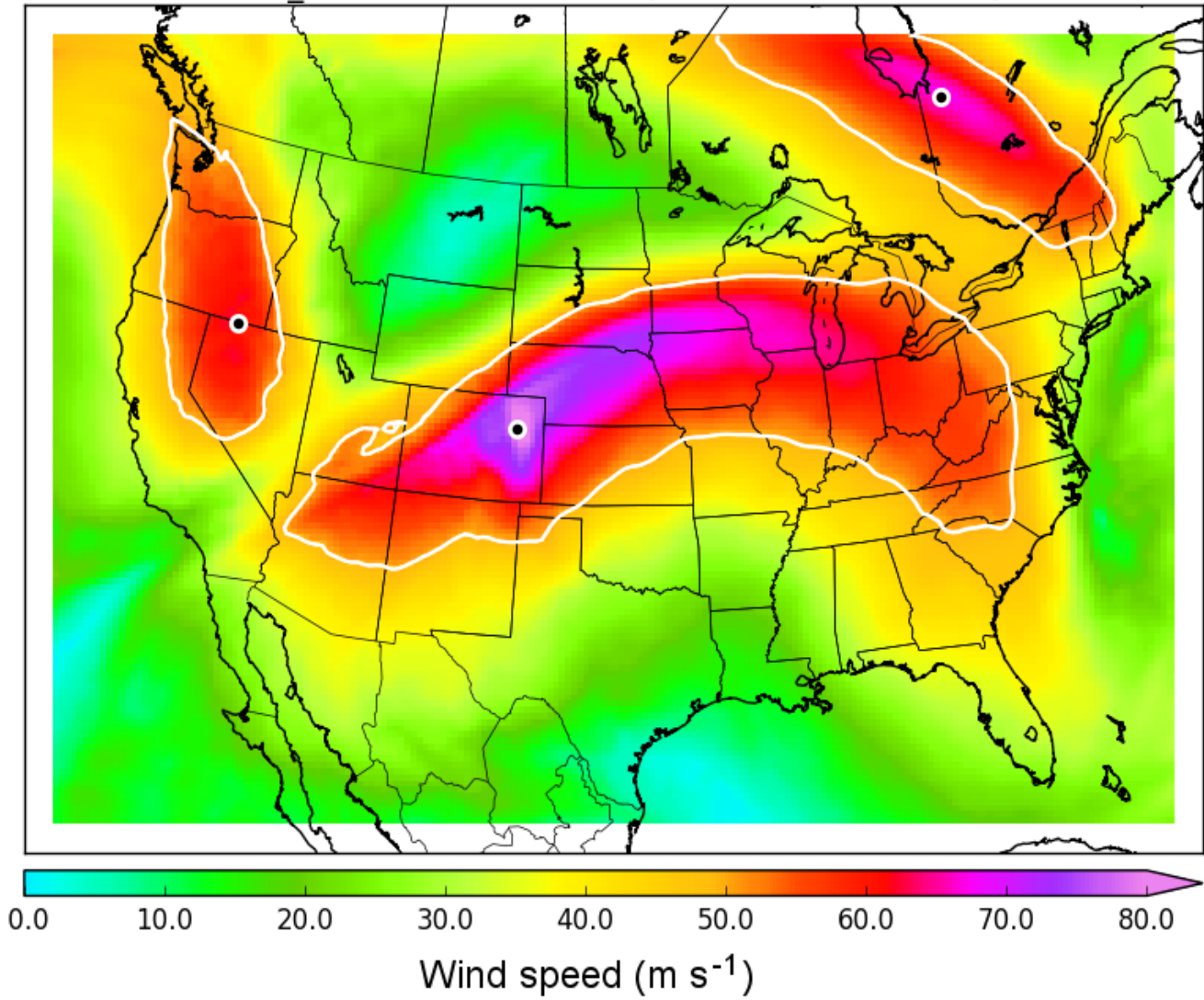
757 Figure 5. The distribution/boxplot (across 120 forecast hours) of p-values testing the
758 significance of the main effects for the four factors (Day, Rep, Par1, Par2) on the five
759 responses (panels a-e).

760

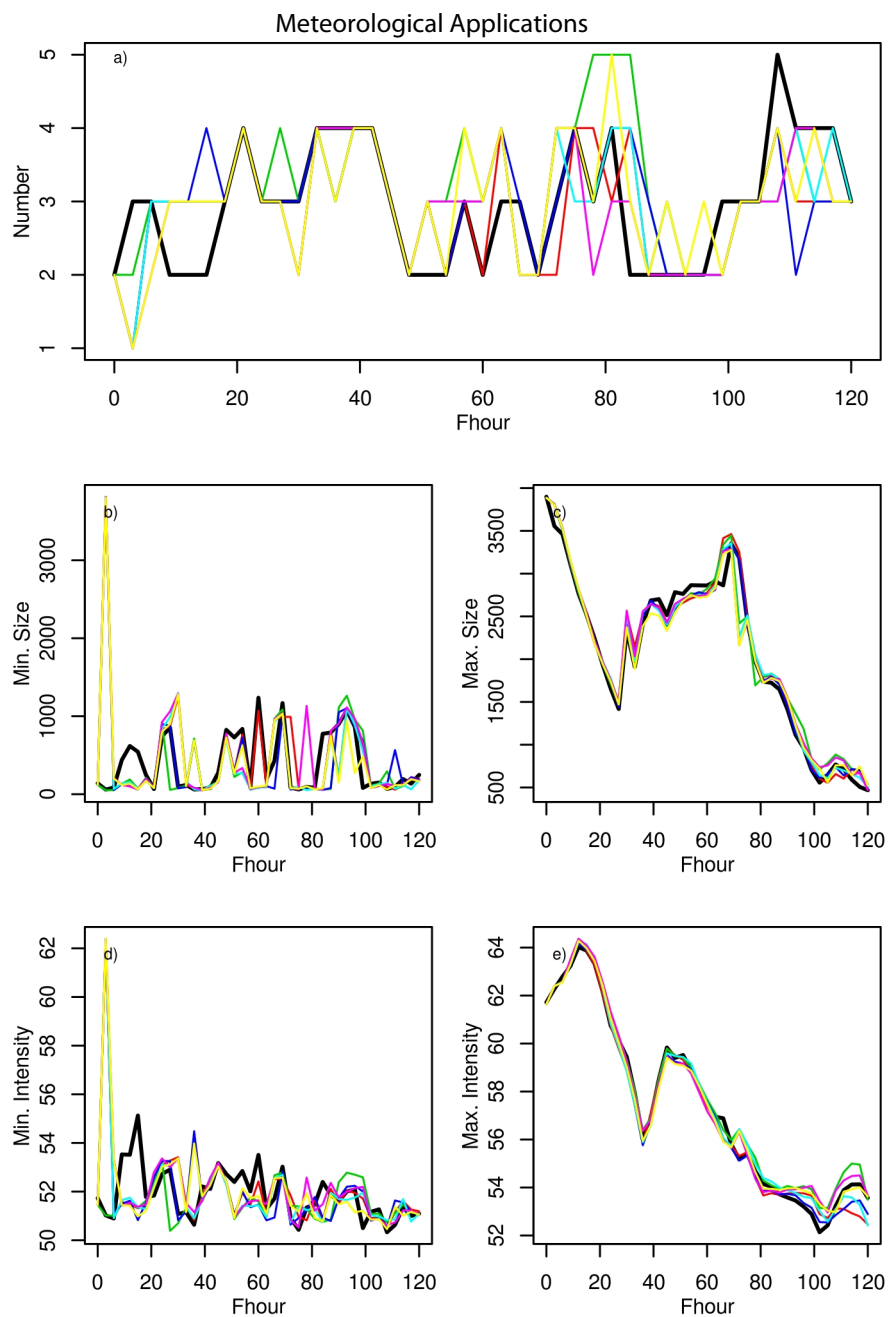
761 Figure 6. The 95% confidence intervals for the intraclass correlation ρ versus forecast
762 time, displaying the effect of the four factors, Day, Rep, Par1 and Par2 (top 4 rows) on the
763 five responses - number of objects (left column), minimum (black) and maximum (red)
764 size (middle column), and minimum (black) and maximum (red) intensity (right column).
765 The last row shows ρ_ϵ , the proportion of total variance in the response not explained by
766 the four factors.



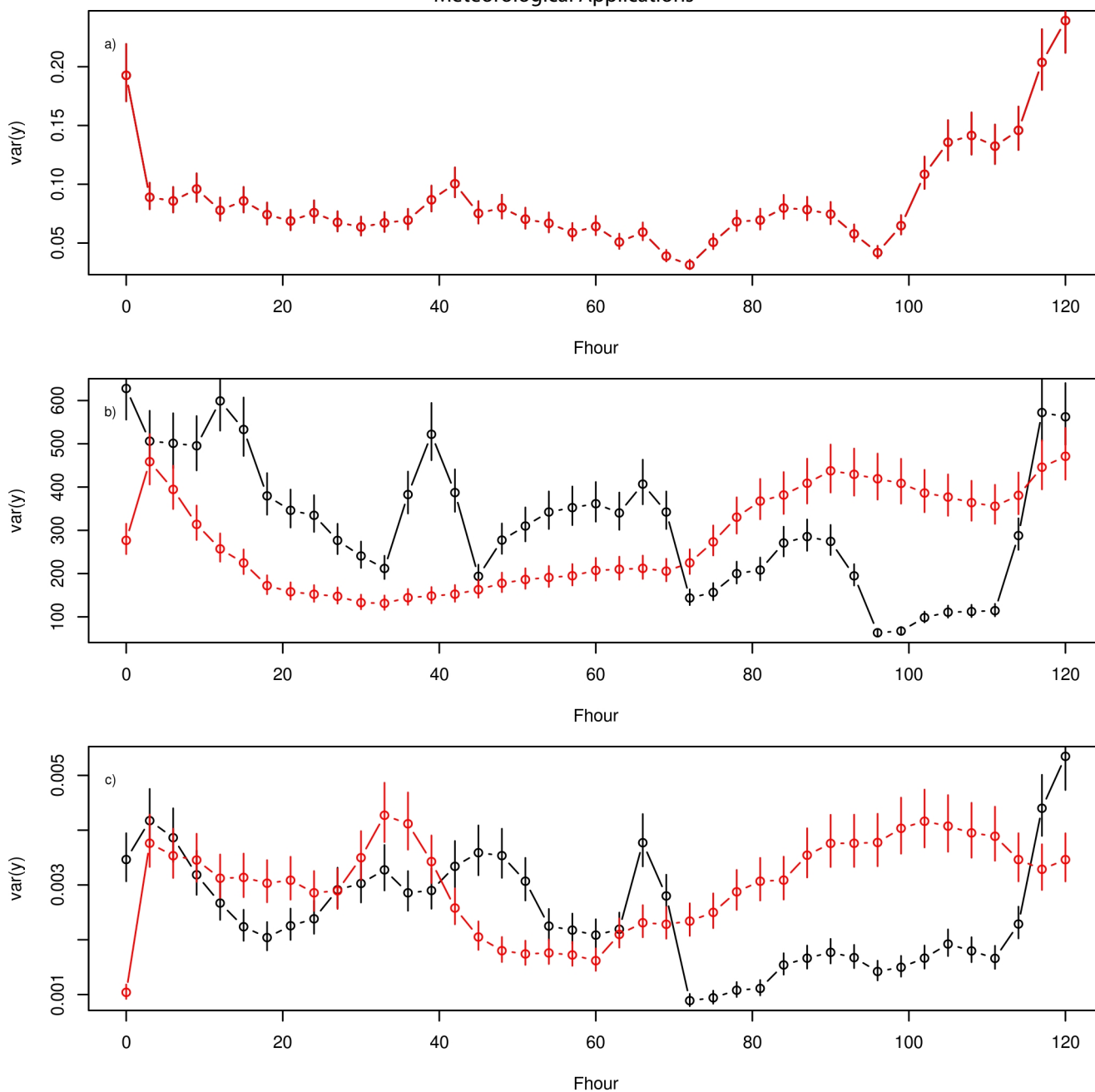
2 Figure 1. The histogram of a) the number of precipitation objects, and their b) size and c)
3 intensity. Panel d) shows the scatterplot of size versus the number of objects.



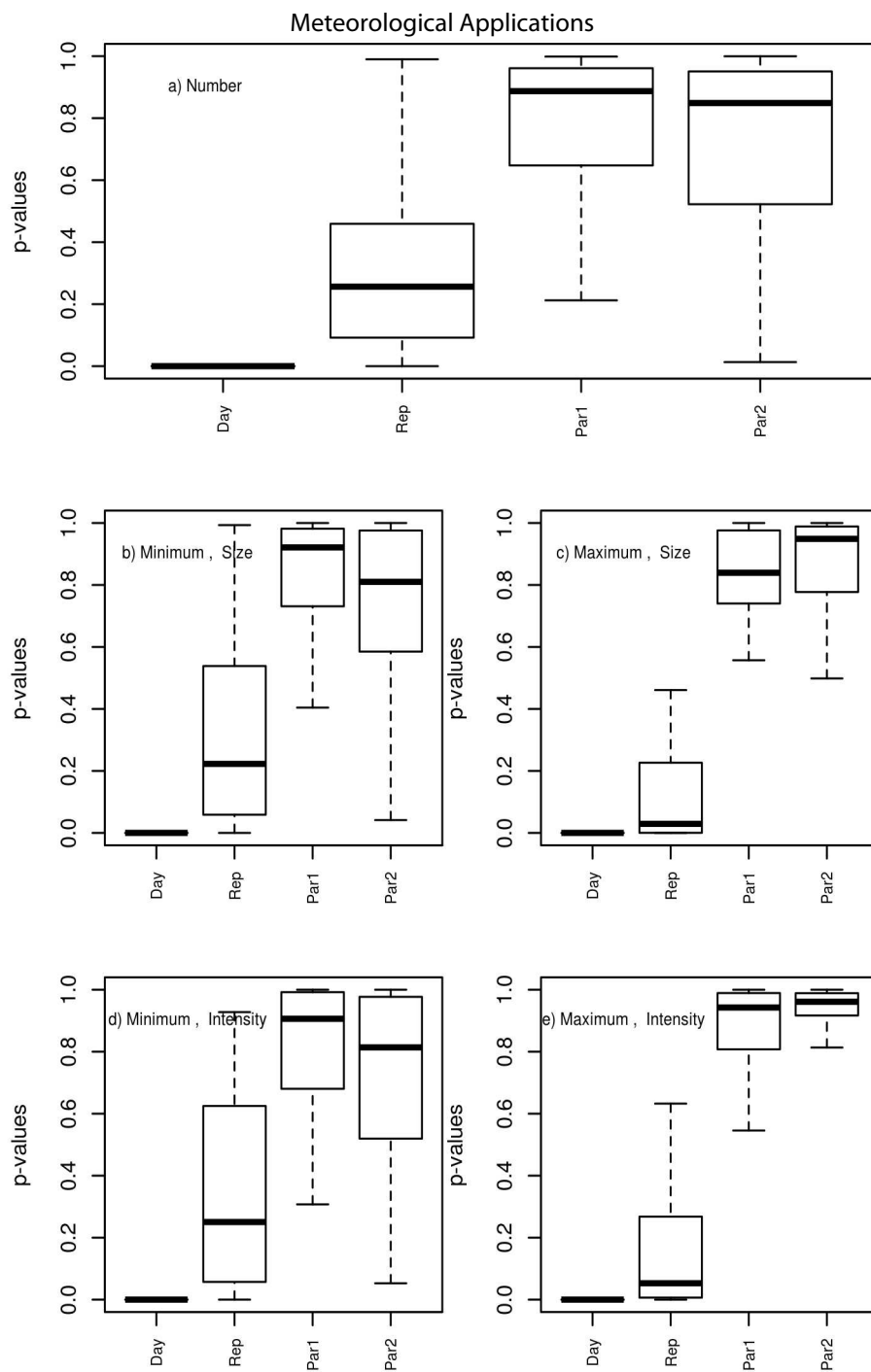
2 Figure 2. Jet streak objects identified within the WRF wind field at 250 hPa for a 42-hour
 3 forecast initialized at 00 UTC on February 9 2015. Jet streaks are identified by white contour
 4 lines, and the location of the maximum wind speed within each object is identified by the
 5 white-contoured black dot.



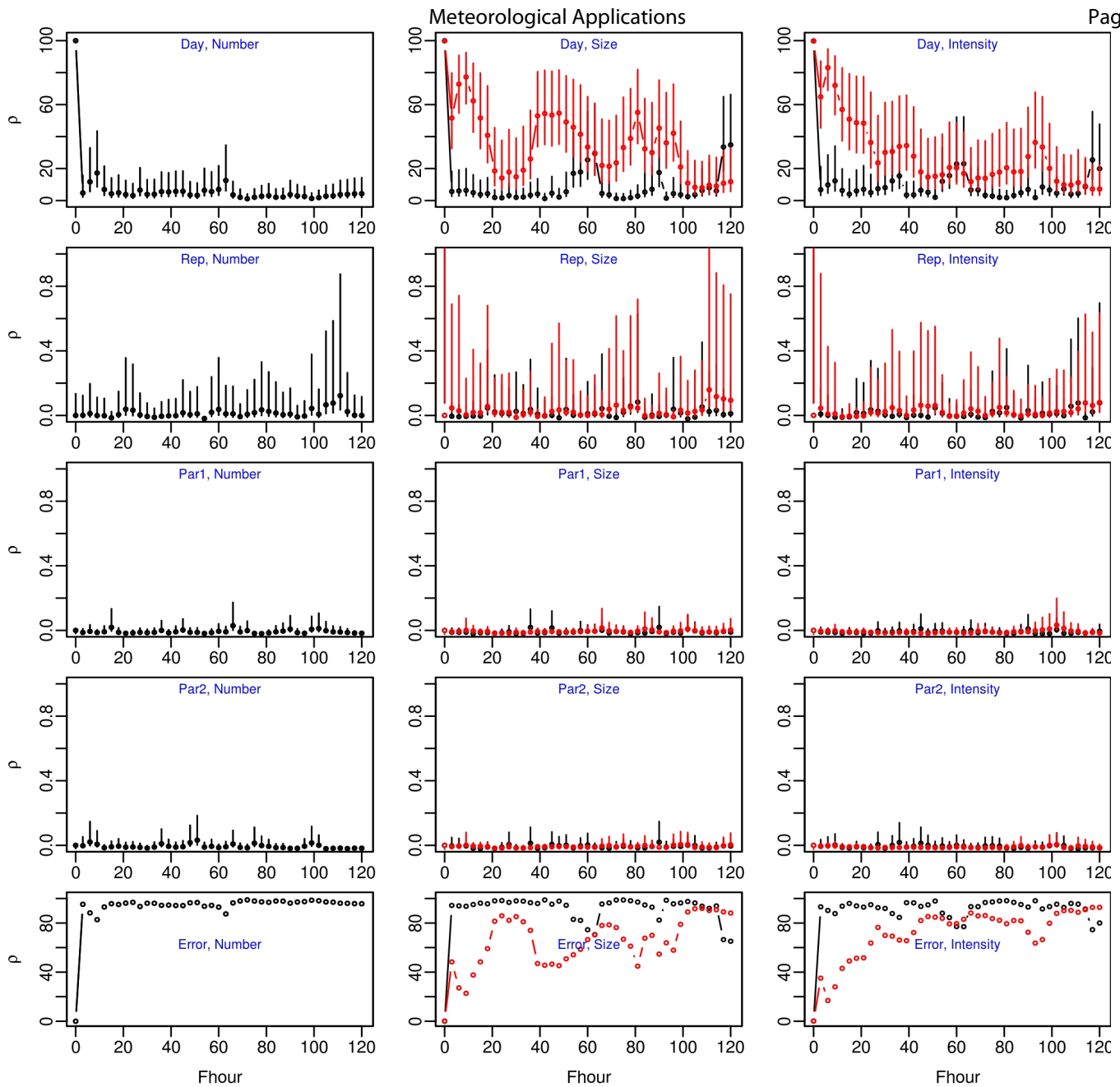
2 Figure 3. The “time series of the five response variables: The Number of objects (a), the
 3 size of the smallest (b) and largest (c) objects, and the intensity of the weakest (d) and
 4 strongest (e) objects. The colors correspond to the six replications, and the thick/black line
 5 corresponds to a control run wherein SKEBS has been turned off. Par1 and Par2 are set to
 6 their default SKEBS values (10^{-5} and 10^{-6} , respectively). Size refers to the number of grid
 7 points in an object, and intensity is measured in m/s .



2 Figure 4. The variance (across all factors - Day, Rep, Par1, and Par2) of the five response
 3 variables - Number of objects (top), minimum size (black) and maximum size (red) (middle
 4 panel), and minimum intensity (black) and maximum intensity (red) (bottom panel). The
 5 vertical lines are 95% confidence intervals, displaying a sense of the uncertainty in these
 6 variance estimates. Size refers to the number of grid points in an object, and intensity is
 7 measured in m/s .



2 Figure 5. The distribution/boxplot (across 120 forecast hours) of p-values testing the signif-
 3 icance of the main effects for the four factors (Day, Rep, Par1, Par2) on the five responses
 4 (panels a-e).



2 Figure 6. The 95% confidence intervals for the intraclass correlation ρ versus forecast time,
 3 displaying the effect of the four factors, Day, Rep, Par1 and Par2 (top 4 rows) on the five
 4 responses - number of objects (left column), minimum (black) and maximum (red) size
 5 (middle column), and minimum (black) and maximum (red) intensity (right column). The
 6 last row shows ρ_ϵ , the proportion of total variance in each response not explained by the
 7 four factors.

1 Table 1. An example of an LSD involving three factors A, B, C, each taking three values
2 (denoted by the indices 1, 2, 3).

3

	A1	A2	A3
B1	C1	C2	C3
B2	C2	C3	C1
B3	C3	C1	C2

4

For Peer Review

A Methodology for Sensitivity Analysis of Spatial Features in Forecasts: The Stochastic Kinetic Energy Backscatter Scheme

Caren Marzban*, Robert Tardif, Natalia Hryniw, Scott Sandgathe

All numerical models have parameters whose values are often set in an *ad hoc* fashion, and so, it is important to assess how these parameters affect the output of the model. The output of many models often contain “objects” examples of which are shown in the figure below. This paper proposes a methodology for assessing how the model parameters affect specific features of such objects.

