

1 A Quantile-Conserving Ensemble Filter Framework. Part III: Data Assimilation for Mixed  
2 Distributions with Application to a Low-Order Tracer Advection Model  
3 Jeffrey Anderson<sup>a</sup>, Chris Riedel<sup>b</sup>, Molly Wieringa<sup>c</sup>, Fairuz Ishraque<sup>d</sup>, Marlee Smith<sup>a</sup>, Helen  
4 Kershaw<sup>a</sup>

5  
6 <sup>a</sup> NCAR/CISL/TDD/DAReS, Boulder, Colorado, jla@ucar.edu

7 <sup>b</sup> Cooperative Programs for the Advancement of Earth System Science, University  
8 Corporation for Atmospheric Research

9 <sup>c</sup> University of Washington, Department of Atmospheric Sciences, Seattle, Washington

10 <sup>d</sup> Department of Geosciences, Princeton University, Princeton, New Jersey

11

12 Submitted to Monthly Weather Review, November 2023

13

14 **Abstract**

15

16 The uncertainty associated with many observed and modeled quantities of interest in Earth  
17 system prediction can be represented by mixed probability distributions that are neither  
18 discrete nor continuous. For instance, a forecast probability of precipitation can have a finite  
19 probability of zero precipitation, consistent with a discrete distribution. However, nonzero  
20 values are not discrete and are represented by a continuous distribution; the same is true for  
21 rainfall rate. Other examples include snow depth, sea ice concentration, amount of a tracer or  
22 the source rate of a tracer. Some Earth system model parameters may also have discrete or  
23 mixed distributions. Most ensemble data assimilation methods do not explicitly consider the  
24 possibility of mixed distributions. The Quantile Conserving Ensemble Filtering Framework  
25 (Anderson 2022, 2023) is extended to explicitly deal with discrete or mixed distributions. An  
26 example is given using bounded normal rank histogram probability distributions applied to  
27 observing system simulation experiments in a low-order tracer advection model. Analyses of  
28 tracer concentration and tracer source are shown to be improved when using the extended  
29 methods. A key feature of the resulting ensembles is that there can be ensemble members with  
30 duplicate values. An extension of the rank histogram diagnostic method to deal with potential  
31 duplicates shows that the ensemble distributions from the extended assimilation methods are  
32 more consistent with the truth.

33

34 SIGNIFICANCE STATEMENT: Data assimilation is a statistical method that is used to combine  
35 information from computer forecasts with measurements of the Earth system. The result is a  
36 better estimate of what is occurring in the physical system. As an example, data assimilation is  
37 used for making weather predictions. Some Earth system quantities, like precipitation, have  
38 special values that can occur very frequently. For instance, zero rainfall is quite common, while  
39 any other specific amount of rainfall, say 0.42 inches, is unusual. New data assimilation tools  
40 that work well for quantities like this are introduced and should lead to better estimates and  
41 predictions of the Earth system.

42

43 KEYWORDS: Data assimilation, Ensembles, Uncertainty, Atmospheric Chemistry

44

## 45 **1. Introduction**

46

47 Ensemble data assimilation methods have been widely applied across Earth system  
48 applications. The input to the assimilation method is an ensemble of forecasts that is assumed  
49 to be a random sample of the uncertainty of a model state vector. Atmospheric data  
50 assimilation for numerical weather prediction remains the most common application  
51 (Houtekamer and Zhang 2016). In this case, the uncertainty distributions for many variables like  
52 temperature, velocity components, and surface pressure are expected to be approximately  
53 normal. Many existing ensemble filter algorithms implicitly assume normality (Burgers et al.  
54 1998, Houtekamer and Mitchell 1998, Pham 2001, Anderson 2001) and are very successful for  
55 weather prediction applications.

56

57 Other types of continuous distributions may be more appropriate for the uncertainty of other  
58 variables (Bocquet et al. 2010). For instance, log-normal (Fletcher and Zupanski 2006), gamma  
59 and inverse gamma distributions might be more appropriate for variables that are bounded like  
60 specific humidity (Bannister et al., 2020). Ensemble filters that can represent gamma and  
61 inverse gamma distributions have been developed (Bishop 2016). Other ensemble methods  
62 have been developed to transform distributions so that they are more normally distributed  
63 (Doron et al. 2013, Kurosawa and Poterjoy 2021), allowing normal ensemble algorithms to work  
64 better (Simon and Bertino 2012). The term Gaussian anamorphosis (Bertino et al. 2003) has  
65 been applied to some of these methods (Beal et al. 2010, Amezcua and Van Leeuwen 2014).  
66 Mixtures of standard continuous distributions like Gaussian kernels (Anderson and Anderson  
67 1999, Grooms 2022) including binormal distributions (Chan et al. 2020) have also been applied.

68

69 The uncertainty for some variables is a mixed probability distribution that includes both  
70 discrete and continuous parts. As an example, the amount of precipitation that falls during a  
71 particular period (Suhaila et al. 2011) might have a discrete probability of being exactly zero in

72 addition to a continuous distribution of being non-zero; the precipitation rate would have a  
73 similar mixed distribution. The amount of sea ice, snow cover, chemical tracer, or water in a  
74 stream also have mixed distributions along with their source and sink rates. Quantities like the  
75 fractional coverage of ice or snow are doubly bounded, and could have a discrete probability of  
76 no cover, a discrete probability of complete coverage, and a continuous distribution for all  
77 intermediate values. A beta distribution might be appropriate for some doubly bounded  
78 quantities.

79

80 Anderson (2003) described a two-step algorithm for computing a variety of ensemble Kalman  
81 filter algorithms and this methodology was extended for more general problems in Grooms  
82 (2022). The input to the first step is an ensemble of estimates of an observed quantity and the  
83 likelihood of the observation, while the output is an ensemble of increments due to the  
84 observation. The second step is a bivariate algorithm that independently computes increments  
85 for each individual model state variable given the increments from step one.

86

87 The first part of this quantile conserving ensemble filter framework (QCEFF) paper sequence  
88 (Anderson 2022; A22 hereafter) describes the use of quantile conserving ensemble filters for  
89 the first step of the two-step algorithm. This allows almost any continuous probability  
90 distribution function (PDF) to be used for the computation of observation increments. The  
91 second part of the QCEFF sequence (Anderson 2023; A23 hereafter) addresses the second part  
92 of the two-step algorithm. It uses a specific variant of anamorphosis, the probit probability  
93 integral (PPI) transform (Amezcuca and Van Leeuwen 2014), to make the bivariate problem  
94 more normal. Again, arbitrary continuous PDF can be used for the probability integral transform  
95 portion of the algorithm. Both QCEFF papers include a description of a particular type of nearly  
96 non-parametric distribution, the bounded normal rank histogram (BNRH) distribution that can  
97 be useful for data assimilation when the details of an appropriate parametric distribution are  
98 not known a priori.

99

100 A22 provides an example using a discrete distribution that is closely related to the particle filter  
101 (Van Leeuwen 2009, Van Leeuwen et al. 2019) and A23 mentions the possibility of using a  
102 similar distribution for the PPI transform. However, neither manuscript provides a detailed  
103 description of the implementation of the discrete distribution and neither explores mixed  
104 distributions. This third part of the QCEFF sequence begins by describing a general framework  
105 for using mixed distributions to represent uncertainty in ensemble filters in section 2. When  
106 ensemble methods are applied for mixed distributions, ensemble members with identical  
107 values for a given state variable are expected to occur. Section 3 extends the results of section  
108 2 to describe a BNRH distribution that works with ensembles with duplicate members. Section 3  
109 also describes an extension of the rank histogram diagnostic tool to ensembles with duplicate  
110 members. Section 4 describes an extension of the low-order Lorenz-96 model to include an  
111 advected tracer and a source. This model is configured to generate ensembles with duplicate  
112 members for both the tracer concentration and source ensemble estimates. Observing system  
113 simulation experiments in Section 5 compare the capabilities of several ensemble filter variants  
114 in this model. Section 6 provides discussion and conclusions.

115

## 116 **2. QCEFF for discrete and mixed probability distributions**

117

118 The QCEFF developed in A22 for the first part of the two-step ensemble DA algorithm requires  
119 finding an appropriate PDF and corresponding cumulative distribution function (CDF) given an  
120 ensemble. It requires multiplying the PDF times a likelihood function to get an analysis  
121 (posterior) PDF and corresponding CDF. It also requires evaluating CDFs and their inverses; this  
122 is also necessary for the probit probability integral (PPI) transforms used for QCEFF  
123 implementations of the second part of the two-step algorithm in A23. A22 includes a brief  
124 discussion of using a particle filter as the prior generalized PDF and provides an example  
125 without carefully defining the algorithm. This section begins by clarifying the application of the  
126 QCEFF for discrete probability distributions (like the particle filter), then extends that to mixed  
127 probability distributions.

128

129 Here, a discrete probability distribution consists of a set of  $K$  real numbers,  $\{x_i, i = 1, \dots, K\}$   
 130 and associated positive real probabilities  $p_i$  with

$$131 \quad \sum_{i=1}^K p_i = 1. \quad (1)$$

132 Suppose a discrete generalized PDF is used as the prior for an observed quantity in data  
 133 assimilation and the observation likelihood is  $L(x)$ . The normalizing constant for the product of  
 134 the prior and the likelihood is

$$135 \quad S = \sum_{i=1}^K L(x_i)p_i. \quad (2)$$

136 An analysis generalized PDF then has the same  $\{x_i\}$  with probabilities

$$137 \quad p_i^a = L(x_i)p_i/S. \quad (3)$$

138

139 To use the QCEFF, it is necessary to evaluate the CDF, and its inverse, corresponding to a  
 140 discrete generalized PDF. Defining the CDF as the integral from  $-\infty$  to  $x$  of the generalized PDF  
 141 leads to discrete jumps at each  $x_i$  so that the CDF is not a function. For the QCEFF, a  
 142 generalized CDF,  $\tilde{F}$ , that is a function is defined by making the value at  $x_i$  the midpoint of the  
 143 jump,

$$144 \quad \tilde{F}(x) = \begin{cases} 0 & \text{if } x < x_1 \\ \sum_{k=1}^i p_k & \text{if } x_i < x < x_{i+1}, \quad i \in \{1, \dots, K-1\} \\ 1 & \text{if } x > x_K \\ \sum_{k=1}^{i-1} p_k + \frac{p_i}{2} & \text{if } x = x_i \end{cases} \quad (4)$$

145 A generalized inverse CDF is defined as

$$146 \quad \tilde{F}^{-1}(y) = \begin{cases} x_1 & \text{if } y \leq p_1 \\ x_i & \text{if } \sum_{k=1}^{i-1} p_k < y \leq \sum_{k=1}^i p_k, \quad i \in \{2, \dots, K\} \end{cases} \quad (5)$$

147 Note that  $x = \tilde{F}^{-1}(\tilde{F}(x))$  but  $\tilde{F}(\tilde{F}^{-1}(y))$  is not necessarily equal to  $y$ . With these definitions,  
 148 it is possible to define a QCEFF that uses any discrete prior, like a particle filter, in observation  
 149 space for the first part of the two-step filter and for the PPI in the regression step.

150

151 As noted in the introduction, mixed distributions are relevant to many geophysical problems.

152 The discrete part of a prior mixed distribution is represented as above except that  $\sum p_i = \alpha$ ;

153 the continuous part of the PDF is  $(1 - \alpha)f_c(x)$ , with  $0 < \alpha < 1$ . The normalizing constant for

154 the product with a likelihood is

155 
$$S = \alpha \sum_{i=1}^K L(x_i) p_i + (1 - \alpha) \int_{-\infty}^{\infty} L(x) f_c(x) dx \quad (6)$$

156 The analysis generalized PDF has discrete part as in (3) and the continuous part

157 
$$(1 - \alpha) f_c(x) L(x) / S. \quad (7)$$

158 A generalized CDF corresponding to a mixed PDF is

159 
$$\tilde{F}_m = (1 - \alpha) \int_{-\infty}^x f_c(x) dx + \alpha \tilde{F}(x) \quad (8)$$

160 where  $\tilde{F}$  is defined in (4). The inverse is clearly defined except at the jumps from the discrete  
161 part of the mixed distribution. Define the bounds of the jumps as

162 
$$J_i^- = \begin{cases} (1 - \alpha) \int_{-\infty}^{x_1} f_c(x) dx & \text{if } i = 1 \\ (1 - \alpha) \int_{-\infty}^{x_i} f_c(x) dx + \sum_{k=1}^{i-1} \alpha p_k, & i \in \{2, \dots, K\} \end{cases} \quad (9)$$

163 and

164 
$$J_i^+ = J_i^- + \alpha p_i, \quad i \in \{1, \dots, K\} \quad (10)$$

165 The inverse can be defined between the jump values as

166 
$$\tilde{F}_m^{-1}(y) = x_i \quad \text{for } J_i^- \leq y \leq J_i^+ \quad (11)$$

167

### 168 3. Tools for data assimilation with duplicate ensemble members

169

#### 170 a. Bounded normal rank histogram distribution

171

172 The QCEFF described in A22 and A23 requires a CDF to compute observation increments and to  
173 do the regression of those increments onto model state variables. The bounded normal rank  
174 histogram (BNRH) distribution is an extension of the rank histogram filter distribution for  
175 observation space increments (Anderson 2010). A BNRH distribution is particularly useful when  
176 the appropriate distribution family is unknown.

177

178 A23 describes the PDF,  $f(x)$ , associated with a BNRH when there are no duplicate ensemble  
179 members. An N-member ensemble partitions the real line into N+1 intervals. The interior  
180 intervals are bounded on both sides; the intervals on the tails can be bounded on one side only  
181 if the quantity itself is not bounded, or bounded on both sides if the quantity is bounded. The  
182 BNRH PDF assigns  $1/(N+1)$  probability to each interval. The probability is uniformly distributed

183 over the range of an interior interval. For intervals on the tails, the probability density is part of  
 184 a normal distribution. The DA algorithms in A22 and A23 require the CDF which is defined in the  
 185 standard fashion as  $F(x) = \int_{-\infty}^x f(x)dx$ . An example of a BNRH CDF is shown in Figure 1a  
 186 (reproduced from A23) for a 5-member ensemble.

187  
 188 The definition of the BNRH CDF is extended here for the case when there are ensemble  
 189 members with duplicate values or when one or more ensemble members have the same value  
 190 as the upper or lower bound of  $x$ . Suppose that possible values of  $x$  are bounded below by  
 191  $B_l \geq -\infty$  and above by  $B_u \leq \infty$ . Given an  $N$ -member ensemble of  $x$  with members not  
 192 necessarily unique, there is at least one ordering of the ensemble values so that  $x_i \leq x_{i+1}$  for  
 193  $i \in \{1, \dots, N - 1\}$ . Given such an ordering, define the CDF as:

$$194 \quad F(x) = \begin{cases} 0 & \text{if } x < B_l \\ C(B_l)/[2(N + 1)] & \text{if } x = B_l \\ A_l \Phi(\mu_l, \sigma^2; x) - A_l \Phi(\mu_l, \sigma^2; B_l) & \text{if } B_l < x < x_1 \\ [i + (x - x_i)/(x_{i+1} - x_i)]/(N + 1) & \text{if } x_i < x < x_{i+1}, \quad i \in \{1, \dots, N - 1\} \\ i/(N + 1) + [C(x) - 1]/[2(N + 1)] & \text{for min } i \text{ with } x = x_i, B_l < x < B_u, \quad i \in \{1, \dots, N\} \\ A_u \Phi(\mu_u, \sigma^2; x) - A_u \Phi(\mu_u, \sigma^2; B_u) + 1 & \text{if } x_N < x < B_u \\ 1 - C(B_u)/[2(N + 1)] & \text{if } x = B_u \\ 1 & \text{if } x > B_u \end{cases}$$

196 (12)

197  $C(x)$  is a function with unbounded real domain and range the whole numbers less than or  
 198 equal to  $N$ , defined as the number of ensemble members with value  $x$ .  $\Phi(\mu, \sigma^2; x)$  is the CDF  
 199 of a normal with mean  $\mu$  and variance  $\sigma^2$  evaluated at  $x$ , and  $\sigma^2$  is the sample variance of the  
 200 ensemble. The means and amplitudes of the normal portions are defined as in A23 so that  
 201  $1/(N + 1)$  probability lies between the outermost ensemble member and the bounds. The  
 202 means are selected so that

$$203 \quad \Phi(\mu_l, \sigma^2; x_1) = \frac{1}{N+1} \quad (13)$$

$$204 \quad \Phi(\mu_u, \sigma^2; x_N) = \frac{N}{N+1} \quad (14)$$

205 and the amplitudes are



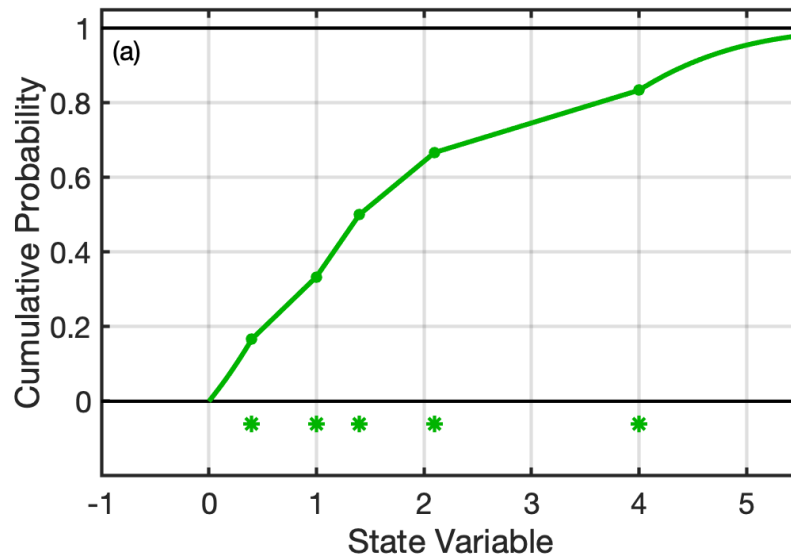
206 
$$A_l = \frac{1}{(N+1)[\Phi(\mu_l, \sigma^2; x_1) - \Phi(\mu_l, \sigma^2; B_l)]} \quad (15)$$

207 
$$A_u = \frac{1}{(N+1)[\Phi(\mu_u, \sigma^2; B_u) - \Phi(\mu_u, \sigma^2; x_N)]} \quad (16)$$

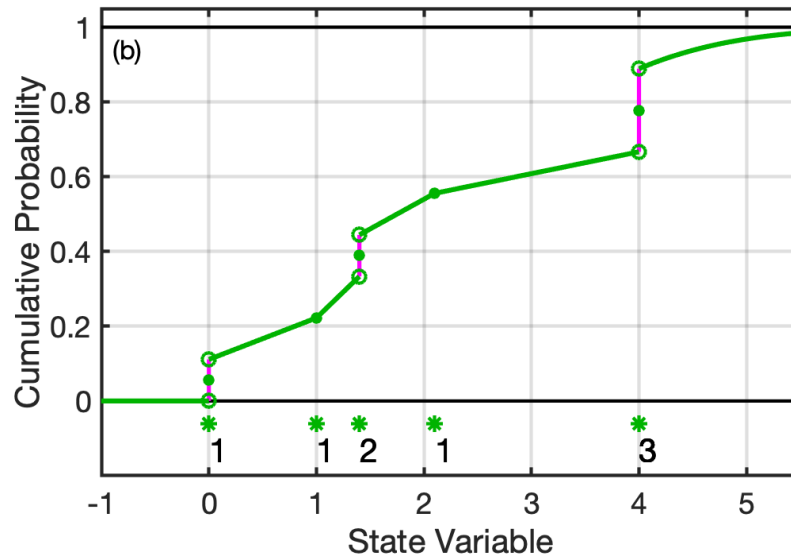
208

209 When there are no duplicate ensemble values,  $C(x_i) = 1 \forall x_i$ , and no ensemble values equal to  
 210 the bounds,  $C(B_l) = C(B_u) = 0$ , the BNRH CDF is equal to the integral from  $-\infty$  to  $x$  of the  
 211 BNRH PDF defined in appendix C of A23 and  $F$  is invertible. However, where  $C(x) > 1$ , or if  
 212  $C(B_l) > 0$  or  $C(B_u) > 0$ , there is a discrete probability, the derivative  $dF(x)/dx$  is undefined,  
 213 and  $F$  is not invertible. It is necessary to define a generalized inverse following the procedure  
 214 for mixed probability distributions in section 2 (A22 notes the need for a generalized inverse for  
 215 some other distribution families in which the PDF is 0 over a bounded range of  $x$ ).

216



217



218

219

220 Figure 1: Cumulative distribution functions (green) for a BNRH distribution for a 5-member  
 221 ensemble (green asterisks) for a variable that is bounded below at zero (a) and for an 8-  
 222 member ensemble with duplicate values and a member with a value at the bound of zero (b).  
 223 The number of duplicates is given by the integer next to asterisk. The vertical magenta lines  
 224 indicate the inverse cumulative distribution function (the quantile function) used for the BNRH.  
 225 Panel a is reproduced from figure C1a in A23.  
 226

226

227 An example CDF for an 8-member ensemble with  $B_l = x_1 < x_2 < x_3 = x_4 < x_5 < x_6 = x_7 =$   
 228  $x_8$  and  $B_u = \infty$  is shown in green in Figure 1b. The interval on the upper tail is a portion of a  
 229 normal CDF.  $1/(N+1)$  probability is uniformly distributed in each interior interval. In non-zero  
 230 range interior intervals, the CDF is piecewise linear. In the case of the duplicate ensemble  
 231 members, the range of the interval between them can be thought of as zero and the  
 232 distribution is discrete. At the point  $x_3$  where there are two ensemble members, there is  
 233  $1/(N+1)$  probability while at the point  $x_6$  with three ensemble members, there is  $2/(N+1)$   
 234 probability. Generalizing, at any point with D duplicate ensemble members, there is  $(D-1)/(N+1)$   
 235 discrete probability. Consistent with section 2 and eq. 12, the BNRH CDF at a point with  
 236 duplicate ensemble members is set to the ‘midpoint’ of the discontinuous jump in the integral  
 237 of the PDF. For instance, at  $x_3$  the CDF is defined as

238 
$$F(x_3) = \left[ \frac{3}{(N+1)} + \frac{4}{N+1} \right] / 2. \quad (17)$$

239 With this extended definition of the CDF, the quantile computed for ensemble members that  
240 share a value is the same. The inverse of the CDF is also needed for the QCEFF algorithms, and it  
241 is not uniquely defined with duplicate ensemble members. The method in section 2 leads to  
242 defining the inverse as the magenta lines in Fig. 1b.

243

#### 244 *b. Rank histograms*

245

246 Consider a sample of  $N + 1$  numbers composed of an  $N$ -member ensemble estimate of a scalar  
247 quantity and an additional value, called the verification here. If there are no duplicate values in  
248 the sample, the rank of the verification is uniquely defined with an integer value in  
249  $\{1, 2, \dots, N + 1\}$ . Define a rank weight vector,  $W_n$ ,  $n = 1, \dots, N + 1$  as

$$250 \quad W_n = \begin{cases} 1 & \text{if } \text{rank}(\text{verification}) = n \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

251 Define the sum of the rank weight vector for a collection of  $M$  ensembles with verifications as

$$252 \quad S_n = \sum_{m=1}^M W_n^m \quad (19)$$

253 A histogram of the vector  $S$ , commonly called the rank histogram (Anderson 1996, Hamill 2001)  
254 is a diagnostic tool for evaluating the consistency of ensemble predictions. If the verification for  
255 each ensemble is drawn from the same distribution as the ensemble, the histogram is expected  
256 to be statistically uniform. Histograms that are not uniform can provide information about the  
257 differences between ensembles and verification. For instance, a U-shaped histogram can  
258 indicate under dispersive ensembles (Wilks 2019).

259

260 For state variables in many common Earth system DA applications, the probability that the  
261 verification duplicates one or more ensemble members is very small, and most discussions of  
262 rank histograms have ignored the possibility. However, this is no longer the case for some types  
263 of bounded state variables which have mixed probability distributions like the examples  
264 discussed in Section 1. If the verification duplicates one or more ensemble members, its rank is  
265 no longer uniquely defined by (18). Suppose that  $D$  ensemble members have the same value as  
266 the verification. When these are removed from the ensemble, the rank of the verification in the  
267  $N + 1 - D$  remaining numbers is uniquely defined, even if there are other duplicate values in

268 the remaining ensemble; let that rank be  $R$ . The actual rank in the full ensemble could range  
269 from  $R$  to  $R + D$  since the order of the verification and its duplicates is not uniquely defined.  
270 Essentially, there is a  $1/(D + 1)$  probability that the rank of the verification is any of these  
271 values. In this case, define the weight vector as

$$272 \quad W_n = \begin{cases} 1/(D + 1) & \text{if } R \leq n \leq R + D \\ 0 & \text{otherwise} \end{cases} \quad (20)$$

273 A sum of rank weight vector can be defined for a collection of ensembles as before with (19),  
274 and the histogram should be uniform if the verifications are drawn from the same distribution  
275 as the ensemble members. Another possible way to define rank histograms for duplicate values  
276 is to randomly select one of the ranks between  $R$  and  $R+D$  and give it the weight of one,  
277 however this generates unnecessary random noise compared to the solution in (20).

278  
279 This treatment of duplicates for rank histograms is essential for application to state variables or  
280 true observations in an OSSE like the one in section 4. When rank histograms are used for  
281 verifications that are real observed quantities, it is important to account for observational error  
282 when generating an appropriate ensemble (Anderson 1996). One way to do this is to add a  
283 random sample from an observational error distribution to each ensemble member generated  
284 by applying a forward operator to the model state. In many cases, adding in this observation  
285 error component would eliminate duplicate values like those that result from bounded state  
286 variables in state space. However, if the error distribution is also mixed, duplicates are still  
287 expected. Note that a deterministic method similar to (20) can also be developed to account for  
288 observational error in the rank histogram.

289

#### 290 **4. A tracer advection extension of the Lorenz-96 Model: L96-T**

291

##### 292 *a. Model description*

293

294 A low-order model with sensitive dependence on initial conditions, low computational cost, and  
295 bounded state variables is useful for testing DA algorithms. The traditional Lorenz-96 model  
296 (Lorenz and Emanuel 1998) has been used in many ensemble DA studies including (A22). Here,

297 the Lorenz-96 model is extended to include two additional types of  $M$  variables that are  
 298 collocated with the standard variables,  $x_m, m = 0, \dots, M - 1$ , on the standard periodic domain.  
 299 The first type,  $q_m$ , represents concentrations of a dimensionless tracer. The second type,  $s_m$ ,  
 300 represents a source rate of the tracer with units of tracer amount per time. A function of the  
 301 standard  $x$  variables is treated as a wind field that passively advects the tracer. The velocity at  
 302 the model grid points at the current time is defined as  $v_m = \bar{V} + \tilde{V}x_m$  where  $\bar{V}$  is a specified  
 303 constant mean velocity,  $\tilde{V}$  is a specified multiplying constant that controls the average  
 304 magnitude of wind perturbations, and  $\tilde{V}x_m$  is an anomalous velocity at gridpoint  $m$ . Velocities  
 305 are expressed with units of nondimensional distance per nondimensional time. A  
 306 nondimensional location is assigned to each grid point in the model so that the distance  
 307 between neighboring grid points is 1 (note that this is different from many previous Lorenz-96  
 308 studies where the distance between grid points is defined as  $1/M$ ).

309  
 310 The time evolution of the standard variables,  $x_m$ , is identical to that used in the basic Lorenz-96  
 311 model (Lorenz and Emanuel 1996). The time evolution of the nonnegative tracer concentration  
 312 used here is:

$$313 \quad q_m^+ = \max[(q_m^{adv} + s_m \Delta t)e^{-E\Delta t} - C\Delta t, 0] \quad (21)$$

314 where  $q_m^+$  is the tracer concentration at grid point  $m$  at the next time step,  $q_m^{adv}$  is the advected  
 315 concentration,  $s_m$  is the source rate at grid point  $m$  at the current time,  $E$  is an exponential  
 316 damping time,  $C$  is a constant sink rate, and  $\Delta t$  is the timestep.

317  
 318 The advection of tracer is computed using an upstream semi-Lagrangian method. The  
 319 computation of  $q_m^{adv}$ , the advected concentration at the next time at grid point  $m$  given the  
 320 wind field at the current time,  $v_m$ , and the concentrations at the current time,  $q_m$ , proceeds as  
 321 follows:

- 322 1. A preliminary upstream target location is defined as  $T = m - v_m \Delta t$ ,
- 323 2. The fractional location of the target between the bounding grid points is  $p = T - [T]$   
 324 where the brackets indicate the floor,

- 325 3. The indices of the grid points bounding the target location are computed as  $L =$   
326  $\text{mod}(\lfloor T \rfloor, M)$  and  $U = \text{mod}(L + 1, M)$ ,  
327 4. The advected concentration is  $q_m^{adv} = (1 - p)q_L + pq_U$

328

329 The specified source is a function of grid point and model time with units of amount per time.  
330 For experiments here, there is a time constant source with rate 5 at grid point 1 and all other  
331 grid points have zero source at all times

$$332 s_m = \begin{cases} 5 & \text{if } m = 1 \\ 0 & \text{otherwise} \end{cases}$$

333

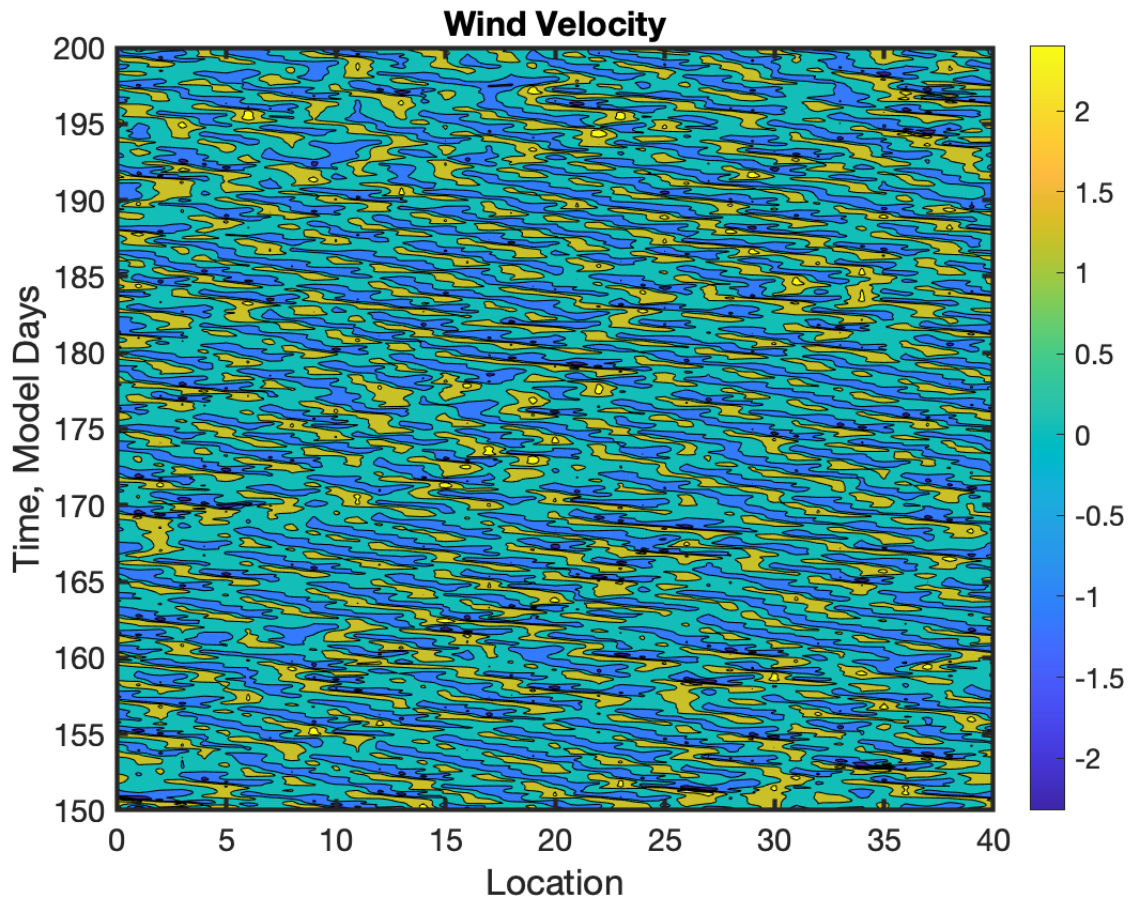
334 *b. L96-T example*

335

336 All results here use the standard 4<sup>th</sup> order Runge-Kutta time stepping algorithm, the  
337 nondimensional  $\Delta t = 0.05$  with an associated dimensional time step of 3600s as done in many  
338 previous studies, and  $M = 40$  grid points. The L96 forcing parameter  $F = 8$ . The mean velocity  
339  $\bar{V} = 0$  and the velocity perturbation multiplier  $\tilde{V} = 5$ , while the constant sink  $C = 0.1$ , and the  
340 exponential sink  $E = 0.25$ .

341

342 Figure 2 shows a time series of the wind field,  $v_m$ , as a function of the model grid point; since  
343  $\bar{V} = 0$ , this is just  $\tilde{V} = 5$  times the standard L96 state variables,  $x_m$ . The well-known group and  
344 phase velocity of the L96 model can be seen.



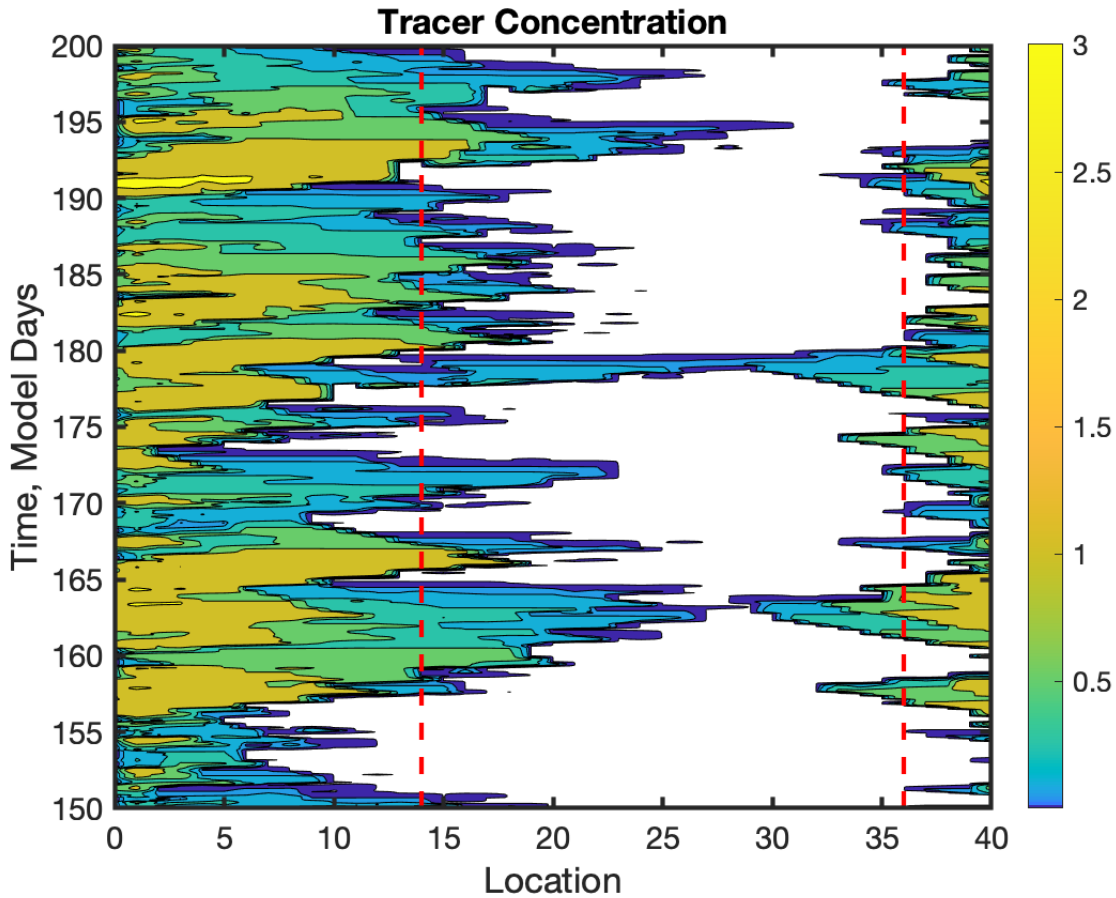
345

346 Figure 2: Wind velocity from the L96 Tracer Advection model for times between day 150 and  
 347 200 in the truth run as a function of model grid point. Units are distance  $\text{hr}^{-1}$ .

348

349 Figure 3 shows the tracer concentration corresponding to the wind field. Plumes of tracer are  
 350 advected away from the source at grid point 1. The velocity is positive more often than  
 351 negative, so plumes extend more frequently and further to the right. However, the wind field is  
 352 sometimes negative leading to shorter plumes extending to the left. The white areas in the plot  
 353 have zero tracer concentration, so a mixed distribution is appropriate. It is rare for plumes to  
 354 extend clear across the domain with this only happening twice in the figure. This behavior is  
 355 roughly analogous to what one might see with a point source in the midlatitudes. It is possible  
 356 to get a variety of other behaviors for the tracer by changing the model parameter values.

357



358

359 Fig. 3: As in figure 2 but for the tracer concentration (nondimensional). The red dashed lines  
 360 highlight grid points with additional diagnostics presented in figures 4, 5 and 6. White areas  
 361 have zero concentration.

362

### 363 5. Data assimilation experiments

364

365 The model integration described in the previous section is used as the truth run for a series of  
 366 observing simulation system experiments (OSSEs). The L96-T model is first integrated for 16500  
 367 hrs (5500 3-hour advances) starting from a default initial state to generate a tuning initial state.  
 368 The default initial state has  $x_1 = 1$  and all other  $x$  state variables are 0; all concentration  
 369 variables are 0. The model is integrated for an additional 16500 hrs from the tuning initial state  
 370 generating synthetic observations every 3 hrs. Forty randomly located observing sites are  
 371 selected for the L96 standard state, and a different set of 40 randomly located sites for tracer  
 372 concentration observations (see Figs 7d and 8d). Observations are taken by linearly



373 interpolating to the site location from the two nearest grid points. For the standard state  
374 observations error is simulated by adding a random draw from a normal distribution with mean  
375 0 and variance 10. For tracer observations error is simulated by adding a random draw from a  
376 truncated normal distribution with variance 0.1 and lower bound of 0 (A23, appendix D).

377

378 Three different observing networks are explored: assimilating only standard state observations,  
379 assimilating only tracer concentration observations, and assimilating both standard and tracer  
380 observations. Two different model configurations are evaluated. In the first, every ensemble  
381 member has the true value of the tracer source variables. In the second, the tracer source  
382 variables are unknown, and every ensemble member has its own (not necessarily unique) time  
383 evolving estimate.

384

385 All assimilation experiments use the adaptive inflation algorithm of Gharamti (2018) with an  
386 inflation damping of 0.9. All experiments also use a Gaspari Cohn localization with the same  
387 constant halfwidth for all observations. Seven halfwidth tuning assimilation experiments are  
388 done for each case, where a case is defined by the observing network, whether the source is  
389 known or unknown, and the ensemble size (20, 40, 80 or 160). As in A23, the halfwidths tested  
390 are  $\{0.075, 0.1, 0.125, 0.175, 0.2, 0.4, \infty\}$ . These tuning assimilations start from the tuning initial  
391 condition and assimilate for 5500 3-hour intervals. Initial ensembles for the standard state  
392 variable are generated by adding a random draw from a normal distribution with mean 0 and  
393 standard deviation 0.01 to the truth value for each variable. Initial ensemble members for the  
394 tracer variables are all equal to the truth. For the case with known sources, all ensemble  
395 members for the source variables are equal to the truth. For the case with unknown sources,  
396 ensemble members for the source are set to a random draw from a normal distribution with  
397 mean 2.5 and standard deviation 2.5; if the draw is less than 0 the source is set to 0 so that the  
398 resulting ensembles are generally mixed distributions with several members that are 0. Results  
399 from the first 500 assimilation steps are discarded and the prior ensemble mean, time mean  
400 RMSE is computed for the standard state and tracer variables for the remaining 5000 steps. For

401 the state only observing network, the localization that minimizes the state RMSE is selected; for  
 402 the other observing networks, the localization that minimizes the tracer RMSE is selected.

403  
 404 The model truth is then integrated for an additional 16500 hours from the end of the tuning  
 405 integration with synthetic observations generated in the same way. Initial conditions for  
 406 ensembles are also generated in the same way as for the tuning experiments. Assimilation  
 407 experiments are performed for each case using the tuned localization and assimilating every 3  
 408 hours. The first 500 steps are discarded, and results are available for the final 5000 assimilation  
 409 steps. The spread for all quantities appears to be spun up after fewer than 100 assimilation  
 410 steps for all experiments.

411  
 412 Four different assimilation algorithms are applied to each case using the QCEFF. As noted in  
 413 A23, a complete description of a QCEFF assimilation algorithm requires information about the  
 414 first step where increments are computed for observed variables and the second step where  
 415 those increments are regressed onto state variables. The QCEFF uses a probit probability  
 416 integral transform (PPI) before doing the regression (A23). Table 1 specifies the details of the  
 417 four algorithms which are referred to as an EAKF, RHF, PQBNRH, and DUAL. Note that the  
 418 normal likelihood used for the  $q$  variable in the EAKF is a normal with the same variance as the  
 419 truncated normal observation error distribution for  $q$ . As noted in A23, using a normal for the  
 420 PPI transform is equivalent to no transform. The BNRH CDFs all have a lower bound of 0 and no  
 421 upper bound, consistent with the nature of the tracer concentration and source variables.

422

	EAKF	RHF	PQBNRH	DUAL
x obs. CDF	Normal	RH	RH	Normal
x likelihood	Normal	Normal	Normal	Normal
x PPI CDF	None	None	RH	None
q obs. CDF	Normal	BNRH	BNRH	BNRH
q likelihood	Normal	Truncated Normal	Truncated Normal	Truncated Normal
q PPI CDF	None	None	BNRH	BNRH

s PPI CDF	None	None	BNRH	BNRH
-----------	------	------	------	------

423

424 Table 1: Assimilation settings for each of the four algorithms explored. For the  $x$  and  $q$   
 425 variables, the continuous CDF and form of the likelihood used for computing observation space  
 426 increments are listed with RH referring to a rank histogram distribution without bounds and  
 427 BNRH referring to a bounded normal rank histogram distribution with a lower bound at zero.  
 428 The continuous distribution used as part of the PPI transform used when regressing observation  
 429 increments onto state variable increments is also listed.

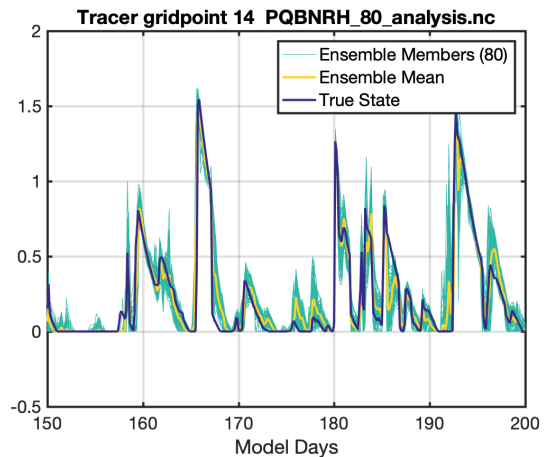
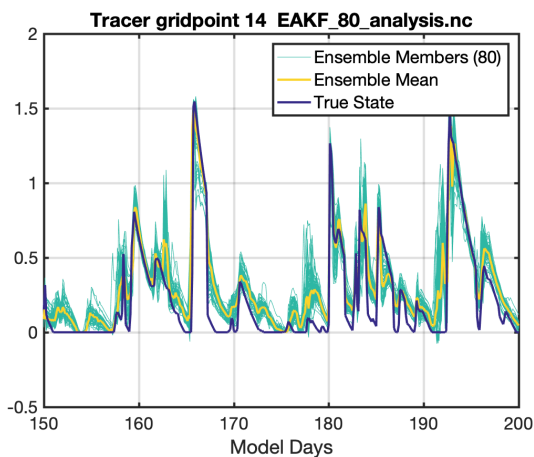
430

431 *a. Known source results*

432

433 Unless otherwise noted, all results shown are for analysis, rather than forecast, variables. Also,  
 434 results shown are for the network observing both standard state and tracer observations unless  
 435 otherwise noted. Figure 4 shows a time series from the EAKF and PQBNRH algorithm 80-  
 436 member assimilations for tracer at grid point 14, which is highlighted by a red dashed line in Fig.  
 437 3. The EAKF ensemble in Fig. 4a represents all the plumes that occur, but also represents two  
 438 plumes between days 150 and 160 that are not real. The ensemble is strongly biased towards  
 439 larger values at some times, in particular around days 168, 173, and 178. The PQBNRH results in  
 440 Fig. 4b also capture all real plumes with smaller values for the two false plumes, but do not  
 441 have the strongly biased periods.

442

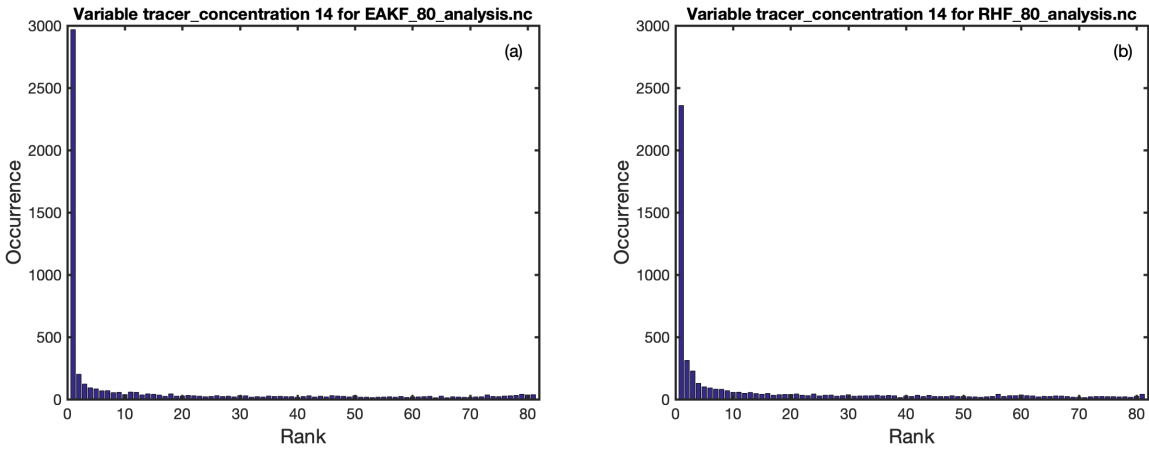


443

444 Fig. 4: Time series of the tracer at grid point 14. Dark curve is the truth and is the same in both  
445 panels. The dark green curves are the 80 analysis ensemble members, and their mean is in  
446 yellow, for an EAKF (left) and a PQBNRH (right); the tracer is nondimensional.

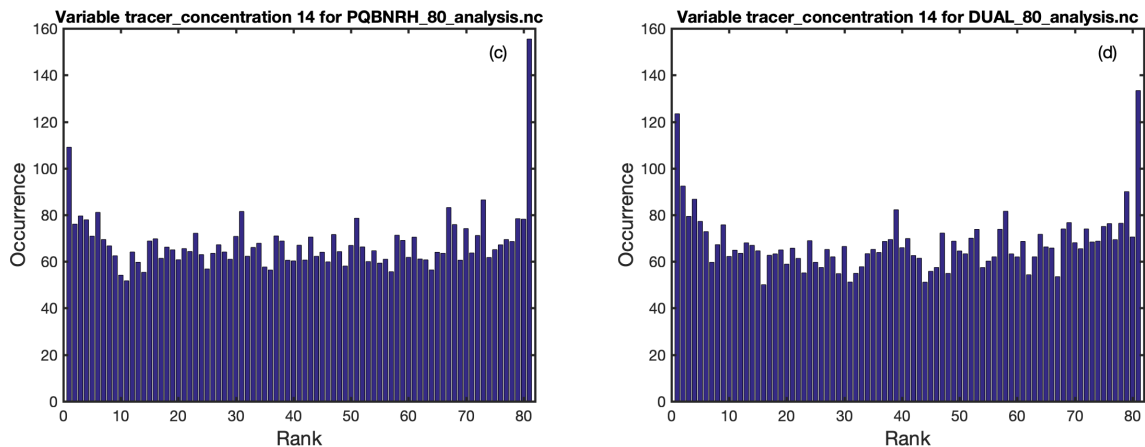
447  
448 Figure 5 shows rank histograms over all 5000 assimilation steps for concentration at grid point  
449 14. The EAKF and RHF algorithms result in very strongly biased histograms with the truth very  
450 often less than the smallest ensemble member. The results for the PQBNRH and DUAL are  
451 radically different. Both have histograms that are nearly uniform except for the two outermost  
452 bins. The PQBNRH has more cases where the truth is larger than the largest ensemble member  
453 while the DUAL algorithm has more cases where the truth is smaller than the smallest member;  
454 however, it is difficult to evaluate whether these differences are statistically significant.

455



456

457

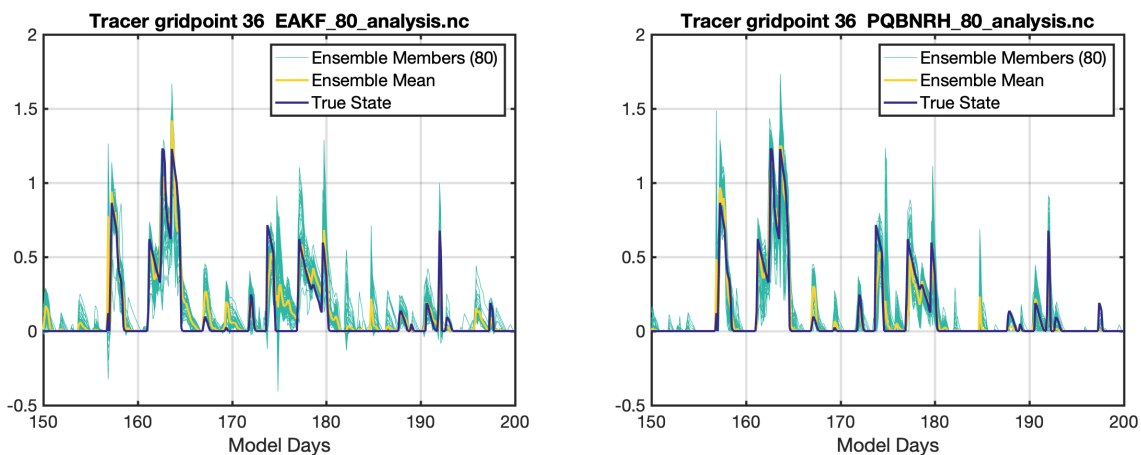


458

459

460 Fig. 5: Rank histograms for 80-member analysis concentration at grid point 14 for an EAKF (a),  
461 RHF (b), PQBNRH (c) and a DUAL filter with an EAKF for the wind and a BNRH for the  
462 concentration (d). Note the different vertical axes in the top and bottom rows.  
463

464 Fig. 6 shows time series of the EAKF and PQBNRH assimilation results for grid point 36 which is  
465 also highlighted in Fig. 3. At this grid point, plumes are less frequent, primarily arriving from the  
466 right. There are extended periods when the true concentration is 0. The EAKF represents all  
467 true plumes, however, there are several instances where the ensemble is strongly biased  
468 towards larger concentration than the truth, and several times when negative ensemble  
469 members occur; this cannot happen with the PQBNRH. The EAKF never has any ensemble  
470 members that are exactly zero and never has duplicate ensemble members. The PQBNRH also  
471 captures all real plumes and has fewer instances of false plumes. The PQBNRH has several  
472 periods when many ensemble members are exactly 0 and some periods where all members are  
473 zero, all at times when the truth is also zero. Results for the RHF are similar to those for the  
474 EAKF, and results for DUAL are similar to those for the PQBNRH in figures 4 and 6 so these are  
475 not displayed.  
476

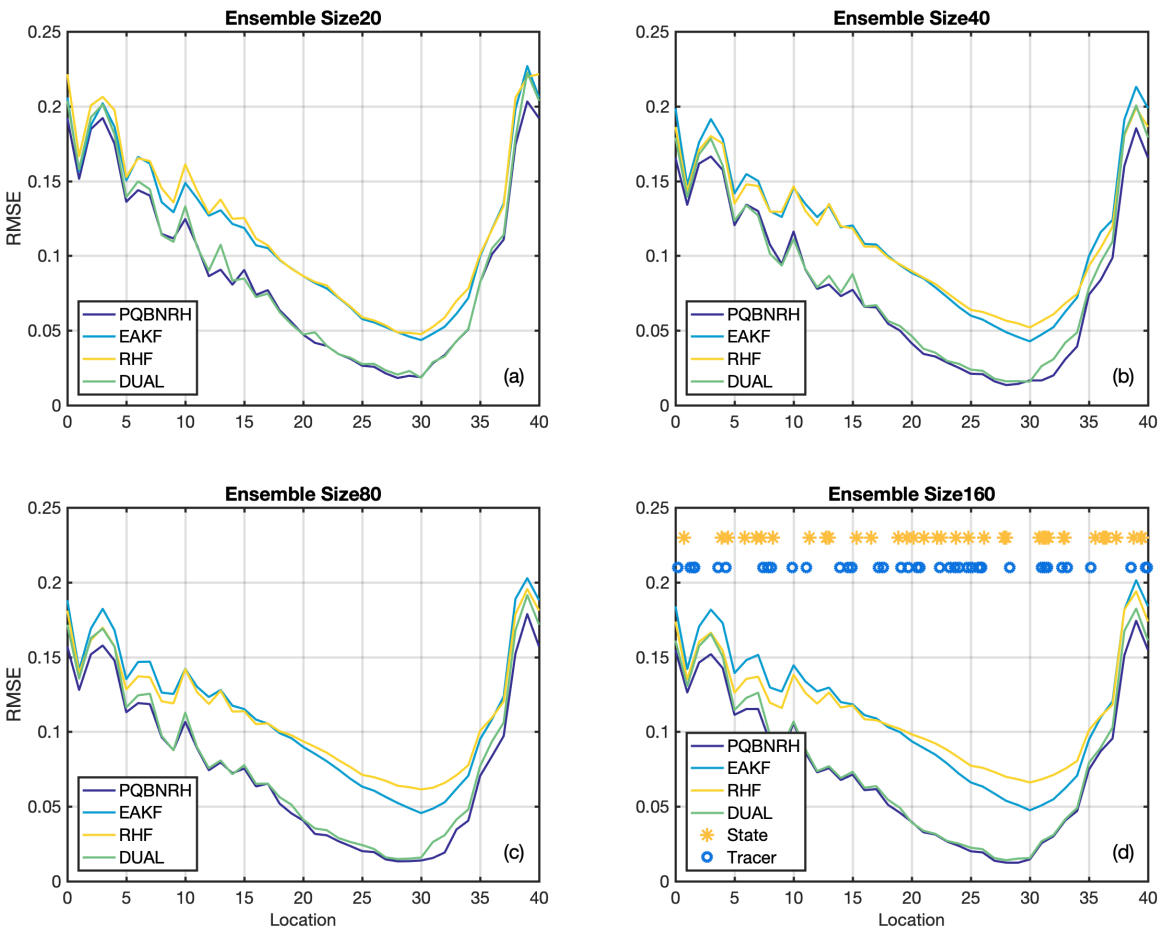


477  
478 Fig. 6: As in figure 4 but for grid point 36.

479  
480 Fig. 7 shows summary results for ensemble mean tracer RMSE over all 5000 assimilation steps  
481 for the four algorithms and four ensemble sizes studied. In general, the results for the PQBNRH  
482 and DUAL algorithms are statistically indistinguishable. The same is true for the EAKF and RHF

483 algorithms. However, in general the PQBNRH/DUAL algorithms have lower RMSE. The RMSE is  
 484 largest to the left of the source at grid point 0 where the true concentration is most variable,  
 485 and smaller far from the source where concentration is smaller. There are not large differences  
 486 as a function of ensemble size; larger ensembles generally have only slightly smaller RMSE. It is  
 487 unclear why ensemble size is not more important here.

488



489

490 Fig. 7: Ensemble mean, time mean RMSE as a function of grid point for the analysis tracer  
 491 concentration for four filter algorithms for ensemble size 20 (a), 40 (b), 80 (c) and 160 (d). The  
 492 locations of the 40 observing stations are shown in (d) for state (yellow circles) and tracer  
 493 concentration (blue asterisks).  
 494

495 It is obvious that assimilating standard state observations that improve the estimate of the  
 496 winds will result in improved estimates of the tracer concentrations. However, the impact of  
 497 tracer observations on the standard state variables is less clear. Assimilations for the network  
 498 observing only tracer produced tracer analysis estimates that have much larger RMSE than

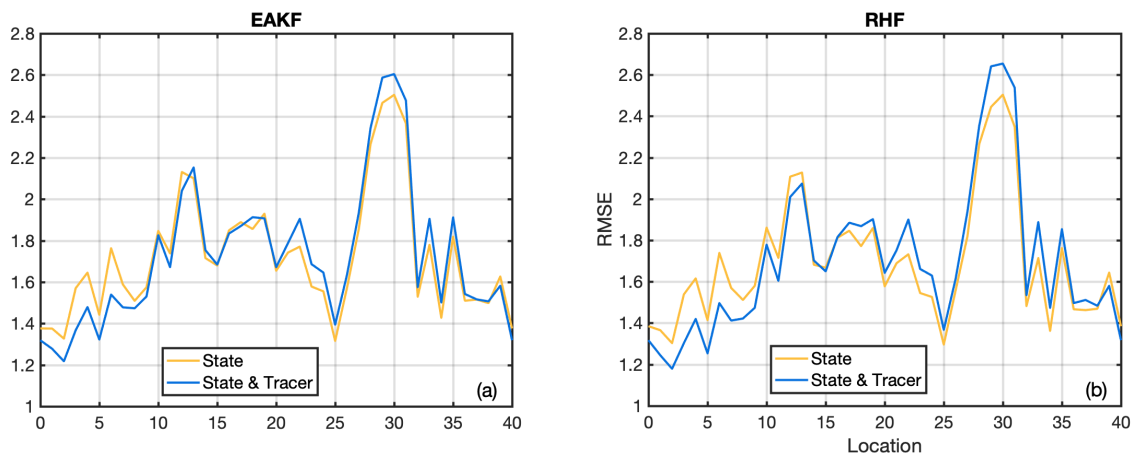
499 those just discussed, although smaller than the RMSE from an unconstrained control ensemble  
500 run. The tracer only network resulted in standard variable RMSE that was only slightly smaller  
501 than the RMSE from an unconstrained control.

502

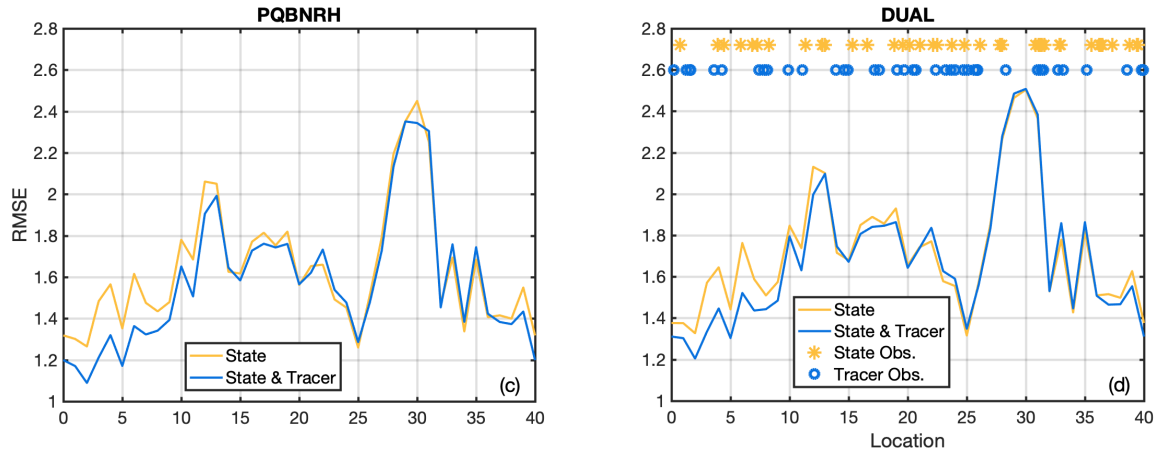
503 A comparison of the standard variable RMSE from the observing network with only standard  
504 state observations to the network with both standard and tracer observations is shown in Fig. 8  
505 for the four algorithms. The RMSE for the standard observation only network has larger RMSE  
506 near grid point 30 and smaller RMSE near grid points 25 and 1. This is due to the random  
507 observing site locations (Fig. 8d). The RMSE is smaller for the PQBNRH than for any of the other  
508 algorithms; note that the EAKF and DUAL are identical for the standard observation network.

509

510 When tracer observations are added in, all four algorithms produce reduced RMSE for the left  
511 part of the domain. The EAKF and RHF produce increased RMSE in the right part of the domain.  
512 The PQBNRH and DUAL produce roughly equivalent RMSE in the right part of the domain. In the  
513 left part of the domain, plumes with large spatial and temporal gradients occur near the source.  
514 These provide information about the flow field that is advecting the plume and lead to the  
515 reduced RMSE for the standard state. Because there is often very little or no tracer in the right  
516 part of the domain, observations of the tracer are expected to provide very little additional  
517 information. The increase in error in the EAKF and RHF suggests that deficiencies in these  
518 algorithms cause the use of these low information observations to degrade the ensemble  
519 estimate.



520



521  
522

523 Fig 8: Ensemble mean, time mean RMSE as a function of grid point for the standard L96 state  
 524 for experiments that assimilate only observations of the standard state, and experiments that  
 525 also assimilate the tracer concentration, shown for an EAKF (a), RHF (b), PQBNRH (c), and a  
 526 DUAL filter with an EAKF for the wind and a BNRH for the concentration (d). The locations of  
 527 the 40 observing stations are shown in (d) for state (yellow asterisks) and tracer concentration  
 528 (blue circles).

529

530 *b. Unknown source results*

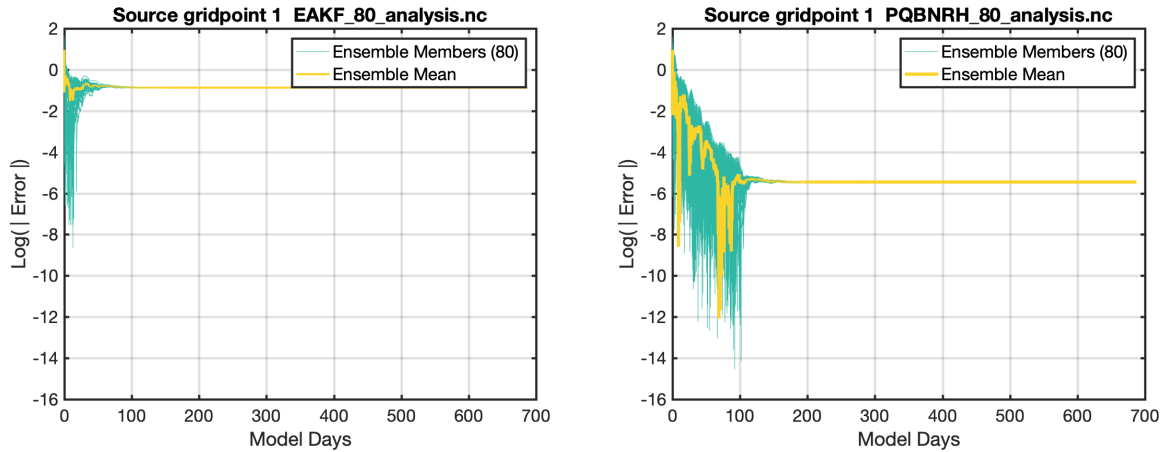
531

532 In these experiments, the source is not known and is estimated by the assimilation algorithms.  
 533 Results are only discussed for the network observing both standard state and tracer  
 534 observations. There is no time tendency model for the tracer. The prior ensembles can have  
 535 their spread increased by the adaptive inflation. Nevertheless, in all experiments, the spread  
 536 becomes increasingly small for the source at all grid points. The source variables are only  
 537 impacted by concentration observations since the source and the state field should not be  
 538 meaningfully correlated.

539

540 Figure 9 shows the natural logarithm of the absolute value of the error for each ensemble  
 541 member and the ensemble mean error at the grid point with the nonzero source in the truth  
 542 for the EAKF and the PQBNRH. Both reduce the ensemble mean error, but the reduction is  
 543 much larger for the PQBNRH. Because of the collapse of spread, both algorithms eventually  
 544 have strongly biased estimates and would produce corresponding rank histograms.





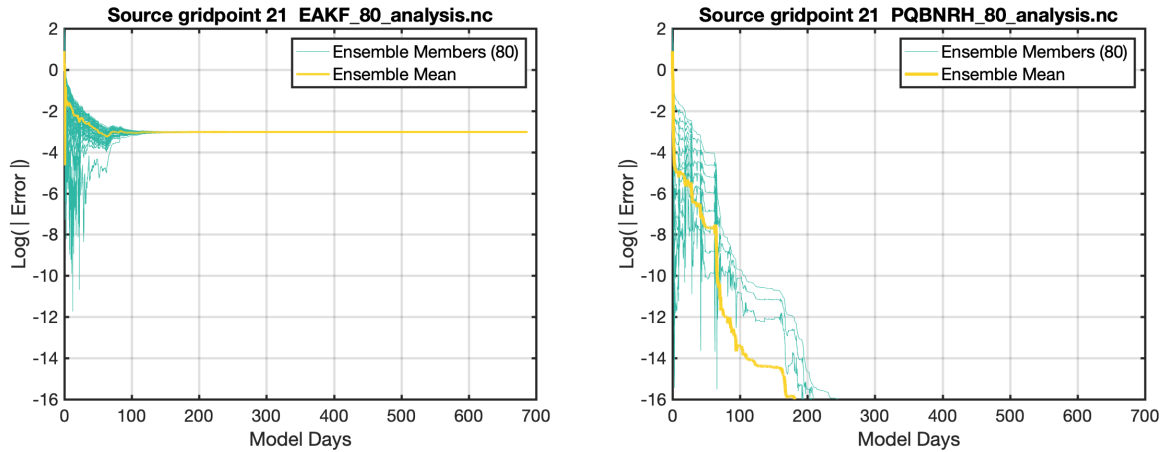
545  
546

547 Fig. 9: Spatial mean of the natural logarithm of the absolute value of the error of the ensemble  
548 mean (yellow) and each of the 80 ensemble members (green) as a function of time for the  
549 source at grid point 1 which has a true value of 5 (units  $\text{hr}^{-1}$ ) for the EAKF (left) and the  
550 PQBNRH (right).

551

552 Figure 10 shows the evolution of the RMSE for grid point 21 which has zero true source. The  
553 RMSE for the EAKF is smaller than it was for grid point 1. The error for the PQBNRH decreases  
554 throughout the 5000 assimilation steps. As the assimilation continues, more and more  
555 ensemble members have values of exactly zero; eventually all ensemble members are zero and  
556 the error of the ensemble mean, and all individual ensembles is zero. At both grid points 1 and  
557 21, the RMSE for the standard state and concentration variables for the PQBNRH are nearly  
558 identical to those for the known source experiments since the source is so accurately  
559 estimated. Results are somewhat degraded for the EAKF which has larger errors in its source  
560 estimates.

561



562

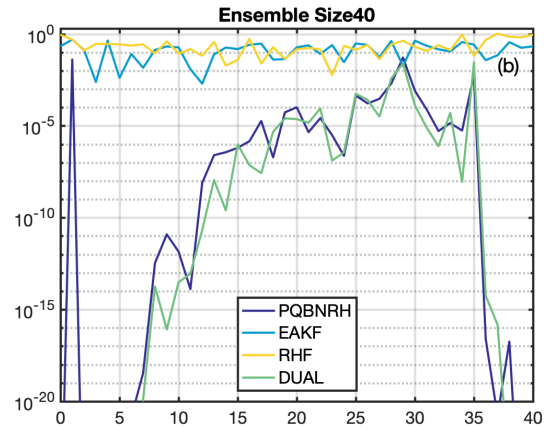
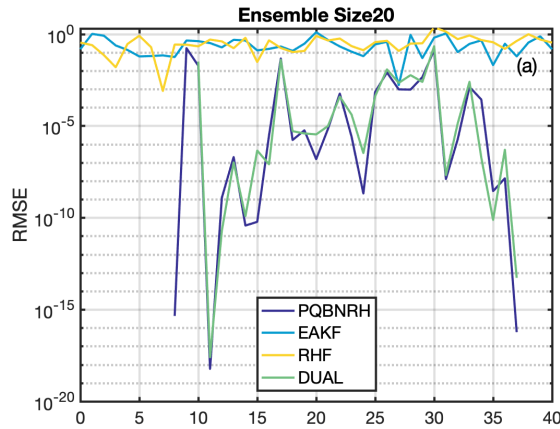
563 Fig. 10: As in 9 for grid point 21 which has zero true source.

564

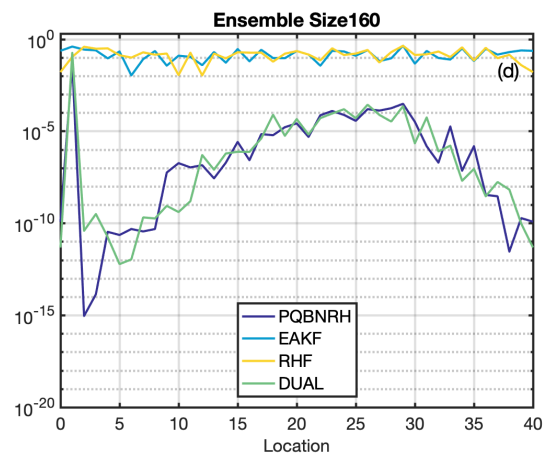
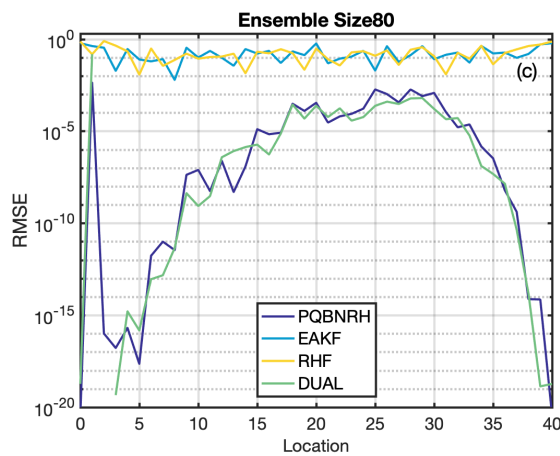
565 Figure 11 shows the RMSE for the source as a function of grid point for each of the four  
 566 algorithms and all four ensemble sizes. The EAKF and RHF produce roughly comparable results  
 567 that have a small dependence on ensemble size. The errors do not have a strong dependence  
 568 on the grid point. The PQBNRH and DUAL are also very similar but have more dependence on  
 569 both ensemble size and grid point. The smallest errors occur for grid points close to the non-  
 570 zero source at grid point 1. The RMSE actually increases with ensemble size in these areas. This  
 571 is due to the rate at which ensemble members become exactly zero which appears to be similar  
 572 across ensembles so that the fraction of nonzero members at a given time increases with  
 573 ensemble size. Larger errors are found for the source point itself and for points far from the  
 574 source. The estimate at point 1 varies little with ensemble size. The RMSE for points remote  
 575 from the source gets smaller and less noisy with increasing ensemble size.

576

577



578



579 Fig. 11: Ensemble mean time mean RMSE as function of grid point for tracer source for four  
580 algorithms for ensemble size 20 (a), 40 (b), 80 (c) and 160 (d). Values that are not plotted for  
581 the PQBNRH and DUAL algorithms are less than  $10^{-20}$  including many that are exactly zero.

582

583

## 584 6. Discussion and conclusions

585

586 The QCEFF has been extended to deal with model and observed variables with mixed  
587 probability distributions. This capability is especially relevant for bounded quantities like  
588 precipitation (Lien et al. 2013), tracer concentrations and sources, and areal coverage (Wieringa  
589 et al., 2023 in press; Riedel and Anderson 2023 in press). It may also be useful for estimating  
590 model parameters with data assimilation (Gharamti et al. 2016); the tracer source in the L96-T  
591 version used here is essentially equivalent to a model parameter.

592

593 The nearly non-parametric BNRH distribution has also been extended to handle duplicate  
594 ensemble members that are expected to occur for variables with mixed distributions. The rank  
595 histogram diagnostic tool was also extended to deal with duplicate ensemble members. An  
596 extension of the Lorenz-96 low-order dynamical system that includes an idealized advective  
597 tracer with local sources was developed to test the new algorithms. This L96-T model should  
598 also provide challenging tests for other data assimilation algorithms including variational  
599 methods and particle filters.

600

601 Results show that the BNRH works better than the EAKF or RHF for an OSSE with the L96-T  
602 model. The RMSE is smaller for the bounded tracer concentration and source variables when  
603 they are close to the bounds as might be expected. Results are also better when these variables  
604 are not close to the bounds and for the unbounded standard state variables from L96. The RH  
605 and PQBNRH algorithms use the BNRH distribution to compute observation increments.  
606 However, the RH uses standard linear regression when updating state variables while the  
607 PQBNRH includes the PPI transform using the BNRH distribution for the probability integral  
608 transform. The RH results are similar to the EAKF results in this case, while the PQBNRH is  
609 better for all variables and locations demonstrating that the PPI is a crucial part of the improved  
610 performance. The DUAL case uses an EAKF for the L96 state which has no bounds and is  
611 expected to be approximately normal. There is some indication that the PQBNRH is slightly  
612 better than the DUAL algorithm, but differences are not quantitatively significant. This suggests  
613 a strategy of using the BNRH distribution for bounded variables but a normal distribution for  
614 other variables may be useful for large model applications.

615

616 The BNRH as described allows duplicate ensemble members and the data assimilation process  
617 can create additional duplicates; this happened for both concentration and source variable  
618 ensembles in the OSSEs here. However, the assimilation process cannot eliminate duplicates. It  
619 can change the value of ensemble members that are exactly at a bound in the prior ensemble.  
620 This means that the model must eliminate duplicates if appropriate. That happens in  
621 experiments here and is most clearly seen in figure 6b where all ensemble members are zero at

622 some times but not at subsequent times. Further investigation into methods that would allow  
623 the assimilation to eliminate duplicates is warranted but would require making a priori  
624 assumptions about the expected errors associated with a given ensemble size.

625

626 All the OSSEs here were performed using the Data Assimilation Research Testbed (DART:  
627 Anderson et al. 2009) which implements the QCEFF including the BNRH; the parallel algorithm  
628 of Anderson and Collins (2007) was used. The results here only examined the use of normal or  
629 BNRH distributions. DART software can support arbitrary distributions and currently supports  
630 gamma, inverse gamma, log-normal, beta, and particle filter distributions. Previous work on  
631 assimilation of bounded quantities has proposed using distributions like the log-normal,  
632 gamma, and inverse gamma. However, the L96-T OSSE explored here presents specific  
633 challenges for using these other distributions. The log-normal and inverse gamma distributions  
634 do not have any probability at zero. This is clearly inappropriate for the mixed distributions in  
635 the OSSE where much of the probability can be at 0 at some times. The gamma distribution can  
636 have probability at zero. However, if the likelihood is a gamma distribution, the corresponding  
637 observation error distribution is inverse gamma (Bishop 2016, A22). This means that  
638 observations of the bounded quantities would not be able to have any probability at zero. This  
639 is clearly problematic for the bounded quantities with realistic instruments. Further work on  
640 explicitly using mixed distributions, for instance a combination of a log-normal with a discrete  
641 distribution, for applications like this is beyond the scope of this report.

642

643 The computational cost of the QCEFF algorithms including the BNRH is discussed in detail in  
644 A23. There is almost no additional cost associated with allowing duplicate ensemble members  
645 so the A23 analysis still applies. As noted there, the additional cost of a BNRH compared to an  
646 EAKF can be significant, especially for low-order model applications. As discussed in Anderson  
647 (2019) and A23, much of this cost is associated with the need to sort the ensemble members  
648 for each state variable. However, the sorting order often changes little between assimilation  
649 steps. Caching the sort order and then using sorts that are efficient for nearly sorted data can

650 potentially result in large computational cost reductions, but these methods have not yet been  
651 implemented in DART.

652

653 The low-order model results here suggest that there may be significant improvements when  
654 the BNRH is used for bounded quantities in large Earth system applications. Initial tests in sea  
655 ice (Wieringa et al. 2023 in press) and chemical transport models will be investigated in  
656 subsequent work. Other types of nearly non-parametric distributions, for instance various  
657 kernels (Grooms 2022, Anderson and Anderson 1999) can also be developed in DART and  
658 should be compared to the BNRH results.

659

660 *Acknowledgements.* This material is based upon work supported by the National Center for  
661 Atmospheric Research, which is a major facility sponsored by the National Science Foundation  
662 under Cooperative Agreement 1852977. Any opinions, findings, and conclusions or  
663 recommendations expressed in this publication are those of the author and do not necessarily  
664 reflect the views of the National Science Foundation. Thanks to Ian Grooms, Joseph Chan,  
665 Hristo Chipilski, Ben Gaubert and the DAREs team for helpful discussions about this material.

666

667 Data availability statement

668 The Lorenz-96 results were generated with DART code that can be found at:

669 [https://github.com/NCAR/DART/releases/tag/MWR\\_QCEFF\\_Part3](https://github.com/NCAR/DART/releases/tag/MWR_QCEFF_Part3).

670

671 REFERENCES

672

673 Amezcua, J. and P. J. Van Leeuwen, 2014: Gaussian anamorphosis in the analysis step of the  
674 EnKF: a joint state-variable/observation approach. *Tellus A: Dynamic Meteorology and*  
675 *Oceanography*, 66. DOI: [10.3402/tellusa.v66.23493](https://doi.org/10.3402/tellusa.v66.23493)

676 Anderson, J. L., 1996: A Method for Producing and Evaluating Probabilistic Forecasts from  
677 Ensemble Model Integrations. *J. Climate*, 9, 1518–1530, [https://doi.org/10.1175/1520-0442\(1996\)009<1518:AMFPAE>2.0.CO;2](https://doi.org/10.1175/1520-0442(1996)009<1518:AMFPAE>2.0.CO;2).

678

679 Anderson, J. L., 2001: An ensemble adjustment Kalman filter for data assimilation. *Mon. Wea.*  
680 *Rev.*, **129**, 2884-2903.

681 Anderson, J. L., 2003: A local least squares framework for ensemble filtering. *Mon. Wea. Rev.*,  
682 **131**, 634-642.

683 Anderson, J. L., 2010: A non-Gaussian ensemble filter update for data assimilation. *Mon. Wea.*  
684 *Rev.*, **138**, 4186–4198, <https://doi.org/10.1175/2010MWR3253.1>

685 Anderson, J. L., 2019: A Nonlinear Rank Regression Method for Ensemble Kalman Filter Data  
686 Assimilation. *Mon. Wea. Rev.*, **147**, 2847–2860, [https://doi.org/10.1175/MWR-D-18-](https://doi.org/10.1175/MWR-D-18-0448.1)  
687 [0448.1](https://doi.org/10.1175/MWR-D-18-0448.1).

688 Anderson, J. L., 2022: A Quantile-Conserving Ensemble Filter Framework. Part I: Updating an  
689 Observed Variable. *Mon. Wea. Rev.*, **150**, 1061–1074, [https://doi.org/10.1175/MWR-D-](https://doi.org/10.1175/MWR-D-21-0229.1)  
690 [21-0229.1](https://doi.org/10.1175/MWR-D-21-0229.1).

691 Anderson, J. L., 2023: A Quantile-Conserving Ensemble Filter Framework. Part 2: Updating an  
692 Observed Variable. *Mon. Wea. Rev.*, **151**, 2759–2777, [https://doi.org/10.1175/MWR-D-](https://doi.org/10.1175/MWR-D-23-0065.1)  
693 [23-0065.1](https://doi.org/10.1175/MWR-D-23-0065.1)

694 Anderson J. L., and S. L. Anderson, 1999: A Monte Carlo implementation of the nonlinear  
695 filtering problem to produce ensemble assimilations and forecasts. *Mon. Wea. Rev.*, **127**,  
696 2741–2758.

697 Anderson, J. L., and N. Collins, 2007: Scalable implementations of ensemble filter algorithms  
698 for data assimilation. *J. Atmos. Oceanic Technol.*, **24**, 1452–1463.

699 Anderson, J., T. Hoar, K. Raeder, H. Liu, N. Collins, R. Torn and A. Arellano, 2009: The Data  
700 Assimilation Research Testbed. *Bul. Amer. Met. Soc.*, **90**, 1283-1296.

701 Bannister, R. N., H. G. Chipilski, and O. Martinez-Alvarado, O., 2020. Techniques and  
702 challenges in the assimilation of atmospheric water observations for numerical weather  
703 prediction towards convective scales. *Q J R Meteorol Soc.*, **146**. [https://doi-](https://doi.org/cuucar.idm.oclc.org/10.1002/qj.3652)  
704 [org.cuucar.idm.oclc.org/10.1002/qj.3652](https://doi.org/cuucar.idm.oclc.org/10.1002/qj.3652)

705 Beal, D., P. Brasseur, J. M. Brankart, Y. Ourmieres, and J. Verron, 2010: Characterization of  
706 mixing errors in a coupled physical biogeochemical model of the North Atlantic:

707 Implications for nonlinear estimation using Gaussian anamorphosis. *Ocean Sci.*, **6**, 247–  
708 262.

709 Bertino, L., G. Evensen and H. Wackernagel, 2003: Sequential Data Assimilation Techniques in  
710 Oceanography. *International Statistical Review*, **71**, 223-241.

711 Bishop, C. H., 2016: The GIGG-EnKF: Ensemble Kalman filtering for highly skewed non-negative  
712 uncertainty distributions. *Q. J. R. Meteorol. Soc.*, **142**, 1395-1412. doi:[10.1002/qj.2742](https://doi.org/10.1002/qj.2742)

713 Bocquet, M., C. A. Pires, and L. Wu, 2010: Beyond Gaussian Statistical Modeling in Geophysical  
714 Data Assimilation. *Mon. Wea. Rev.*, **138**, 2997-3023.  
715 <https://doi.org/10.1175/2010MWR3164.1>.

716 Burgers, G., P. J. van Leeuwen, and G. Evensen, 1998: Analysis scheme in the ensemble  
717 Kalman filter. *Mon. Wea. Rev.*, **126**, 1719-1724.

718 Chan, M., J. L. Anderson, and X. Chen, 2020: An Efficient Bi-Gaussian Ensemble Kalman Filter  
719 for Satellite Infrared Radiance Data Assimilation. *Mon. Wea. Rev.*, **148**, 5087–  
720 5104, <https://doi.org/10.1175/MWR-D-20-0142.1>.

721 Doron, M., P. Brasseur, J. M. Brankart, S. N. Losa, and A. Melet, 2013: Stochastic estimation of  
722 biogeochemical parameters from Globcolour ocean colour satellite data in a North  
723 Atlantic 3D ocean coupled physical–biogeochemical model. *J. Marine Systems*, **117**, 81–  
724 95.

725 Fletcher, S., and M. Zupanski, 2006: A data assimilation method for log-normally distributed  
726 observational errors. *Quart. J. Roy. Met. Soc.*, **132**, 2505 - 2519. 10.1256/qj.05.222.

727 Gharamti, M., 2018: Enhanced Adaptive Inflation Algorithm for Ensemble Filters. *Mon. Wea.*  
728 *Rev.*, **146**, 623–640, <https://doi.org/10.1175/MWR-D-17-0187.1>.

729 Gharamti, M. E., A. Samuelsen, L. Bertino, E. Simon, A. Korosov, and U. Daewel, 2016: Online  
730 tuning of ocean biogeochemical model parameters using ensemble estimation  
731 techniques: Application to a one-dimensional model in the North Atlantic. *Journal of*  
732 *Marine Systems*, 168, 1-16. <https://doi.org/10.1016/j.jmarsys.2016.12.003>.

733 Grooms, I., 2022: A comparison of nonlinear extensions to the ensemble Kalman  
734 filter. *Comput Geosci* **26**, 633–650. <https://doi.org/10.1007/s10596-022-10141-x>



735 Hamill, T. M., 2001: Interpretation of Rank Histograms for Verifying Ensemble Forecasts. *Mon.*  
736 *Wea. Rev.*, **129**, 550–560, [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0493(2001)129<0550:IORHFV>2.0.CO;2)  
737 [0493\(2001\)129<0550:IORHFV>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0550:IORHFV>2.0.CO;2).

738 Houtekamer, P. L., and H. L. Mitchell, 1998: Data assimilation using an ensemble Kalman filter  
739 technique. *Mon. Wea. Rev.*, **126**, 796-811.

740 Houtekamer, P. L., and F. Zhang, 2016: Review of the ensemble Kalman filter for atmospheric  
741 data assimilation, *Mon. Wea. Rev.*, **144**, 4489-4532.

742 Kurosawa, K., and J. Poterjoy, 2021: Data Assimilation Challenges Posed by Nonlinear  
743 Operators: A Comparative Study of Ensemble and Variational Filters and Smoothers. *Mon.*  
744 *Wea. Rev.*, **149**, 2369–2389, <https://doi.org/10.1175/MWR-D-20-0368.1>.

745 Lien, G.Y., E. Kalnay, and T. Miyoshi, 2013: Effective assimilation of global precipitation:  
746 Simulation experiments. *Tellus A* 65, DOI: [10.3402/tellusa.v65i0.19915](https://doi.org/10.3402/tellusa.v65i0.19915)

747 Lorenz, E. N. and K. A. Emanuel, 1998: Optimal sites for supplementary weather observations:  
748 Simulation with a small model. *J. Atmos. Sci.*, **55**, 399-414.

749 Pham, D.T., 2001: Stochastic methods for sequential data assimilation in strongly nonlinear  
750 systems. *Mon. Wea. Rev.*, **129**, 1194–1207. Doi:10.1175/1520-0493

751 Riedel, C., 2023: In press

752 Simon, E. , and L. Bertino, 2012: Gaussian anamorphosis extension of the DEnKF for combined  
753 state parameter estimation: application to a 1D ocean ecosystem model. *J. Marine Syst.*,  
754 **89**, 1–18.

755 Suhaila, J., K. Ching-Yee, Y. Fadhilah and F. Hui-Mean, 2011: Introducing the Mixed  
756 Distribution in Fitting Rainfall Data. *Open Journal of Modern Hydrology*, **1**, 11-22.  
757 doi: [10.4236/ojmh.2011.12002](https://doi.org/10.4236/ojmh.2011.12002).

758 Van Leeuwen, P. J., 2009: Particle filtering in geophysical systems. *Mon. Wea. Rev.*, **137**, 4089-  
759 4114. [doi.org/10.1175/2009MWR2835.1](https://doi.org/10.1175/2009MWR2835.1)

760 Van Leeuwen, P. J., H. R. Künsch, L. Nerger, R. Potthast, and S. Reich, 2019: Particle filters for  
761 high-dimensional geoscience applications: A review. *Quart. J. Roy. Met. Soc.*, **149**, 2335-  
762 2365. [doi.org/10.1002/qj.3551](https://doi.org/10.1002/qj.3551)

763 Wilks, D. S., 2019: Indices of Rank Histogram Flatness and Their Sampling Properties. *Mon.*  
764 *Wea. Rev.*, **147**, 763–769, <https://doi.org/10.1175/MWR-D-18-0369.1>.