

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^a, Chris Riedel^b, Molly Wieringa^c, Fairuz Ishraque^d, Marlee Smith^a, Helen
4 Kershaw^a

5
6 ^a NSF NCAR/CISL/TDD/DAReS, Boulder, Colorado, jla@ucar.edu

7 ^b Cooperative Programs for the Advancement of Earth System Science, University
8 Corporation for Atmospheric Research

9 ^c University of Washington, Department of Atmospheric Sciences, Seattle, Washington

10 ^d Department of Geosciences, Princeton University, Princeton, New Jersey

11

12 Submitted to Monthly Weather Review, November 2023

13 Revised March 2024

14

15 **Abstract**

16

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

34

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

43

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

45

46 **1. Introduction**

47

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

57

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

70

71 The uncertainty for some variables is a mixed probability distribution that includes both
72 discrete and continuous parts. As an example, the amount of precipitation that falls during a

73 particular period (Suhaila et al. 2011) might have a discrete probability of being exactly zero in
74 addition to a continuous distribution of being non-zero; the precipitation rate would have a
75 similar mixed distribution. The amount of sea ice, snow cover, chemical tracer, or water in a
76 stream also have mixed distributions along with their source and sink rates. Quantities like the
77 fractional coverage of ice or snow are doubly bounded, and could have a discrete probability of
78 no cover, a discrete probability of complete coverage, and a continuous distribution for all
79 intermediate values. A beta distribution might be appropriate for some doubly bounded
80 quantities.

81

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

88

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

101

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

117

118 **2. QCEFF for discrete and mixed probability distributions**

119

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

130

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

$$133 \quad \sum_{i=1}^K p_i = 1. \quad (1)$$

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

$$137 \quad S = \sum_{i=1}^K L(x_i)p_i. \quad (2)$$

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

$$139 \quad p_i^a = L(x_i)p_i/S. \quad (3)$$

140

141 The CDF, $F(x)$, is defined as the probability that the value of a random variable is less than or
 142 equal to x . The QCEFF computes quantiles for each ensemble member by evaluating the CDF.

143 The CDF for a discrete generalized PDF is

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

145

146 However, use of the standard F can lead to biased analysis ensembles as discussed in Appendix

147 A. To avoid this, a modified CDF, \tilde{F} , is defined as

$$148 \quad \tilde{F}(x) = \begin{cases} 0 & \text{if } x < x_1 \\ \sum_{k=1}^i p_k & \text{if } x_i < x < x_{i+1}, \quad i \in \{1, \dots, K-1\} \\ 1 & \text{if } x > x_K \\ \sum_{k=1}^{i-1} p_k + \frac{p_i}{2} & \text{if } x = x_i \end{cases}. \quad (5)$$

149 This modified CDF differs from the standard definition only at values of x with finite probability
 150 and avoids the bias issue (Appendix A).

151

152 The QCEFF algorithm also requires computing an inverse of the modified CDF. The inverse is
 153 well defined for all values in the range of \tilde{F} . This is sufficient for the first part of the QCEFF

154 algorithm that computes increments for an observed variable (A22); quantiles are strictly

155 conserved. However, the computation of the impact of an observation on another variable in

156 A23 can require the inverse of quantiles with arbitrary values between 0 and 1. A generalized
 157 inverse with domain $[0,1]$ is defined as

$$158 \quad \tilde{F}^{-1}(y) = \begin{cases} x_1 & \text{if } y \leq p_1 \\ x_i & \text{if } \sum_{k=1}^{i-1} p_k < y \leq \sum_{k=1}^i p_k, \quad i \in \{2, \dots, K\} \end{cases} \quad (6)$$

159 Note that $x = \tilde{F}^{-1}(\tilde{F}(x))$ but $\tilde{F}(\tilde{F}^{-1}(y))$ is not necessarily equal to y (see Fig. A1 for an
 160 example). With these definitions, it is possible to define a QCEFF that uses any discrete prior,
 161 like a particle filter, in observation space for the first part of the two-step filter and for the PPI
 162 in the regression step.

163

164 As noted in the introduction, mixed distributions are relevant to many geophysical problems.
 165 The discrete part of a prior mixed distribution is represented as above except that $\sum p_i = \alpha$;
 166 the continuous part of the PDF is $(1 - \alpha)f_c(x)$, with $0 < \alpha < 1$. The normalizing constant for
 167 the product with a likelihood is

$$168 \quad S = \alpha \sum_{i=1}^K L(x_i) p_i + (1 - \alpha) \int_{-\infty}^{\infty} L(x) f_c(x) dx. \quad (7)$$

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

$$170 \quad (1 - \alpha) f_c(x) L(x) / S. \quad (8)$$

171 The corresponding CDF is

$$172 \quad F_m = (1 - \alpha) \int_{-\infty}^x f_c(t) dt + \alpha F(x), \quad (9)$$

173 where F is defined in (4), and the subscript m indicates this is a mixed distribution. Again,
 174 Appendix A shows that the use of this mixed CDF leads to additional bias problems when used
 175 with the QCEFF.

176

177 A modified CDF corresponding to a mixed PDF is

$$178 \quad \tilde{F}_m = (1 - \alpha) \int_{-\infty}^x f_c(t) dt + \alpha \tilde{F}(x), \quad (10)$$

179 where \tilde{F} is defined in (5). The appendix discusses how this modified CDF eliminates issues of
 180 biased analysis ensembles. Note that the modified CDF is equivalent to the standard CDF if
 181 there are no discrete parts to the distribution. A generalized inverse of the modified CDF can be

182 defined in a similar way as for the discrete distribution and $x = \widetilde{F}_m^{-1}(\widetilde{F}_m(x))$ but $\widetilde{F}(\widetilde{F}^{-1}(y))$
183 is not necessarily equal to y (see Figure A2 for an example).

184

185 3. Tools for data assimilation with duplicate ensemble members

186

187 a. Bounded normal rank histogram distribution

188

189 The QCEFF described in A22 and A23 requires a CDF to compute observation increments and to
190 do the regression of those increments onto model state variables. The bounded normal rank
191 histogram (BNRH) distribution is an extension of the rank histogram filter distribution that was
192 developed for observation space increments (Anderson 2010). A BNRH distribution is
193 particularly useful when the appropriate distribution family is unknown.

194

195 A23 describes the PDF, $f(x)$, associated with a BNRH when there are no duplicate ensemble
196 members. An N-member ensemble partitions the real line into N+1 intervals. The interior
197 intervals are bounded on both sides; the intervals on the tails can be bounded on one side only
198 if the quantity itself is not bounded, or bounded on both sides if the quantity is bounded. The
199 BNRH PDF assigns $1/(N+1)$ probability to each interval. The probability is uniformly distributed
200 over the range of an interior interval. For intervals on the tails, the probability density is part of
201 a normal distribution. The DA algorithms in A22 and A23 require the CDF which is defined in the
202 standard fashion as $F(x) = \int_{-\infty}^x f(x)dx$. An example of a BNRH CDF is shown in Figure 1a
203 (reproduced from A23) for a 5-member ensemble.

204

205 The definition of the BNRH CDF is extended here for the case when there are ensemble
206 members with duplicate values or when one or more ensemble members have the same value
207 as the upper or lower bound of x . As noted in Appendix A, using the standard definition of the
208 CDF can lead to biased analysis ensembles. Instead, a modified version that avoids the bias can
209 be defined in the same way as for the mixed distributions in the previous section. Suppose that
210 possible values of x are bounded below by $B_l \geq -\infty$ and above by $B_u \leq \infty$. Given an N-

211 member ensemble of x with members not necessarily unique, there is at least one ordering of
 212 the ensemble values so that $x_i \leq x_{i+1}$ for $i \in \{1, \dots, N - 1\}$. Given such an ordering, define the
 213 modified CDF as:

$$214 \quad \tilde{F}(x) = \begin{cases} 0 & \text{if } x < B_l \\ C(B_l)/[2(N + 1)] & \text{if } x = B_l \\ A_l \Phi(\mu_l, \sigma^2; x) - A_l \Phi(\mu_l, \sigma^2; B_l) & \text{if } B_l < x < x_1 \\ [i + (x - x_i)/(x_{i+1} - x_i)]/(N + 1) & \text{if } x_i < x < x_{i+1}, i \in \{1, \dots, N - 1\} \\ i/(N + 1) + [C(x) - 1]/[2(N + 1)] & \text{for min } i \text{ with } x = x_i, B_l < x < B_u, i \in \{1, \dots, N\} \\ A_u \Phi(\mu_u, \sigma^2; x) - A_u \Phi(\mu_u, \sigma^2; B_u) + 1 & \text{if } x_N < x < B_u \\ 1 - C(B_u)/[2(N + 1)] & \text{if } x = B_u \\ 1 & \text{if } x > B_u \end{cases}$$

216 (11)

217 $C(x)$ is a function with unbounded real domain and range that are the whole numbers less than
 218 or equal to N , defined as the number of ensemble members with value x . $C(x)$ is only nonzero
 219 at values of x that are associated with at least one ensemble member and is used to define the
 220 location of discontinuous jumps in the BNRH CDF. $\Phi(\mu, \sigma^2; x)$ is the CDF of a normal
 221 distribution with mean μ and variance σ^2 evaluated at x , and σ^2 is the sample variance of the
 222 ensemble. The means and amplitudes of the normally distributed portions are defined as in A23
 223 so that $1/(N + 1)$ probability lies between the outermost ensemble member and the bounds.

224 The means are selected so that

$$225 \quad \Phi(\mu_l, \sigma^2; x_1) = \frac{1}{N+1} \quad , \quad (12)$$

$$226 \quad \Phi(\mu_u, \sigma^2; x_N) = \frac{N}{N+1} \quad , \quad (13)$$

227 and the amplitudes are

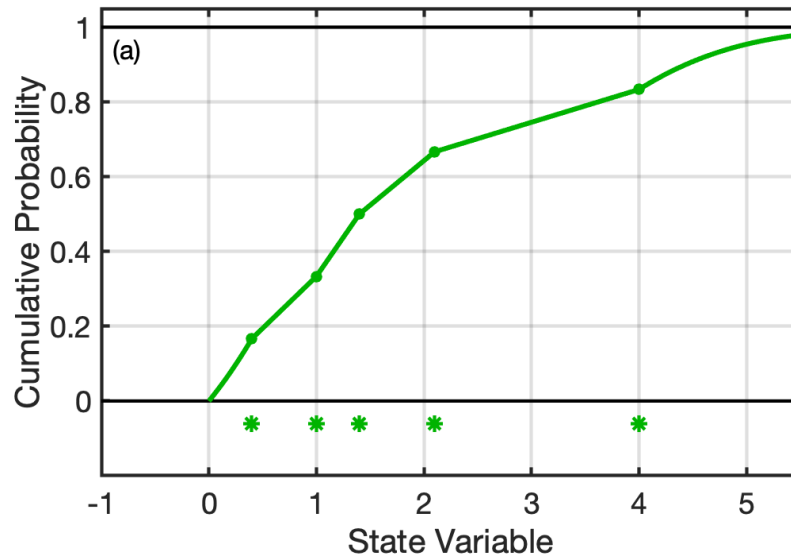
$$228 \quad A_l = \frac{1}{(N+1)[\Phi(\mu_l, \sigma^2; x_1) - \Phi(\mu_l, \sigma^2; B_l)]} \quad , \quad (14)$$

$$229 \quad A_u = \frac{1}{(N+1)[\Phi(\mu_u, \sigma^2; B_u) - \Phi(\mu_u, \sigma^2; x_N)]} \quad . \quad (15)$$

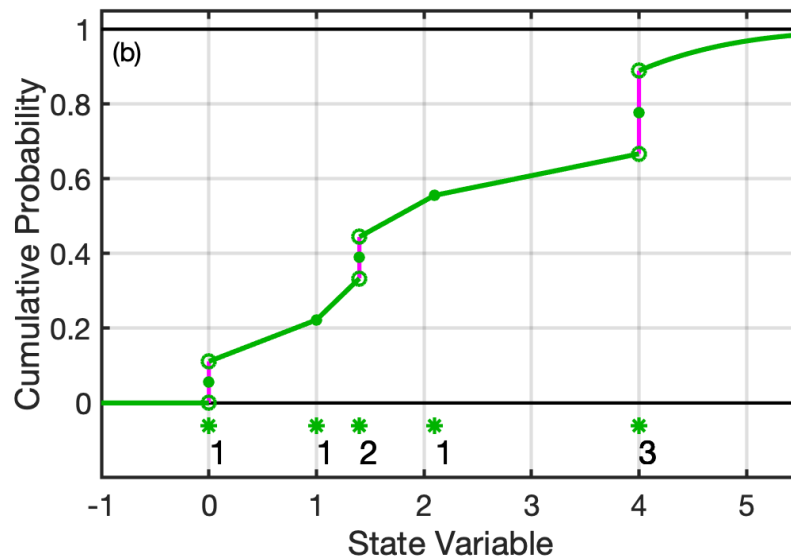
230

231 When there are no duplicate ensemble values, $C(x_i) = 1 \forall x_i$, and no ensemble values equal to
 232 the bounds, $C(B_l) = C(B_u) = 0$, the BNRH modified CDF is the same as the standard CDF and
 233 equal to the integral from $-\infty$ to x of the BNRH PDF defined in appendix C of A23. The

234 generalized inverse of the modified CDF (depicted in figure 1b) is defined in the same way as for
 235 mixed distributions in Section 2.
 236



237



238

239

240 Figure 1: Cumulative distribution functions (green) for a BNRH distribution for a 5-member
 241 ensemble (green asterisks) for a variable that is bounded below at zero (a) and modified CDF
 242 for an 8-member ensemble with duplicate values and a member with a value at the bound of
 243 zero (b). The number of duplicates is given by the integer next to asterisk. The vertical magenta
 244 lines indicate the generalized inverse cumulative distribution function (the quantile function)
 245 used for the BNRH. Panel a is reproduced from figure C1a in A23.

246

247 An example modified CDF for an ensemble with $N = 8$ members, $B_l = x_1 < x_2 < x_3 = x_4 <$
248 $x_5 < x_6 = x_7 = x_8$ and $B_u = \infty$ is shown in green in Figure 1b. The interval on the upper tail is
249 a portion of a normal CDF. $1/(N+1)$ probability is uniformly distributed in each interior interval.
250 In non-zero range interior intervals, the CDF is piecewise linear. In the case of the duplicate
251 ensemble members, the range of the interval between them can be thought of as zero and the
252 distribution is discrete. At the point x_3 where there are two ensemble members, there is
253 $1/(N+1)$ probability while at the point x_6 with three ensemble members, there is $2/(N+1)$
254 probability. Generalizing, at any point with D duplicate ensemble members, there is $(D-1)/(N+1)$
255 discrete probability. Consistent with section 2 and (11), the BNRH CDF at a point with duplicate
256 ensemble members is set to the ‘midpoint’ of the discontinuous jump in the integral of the PDF.
257 For instance, at x_3 the CDF is defined as

$$258 \quad F(x_3) = \left[\frac{3}{N+1} + \frac{4}{N+1} \right] / 2. \quad (16)$$

259 With this modified CDF, the quantile computed for ensemble members that share a value is the
260 same. The appendix discusses other possible modified CDF variants, but those are not explored
261 further here.

262

263 *b. Rank histograms*

264

265 The rank histogram is a standard tool for the verification of ensemble forecasts (Anderson
266 1996, Hamill 2001). This section describes an extension to the standard rank histogram that is
267 useful for evaluating ensemble forecasts that contain possibly duplicated ensemble values. This
268 should not be confused with the earlier discussion of the BNRH which is related to but distinct
269 from the rank histogram.

270

271 Consider a sample of $N + 1$ numbers composed of an N -member ensemble estimate (forecast)
272 of a scalar quantity and an additional value, the verification of the forecast. If there are no
273 duplicate values in the sample, the rank of the verification is uniquely defined with an integer
274 value in $\{1, 2, \dots, N + 1\}$. Define a rank weight vector, W_n , $n = 1, \dots, N + 1$ as

275
$$W_n = \begin{cases} 1 & \text{if } \text{rank}(\text{verification}) = n \\ 0 & \text{otherwise} \end{cases} . \quad (17)$$

276 For example, if the verification is the 5th smallest value in the sorted set of $N+1$ values, then
 277 $W_5 = 1$ while all other elements of W are 0. Define the vector $S_n, n = 1, \dots, N + 1$ to be the
 278 sum of the rank weight vector for a collection of M ensembles with verifications as

279
$$S_n = \sum_{m=1}^M W_n^m . \quad (18)$$

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

286

287 For state variables in many common Earth system DA applications, the probability that the
 288 verification duplicates one or more ensemble members is negligible, and most discussions of
 289 rank histograms have ignored the possibility. However, this is no longer the case for some types
 290 of bounded state variables which have mixed probability distributions like the examples
 291 discussed in Section 1. If the verification duplicates one or more ensemble members, its rank is
 292 no longer uniquely defined by (17). Suppose that D ensemble members have the same value as
 293 the verification. When these are removed from the ensemble, the rank of the verification in the
 294 $N + 1 - D$ remaining numbers is uniquely defined, even if there are other duplicate values in
 295 the remaining ensemble; let that rank be R . The actual rank in the full ensemble could range
 296 from R to $R + D$ since the order of the verification and its duplicates is not uniquely defined.
 297 Essentially, there is a $1/(D + 1)$ probability that the rank of the verification is any of these
 298 values. In this case, define the weight vector as

299
$$W_n = \begin{cases} 1/(D + 1) & \text{if } R \leq n \leq R + D \\ 0 & \text{otherwise} \end{cases} . \quad (19)$$

300 A sum of rank weight vector can be defined for a collection of ensembles as before with (18),
 301 and the histogram should be uniform if the verifications are drawn from the same distribution
 302 as the ensemble members. Another possible way to define rank histograms for duplicate values

303 is to randomly select one of the ranks between R and R+D and give it the weight of one,
304 however this random selection generates unnecessary sampling noise compared to the solution
305 in (19).

306

307 This treatment of duplicates for rank histograms is essential for application to state variables or
308 true observations in an OSSE like the one in section 4. When rank histograms are used for
309 verifications that are real observed quantities, it is important to account for observational error
310 when generating an appropriate ensemble (Anderson 1996). One way to do this is to add a
311 random sample from an observational error distribution to each ensemble member generated
312 by applying a forward operator to the model state. In many cases, adding in this observation
313 error component would eliminate duplicate values like those that result from bounded state
314 variables in state space. However, if the error distribution is also mixed, duplicates are still
315 expected. Note that a deterministic method similar to (19) can also be developed to account for
316 observational error in the rank histogram.

317

318 4. A tracer advection extension of the Lorenz-96 Model: L96-T

319

320 *a. Model description*

321

322 A low-order model with sensitive dependence on initial conditions, low computational cost, and
323 bounded state variables is useful for testing DA algorithms. The traditional Lorenz-96 model
324 (Lorenz and Emanuel 1998) has been used in many ensemble DA studies including (A22). Here,
325 the Lorenz-96 model is extended to include two additional types of M variables that are
326 collocated with the standard variables, $x_m, m = 0, \dots, M - 1$, on the standard periodic domain.
327 The first type, q_m , represents concentrations of a dimensionless tracer. The second type, s_m ,
328 represents a source rate of the tracer with units of tracer amount per time. A function of the
329 standard x variables is treated as a wind field that passively advects the tracer. The velocity at
330 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
331 constant mean velocity, \tilde{V} is a specified multiplying constant that controls the average

332 magnitude of wind perturbations, and $\tilde{V}x_m$ is an anomalous velocity at gridpoint m . Velocities
 333 are expressed with units of nondimensional distance per nondimensional time. A
 334 nondimensional location is assigned to each grid point in the model so that the distance
 335 between neighboring grid points is 1 (note that this is different from many previous Lorenz-96
 336 studies where the distance between grid points is defined as $1/M$).

337

338 The time evolution of the standard variables, x_m , is identical to that used in the basic Lorenz-96
 339 model (Lorenz and Emanuel 1996). The time evolution of the nonnegative tracer concentration
 340 used here is:

$$341 \quad q_m^+ = \max[(q_m^{adv} + s_m \Delta t)e^{-E\Delta t} - C\Delta t, 0] \quad , \quad (20)$$

342 where q_m^+ is the tracer concentration at grid point m at the next time step, q_m^{adv} is the advected
 343 concentration, s_m is the source rate at grid point m at the current time, E is an exponential
 344 damping time, C is a constant sink rate, and Δt is the timestep.

345

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

- 350 1. A preliminary upstream target location is defined as $T = m - v_m \Delta t$,
- 351 2. The fractional location of the target between the bounding grid points is $p = T - [T]$,
 352 where the brackets indicate the floor,
- 353 3. The indices of the grid points bounding the target location are computed as $L =$
 354 $\text{mod}([T], M)$ and $U = \text{mod}(L + 1, M)$,
- 355 4. The advected concentration is $q_m^{adv} = (1 - p)q_L + pq_U$.

356

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

$$360 \quad s_m = \begin{cases} 5 & \text{if } m = 1 \\ 0 & \text{otherwise} \end{cases}$$

361

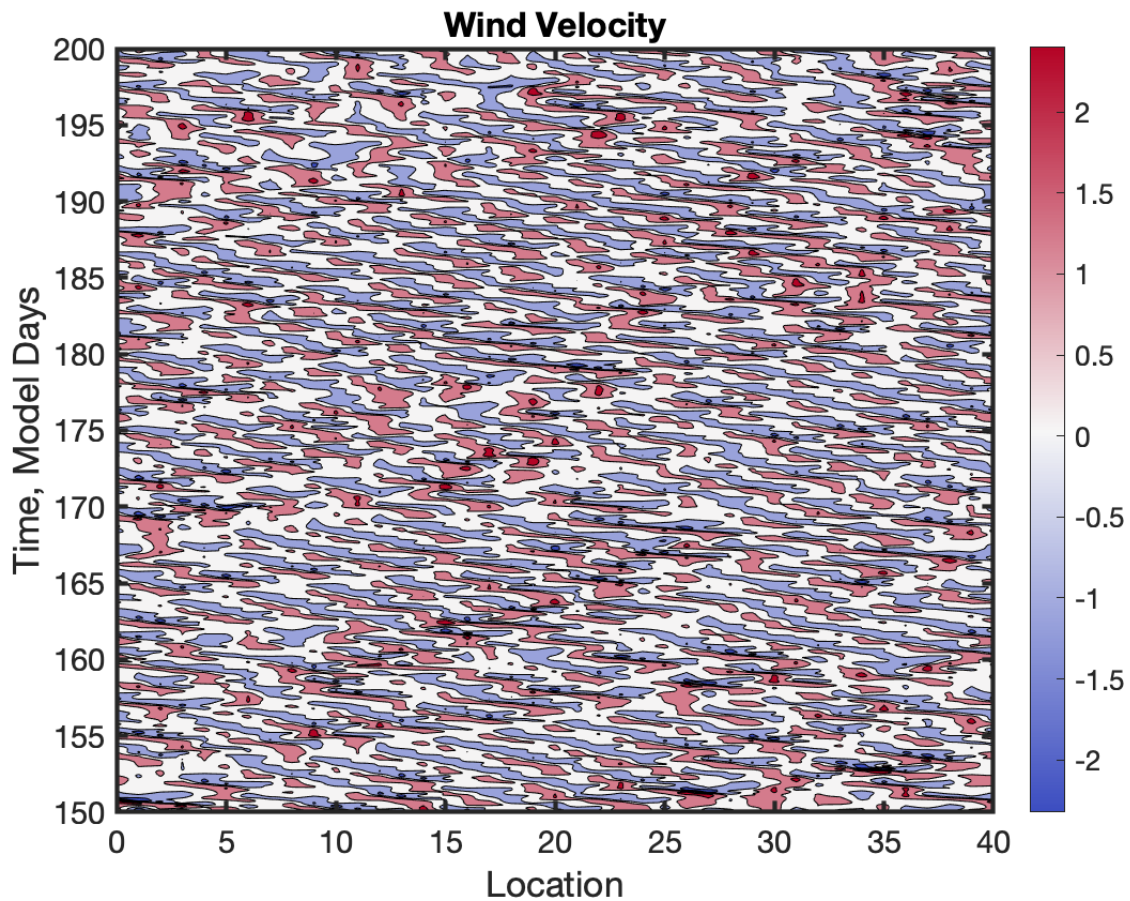
362 *b. L96-T example*

363

364 All results here use the standard 4th order Runge-Kutta time stepping algorithm, the
365 nondimensional $\Delta t = 0.05$ with an associated dimensional time step of 3600s as done in many
366 previous studies, and $M = 40$ grid points. The L96 forcing parameter $F = 8$. The mean velocity
367 $\bar{V} = 0$ and the velocity perturbation multiplier $\tilde{V} = 5$, while the constant sink $C = 0.1$, and the
368 exponential sink $E = 0.25$.

369

370 Figure 2 shows a time series of the wind field, v_m , as a function of the model grid point; since
371 $\bar{V} = 0$, this is just $\tilde{V} = 5$ times the standard L96 state variables, x_m . The well-known group and
372 phase velocity of the L96 model can be seen.

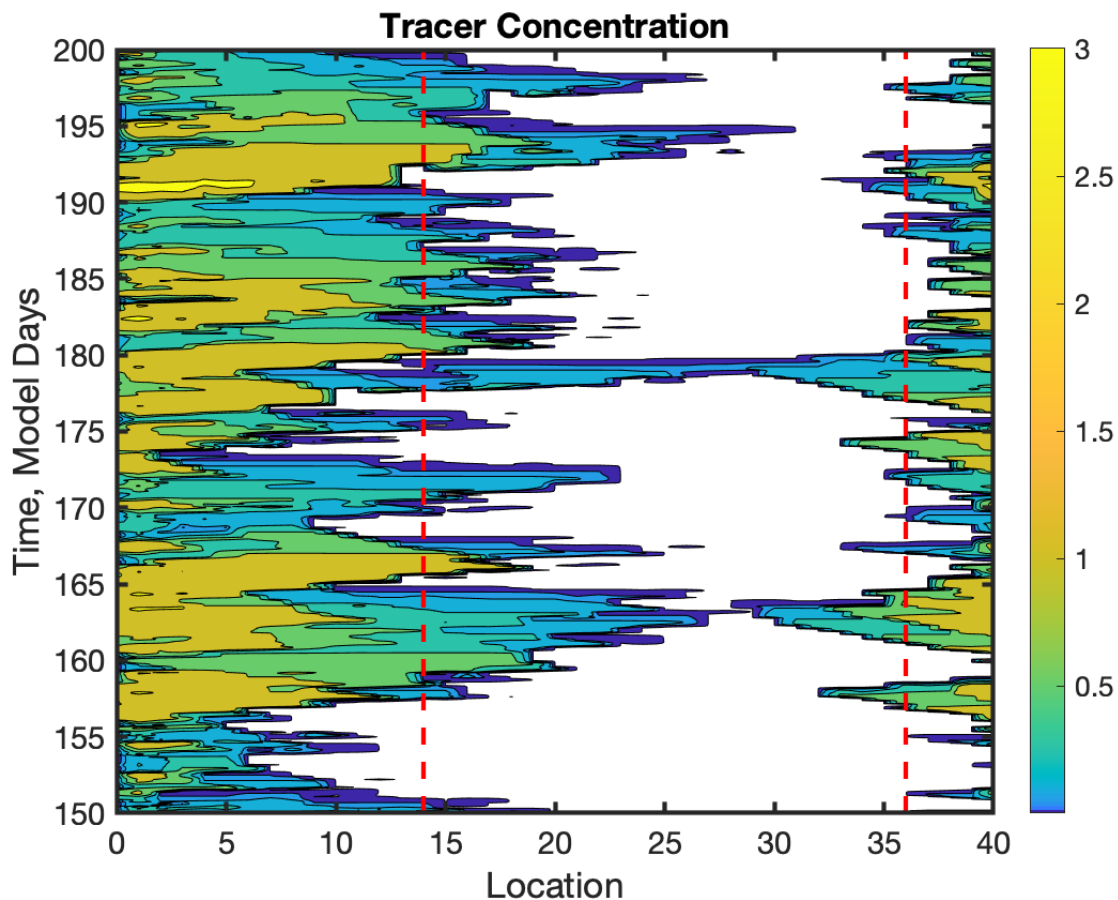


373

374 Figure 2: Wind velocity from the L96 Tracer Advection model for times between day 150 and
375 200 in the truth run as a function of model grid point. Units are distance hr^{-1} .
376

377 Figure 3 shows the tracer concentration corresponding to the wind field. Plumes of tracer are
378 advected away from the source at grid point 1. The velocity is positive more often than
379 negative, so plumes extend more frequently and further to the right. However, the wind field is
380 sometimes negative leading to shorter plumes extending to the left. The white areas in the plot
381 have zero tracer concentration, so a mixed distribution is appropriate. For example, the
382 climatological frequency of zero concentration is greater than 90% for grid points 25 to 34, but
383 less than 20% for grid points 2 to 10. It is rare for plumes to extend clear across the domain
384 with this only happening twice in the figure. This behavior is roughly analogous to what one
385 might see with a point source in the midlatitudes. It is possible to get a variety of other
386 behaviors for the tracer by changing the model parameter values.

387



388

389 Fig. 3: As in figure 2 but for the tracer concentration (nondimensional). The red dashed lines
390 highlight grid points with additional diagnostics presented in figures 4, 5 and 6. White areas
391 have zero concentration.
392

393 5. Data assimilation experiments

394

395 The model integration described in the previous section is used as the truth run for a series of
396 observing simulation system experiments (OSSEs). The L96-T model is first integrated for 16500
397 hrs (5500 3-hour advances) starting from a default initial state to generate a tuning initial state.
398 The default initial state has $x_1 = 1$ and all other x state variables are 0; all concentration
399 variables are 0. The model is integrated for an additional 16500 hrs from the tuning initial state
400 generating synthetic observations every 3 hrs. Forty randomly located observing sites are
401 selected for the L96 standard state, and a different set of 40 randomly located sites for tracer
402 concentration observations (see Figs 7d and 8d). Observations are taken by linearly
403 interpolating to the site location from the two nearest grid points. For the standard state
404 observations error is simulated by adding a random draw from a normal distribution with mean
405 0 and variance 10; the variance of the state is approximately 13. The relatively large
406 observation error variance was selected so that the state observations on their own do not too
407 strongly constrain the state analysis error. This allows the addition of observations of the tracer
408 concentration to make a noticeable reduction to the state analysis error. For tracer
409 observations error is simulated by adding a random draw from a truncated normal distribution
410 with variance 0.1 and lower bound of 0 (A23, appendix D). This observational error variance is
411 small compared to the variance of the tracer concentration near the source, but only about a
412 factor of two smaller than the variance far from the source. The analysis error variance of the
413 tracer concentration when only state observations are assimilated varies between 0.25 and 1.0
414 depending on what assimilation algorithms are used. The observation error variance of 0.1 is
415 small enough so that there is room for significant improvement when observations of tracer
416 concentration are added.

417

418 Three different observing networks are explored: assimilating only standard state observations,
419 assimilating only tracer concentration observations, and assimilating both standard and tracer
420 observations. Two different model configurations are evaluated. In the first, every ensemble
421 member has the true value of the tracer source variables. In the second, the tracer source
422 variables are unknown, and every ensemble member has its own (not necessarily unique) time
423 evolving estimate.

424

425 All assimilation experiments use the adaptive inflation algorithm of Gharamti (2018) with an
426 inflation damping of 0.9. All experiments also use a Gaspari Cohn localization with the same
427 constant halfwidth for all observations. Seven halfwidth tuning assimilation experiments are
428 done for each case, where a case is defined by the observing network, whether the source is
429 known or unknown, and the ensemble size (20, 40, 80 or 160). As in A23, the halfwidths tested
430 are $\{0.075, 0.1, 0.125, 0.175, 0.2, 0.4, \infty\}$. These tuning assimilations start from the tuning initial
431 condition and assimilate for 5500 3-hour intervals. Initial ensembles for the standard state
432 variable are generated by adding a random draw from a normal distribution with mean 0 and
433 standard deviation 0.01 to the truth value for each variable. Initial ensemble members for the
434 tracer variables are all equal to the truth. For the case with known sources, all ensemble
435 members for the source variables are equal to the truth. For the case with unknown sources,
436 ensemble members for the source are set to a random draw from a normal distribution with
437 mean 2.5 and standard deviation 2.5; if the draw is less than 0 the source is set to 0 so that the
438 resulting ensembles are generally mixed distributions with several members that are 0. Results
439 from the first 500 assimilation steps are discarded and the prior ensemble mean, time mean
440 RMSE is computed for the standard state and tracer variables for the remaining 5000 steps. For
441 the state only observing network, the localization that minimizes the state RMSE is selected; for
442 the other observing networks, the localization that minimizes the tracer RMSE is selected.

443

444 The model truth is then integrated for an additional 16500 hours from the end of the tuning
445 integration with synthetic observations generated in the same way. Initial conditions for
446 ensembles are also generated in the same way as for the tuning experiments. Assimilation

447 experiments are performed for each case using the tuned localization and assimilating every 3
448 hours. The first 500 steps are discarded, and results are available for the final 5000 assimilation
449 steps. The spread for all quantities appears to be spun up after fewer than 100 assimilation
450 steps for all experiments.

451

452 Four different assimilation algorithms are applied to each case using the QCEFF. As noted in
453 A23, a complete description of a QCEFF assimilation algorithm requires information about the
454 first step where increments are computed for observed variables and the second step where
455 those increments are regressed onto state variables. The QCEFF uses a probit probability
456 integral transform (PPI) before doing the regression (A23). Table 1 specifies the details of the
457 four algorithms which are referred to as an EAKF, RHF, PQBNRH, and DUAL.

458

459 The EAKF is the standard Ensemble Adjustment Kalman Filter (Anderson 2001, 2003) which
460 assumes normal continuous distributions for all state and observation variables and for the
461 observation likelihood. Note that the likelihood used for the q variable in the EAKF is a normal
462 distribution with the same variance as the truncated normal observation error distribution for
463 q . As noted in A23, using a normal distribution for the PPI transform is equivalent to no
464 transform.

465

466 The RHF is an extension of the original Rank Histogram Filter (Anderson 2010). It uses the rank
467 histogram prior distribution for the scalar update of the observation prior for the Lorenz-96
468 state variable. It uses a BNRH prior for the scalar update of the observation prior for the tracer
469 concentration and source variables. The BNRH CDFs all have a lower bound of 0 and no upper
470 bound, consistent with the nature of the tracer concentration and source variables. Like the
471 EAKF, it uses standard linear regression to compute the increments for state variables from
472 observation increments.

473

474 For consistency with Anderson (2022), PQBNRH (Probit Quantile BNRH) is used to describe the
475 third algorithm. It uses the same rank histogram and BNRH distributions for scalar updates of

476 the observation prior as the RHF. However, it performs a PPI transform with a BNRH
 477 distribution for all variables before computing the increments for state variables via regression.
 478

479 The final algorithm is referred to as DUAL (not an acronym) because it uses a traditional EAKF
 480 for standard L96 state for both observation increments and regression, but a BNRH for both
 481 observation increments and regression for the tracer and source.

482

	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

483

484 Table 1: Assimilation settings for each of the four algorithms explored. For the x and q
 485 variables, the continuous CDF and form of the likelihood used for computing observation space
 486 increments are listed with *normal* referring to a normal distribution, *RH* referring to a rank
 487 histogram distribution without bounds and *BNRH* referring to a bounded normal rank
 488 histogram distribution with a lower bound at zero. The continuous distribution used as part of
 489 the PPI transform used when regressing observation increments onto state variable increments
 490 is also listed.

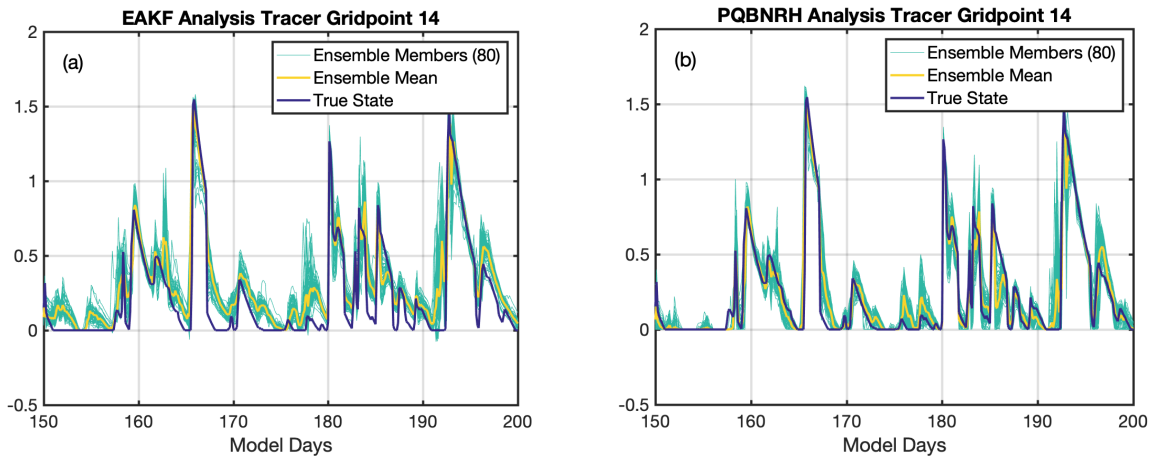
491

492 *a. Known source results*

493

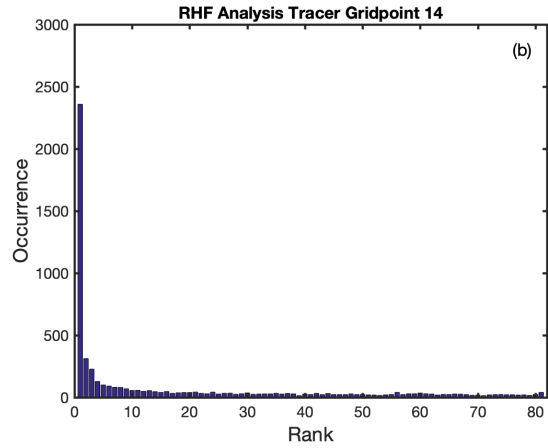
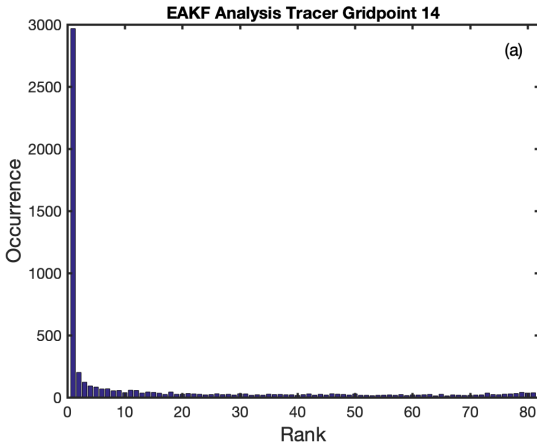
494 Unless otherwise noted, all results shown are for analysis, rather than forecast, variables. Also,
 495 results shown are for the network observing both standard state and tracer observations unless
 496 otherwise noted. Figure 4 shows a time series from the EAKF and PQBNRH algorithm 80-
 497 member assimilations for a tracer at grid point 14, which is highlighted by a red dashed line in
 498 Fig. 3. The EAKF ensemble in Fig. 4a represents all the plumes that occur, but also represents

499 two plumes between days 150 and 160 that are not real. The ensemble is strongly biased
500 towards larger values at some times, in particular around days 168, 173, and 178. The PQBNRH
501 results in Fig. 4b also capture all real plumes with smaller values for the two false plumes.
502

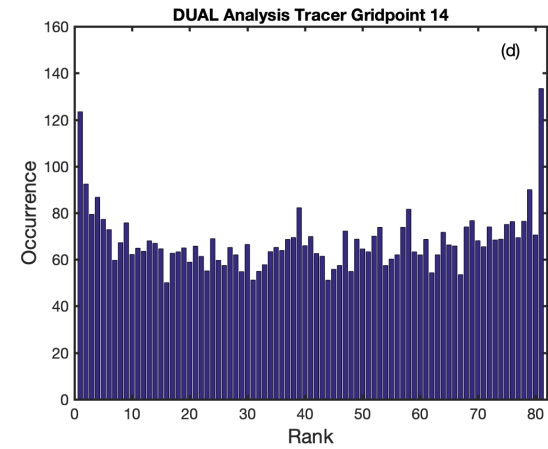
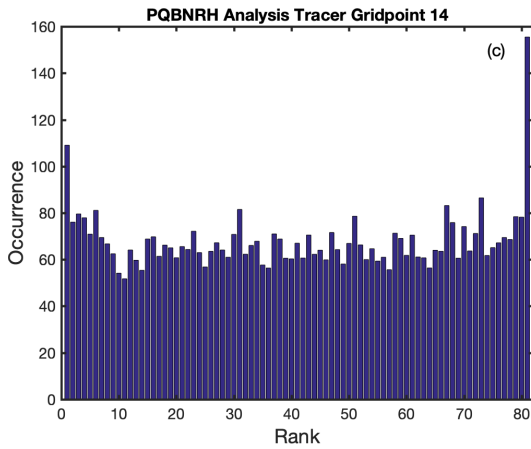


503
504 Fig. 4: Time series of the tracer at grid point 14. Dark curve is the truth and is the same in both
505 panels. The green curves are the 80 analysis ensemble members, and their mean is in yellow,
506 for an EAKF (a) and a PQBNRH (b); the tracer is nondimensional.

507
508 Figure 5 shows rank histograms over all 5000 assimilation steps for concentration at grid point
509 14. The EAKF and RHF algorithms result in very strongly biased histograms with the truth very
510 often less than the smallest ensemble member. The results for the PQBNRH and DUAL are
511 radically different. Both have histograms that are nearly uniform except for the two outermost
512 bins. The PQBNRH has more cases where the truth is larger than the largest ensemble member
513 while the DUAL algorithm has more cases where the truth is smaller than the smallest member;
514 however, it is difficult to evaluate whether these differences are statistically significant.
515



516
517



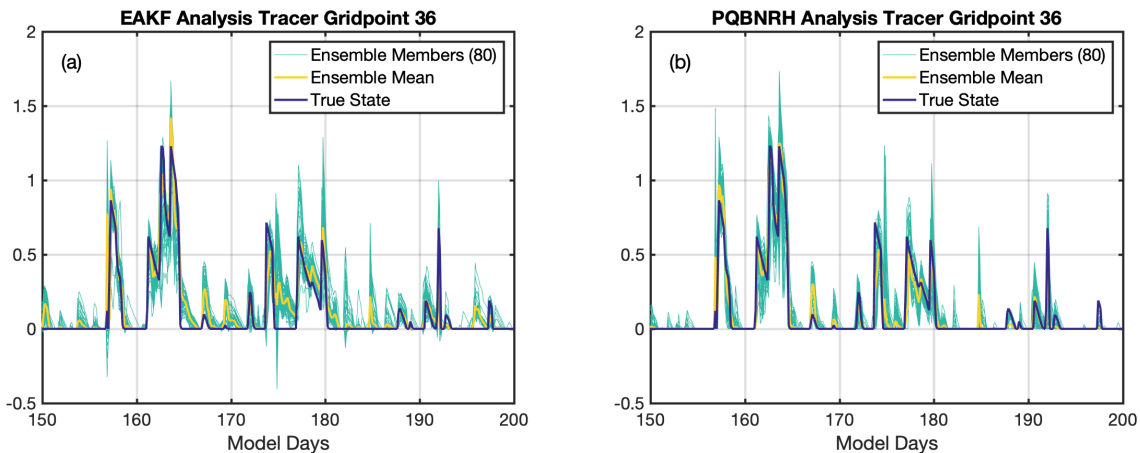
518
519

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

523

524 Fig. 6 shows time series of the EAKF and PQBNRH assimilation results for grid point 36 which is
525 also highlighted in Fig. 3. At this grid point, plumes are less frequent, primarily arriving from the
526 right. There are extended periods when the true concentration is 0. The EAKF represents all
527 true plumes, however, there are several instances where the ensemble is strongly biased
528 towards larger concentration than the truth, and several times when negative ensemble
529 members occur; this cannot happen with the PQBNRH. The EAKF never has any ensemble
530 members that are exactly zero and never has duplicate ensemble members. The PQBNRH also
531 captures all real plumes and has fewer instances of false plumes. The PQBNRH has several

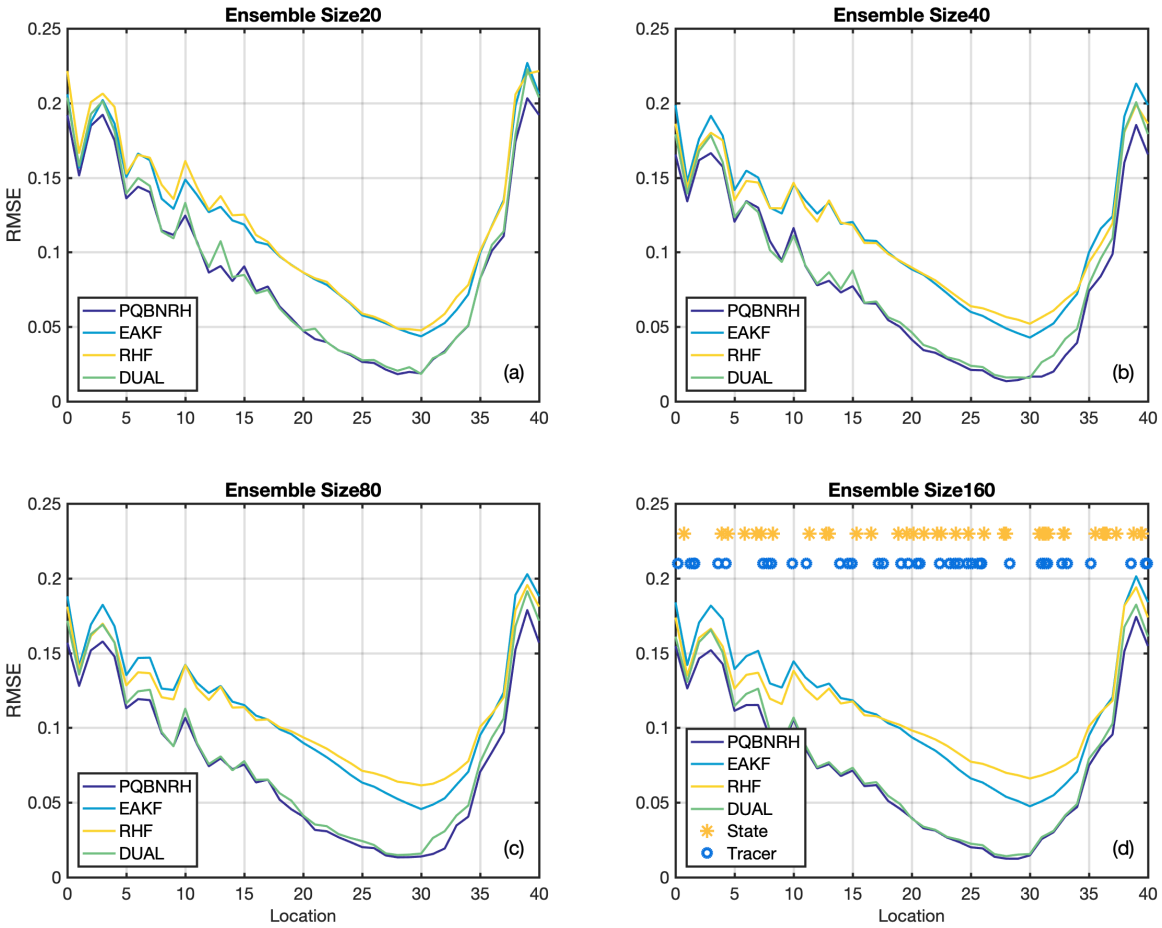
532 periods when many ensemble members are exactly 0 and some periods where all members are
533 zero, all at times when the truth is also zero. Results for the RHF are similar to those for the
534 EAKF, and results for DUAL are similar to those for the PQBNRH in figures 4 and 6 so these are
535 not displayed.
536



537
538 Fig. 6: As in figure 4 but for grid point 36.

539
540 Fig. 7 shows summary results for ensemble mean tracer RMSE over all 5000 assimilation steps
541 for the four algorithms and four ensemble sizes studied. In general, the results for the PQBNRH
542 and DUAL algorithms are statistically indistinguishable. The same is true for the EAKF and RHF
543 algorithms. However, in general the PQBNRH/DUAL algorithms have lower RMSE. The RMSE is
544 largest to the left of the source at grid point 0 where the true concentration is most variable,
545 and smaller far from the source where concentration is smaller. There are not large differences
546 as a function of ensemble size; larger ensembles generally have only slightly smaller RMSE. It is
547 unclear why ensemble size is not more important here and this certainly merits future study.

548

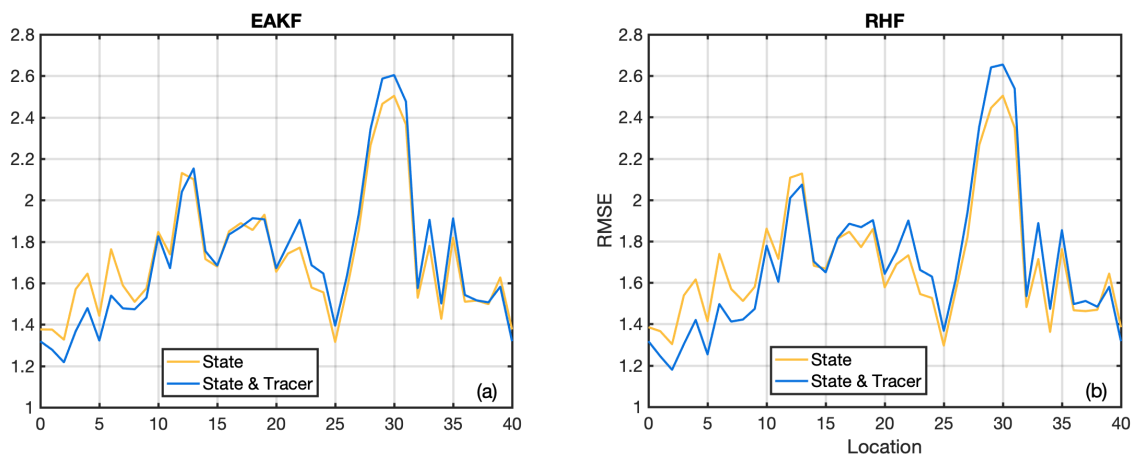


549

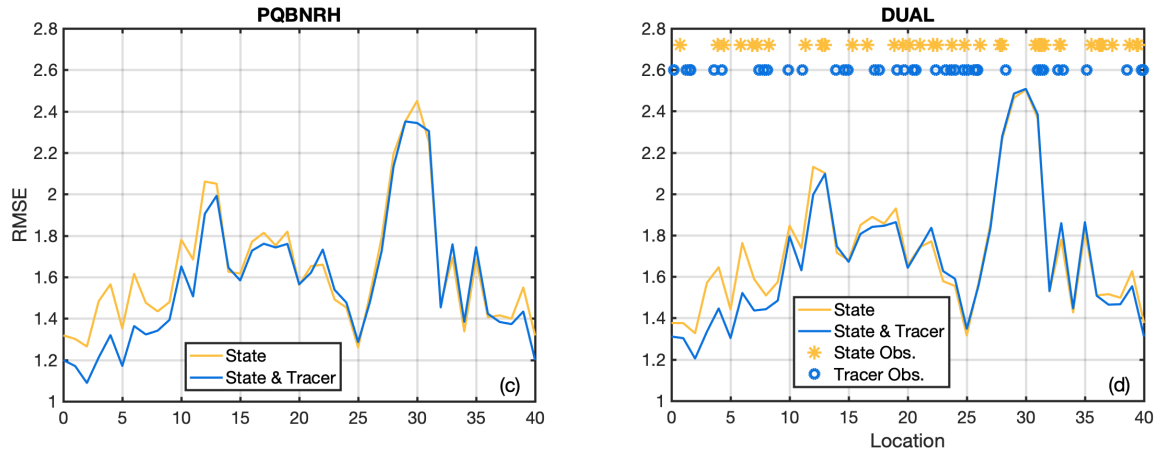
550 Fig. 7: Ensemble mean, time mean RMSE as a function of grid point for the analysis tracer
551 concentration for four filter algorithms for ensemble size 20 (a), 40 (b), 80 (c) and 160 (d). The
552 locations of the 40 observing stations are shown in (d) for state (yellow asterisks) and tracer
553 concentration (blue circles).
554

555 It is obvious that assimilating standard state observations that improve the estimate of the
556 winds will result in improved estimates of the tracer concentrations. However, the impact of
557 tracer observations on the standard state variables is less clear. Assimilations for the network
558 observing only tracer produced tracer analysis estimates that have much larger RMSE than
559 those just discussed, although smaller than the RMSE from an unconstrained control ensemble
560 run. The tracer only network resulted in standard variable RMSE that was only slightly smaller
561 than the RMSE from an unconstrained control.
562

563 A comparison of the standard variable RMSE from the observing network with only standard
 564 state observations to the network with both standard and tracer observations is shown in Fig. 8
 565 for the four algorithms. The RMSE for the standard observation only network has larger RMSE
 566 near grid point 30 and smaller RMSE near grid points 25 and 1. This is due to the random
 567 observing site locations (Fig. 8d). The RMSE is smaller for the PQBNRH than for any of the other
 568 algorithms; note that the EAKF and DUAL are identical for the standard observation network.
 569
 570 When tracer observations are added in, all four algorithms produce reduced RMSE for the left
 571 part of the domain. The EAKF and RHF produce increased RMSE in the right part of the domain.
 572 The PQBNRH and DUAL produce roughly equivalent RMSE in the right part of the domain. In the
 573 left part of the domain, plumes with large spatial and temporal gradients occur near the source.
 574 These provide information about the flow field that is advecting the plume and lead to the
 575 reduced RMSE for the standard state. Because there is often very little or no tracer in the right
 576 part of the domain, observations of the tracer are expected to provide very little additional
 577 information. The increase in error in the EAKF and RHF suggests that deficiencies in these
 578 algorithms cause the use of these low information observations to degrade the ensemble
 579 estimate.



580



581

582 Fig 8: Ensemble mean, time mean RMSE as a function of grid point for the standard L96 state
 583 for experiments that assimilate only observations of the standard state, and experiments that
 584 also assimilate the tracer concentration, shown for an EAKF (a), RHF (b), PQBNRH (c), and a
 585 DUAL filter with an EAKF for the wind and a BNRH for the concentration (d). The locations of
 586 the 40 observing stations are shown in (d) for state (yellow asterisks) and tracer concentration
 587 (blue circles).
 588

588

589 *b. Unknown source results*

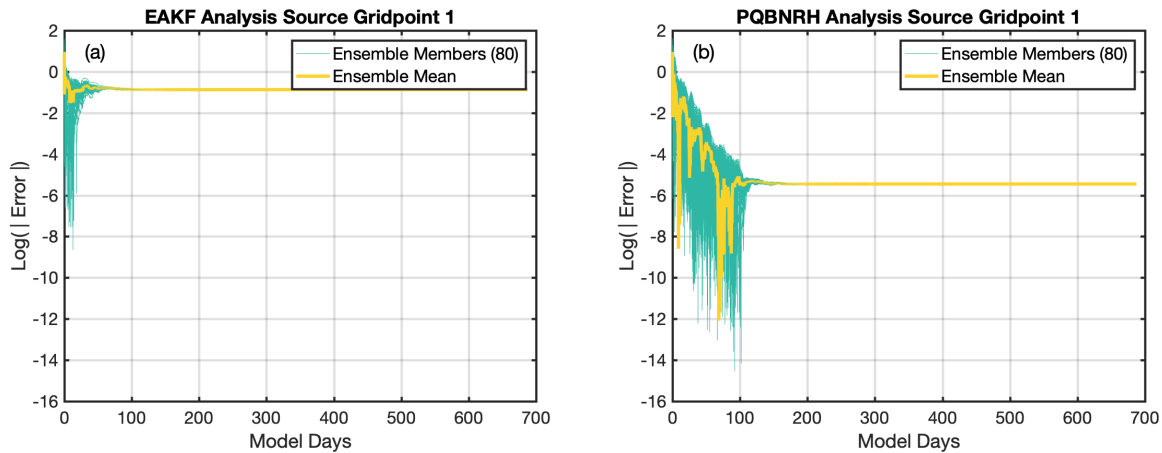
590

591 The true source is constant in time at each gridpoint, with a value of zero everywhere except at
 592 gridpoint 1. In these experiments, the source is estimated by the assimilation algorithms
 593 instead of being specified. Results are only discussed for the network observing both standard
 594 state and tracer observations. There is no time tendency model for the tracer. The prior
 595 ensembles can have their spread increased by the adaptive inflation. Nevertheless, in all
 596 experiments, the spread becomes increasingly small for the source at all grid points. The source
 597 variables are only impacted by concentration observations since the source and the state field
 598 should not be meaningfully correlated; in the truth, the state has no impact on the source and
 599 the source has no impact on the state.

600

601 Figure 9 shows the natural logarithm of the absolute value of the error for each ensemble
 602 member and the ensemble mean error at the grid point with the nonzero source in the truth
 603 for the EAKF and the PQBNRH. Both reduce the ensemble mean error, but the reduction is
 604 much larger for the PQBNRH. Because of the collapse of spread, both algorithms eventually

605 have strongly biased estimates and would produce corresponding rank histograms. As in many
606 source estimation applications, the spread of the source collapses despite the use of adaptive
607 inflation because of the limited correlation between tracer concentration observations and
608 estimates of the source.



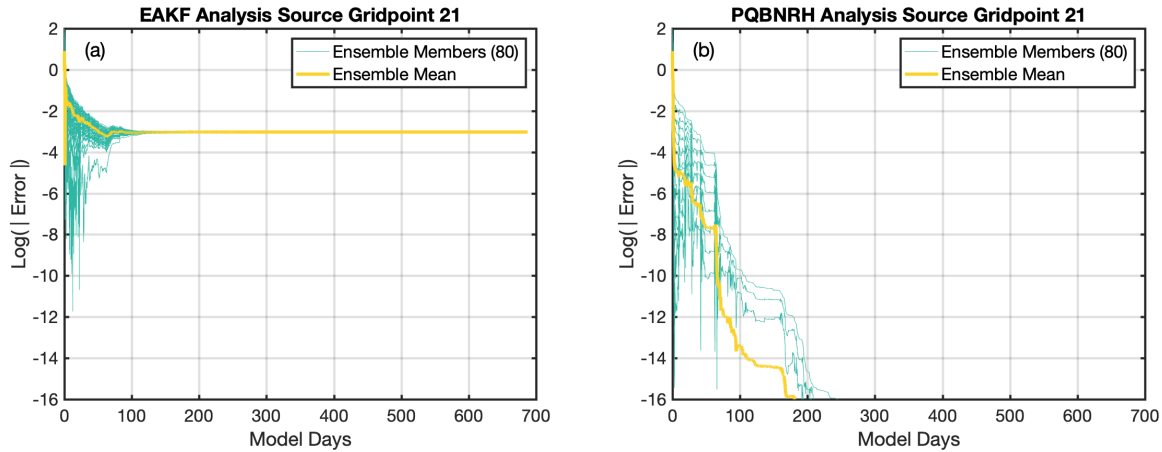
609
610

611 Fig. 9: Spatial mean of the natural logarithm of the absolute value of the error of the ensemble
612 mean (yellow) and each of the 80 ensemble members (green) as a function of time for the
613 source at grid point 1 which has a true value of 5 (units hr^{-1}) for the EAKF (a) and the PQBNRH
614 (b).

615

616 Figure 10 shows the evolution of the RMSE for grid point 21 which has zero true source. The
617 RMSE for the EAKF is smaller than it was for grid point 1. The error for the PQBNRH decreases
618 throughout the 5000 assimilation steps. As the assimilation continues, more and more
619 ensemble members have values of exactly zero; eventually all ensemble members are zero and
620 the error of the ensemble mean, and all individual ensembles is zero. At both grid points 1 and
621 21, the RMSE for the standard state and concentration variables for the PQBNRH are nearly
622 identical to those for the known source experiments since the source is so accurately
623 estimated. Results are somewhat degraded for the EAKF which has larger errors in its source
624 estimates.

625



626

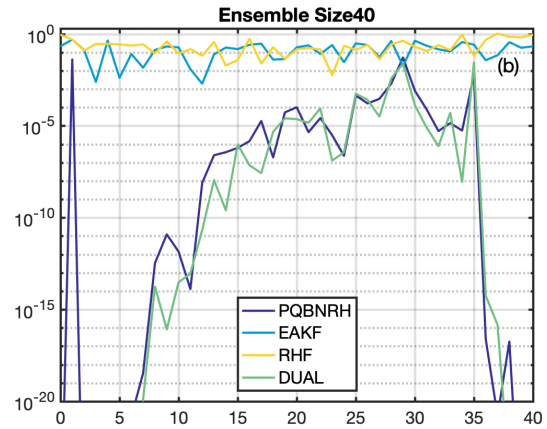
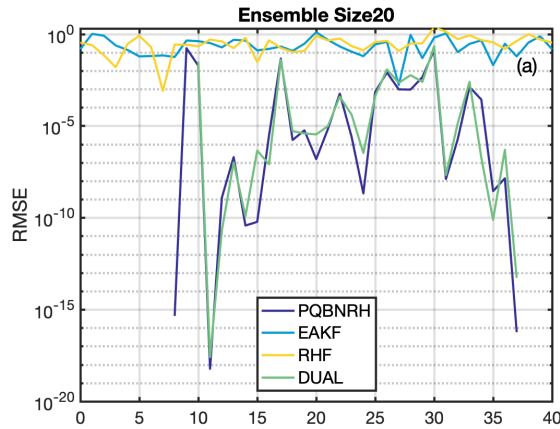
627 Fig. 10: As in 9 for grid point 21 which has zero true source.

628

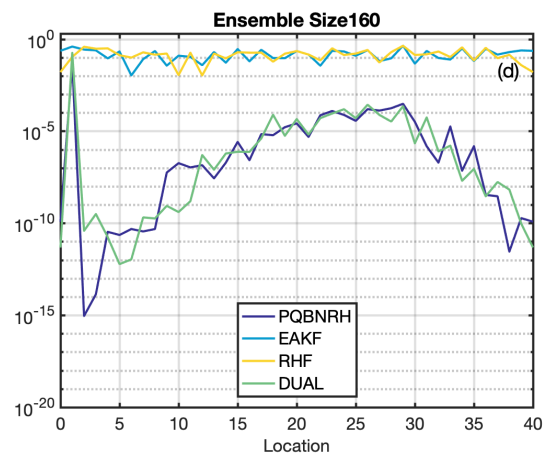
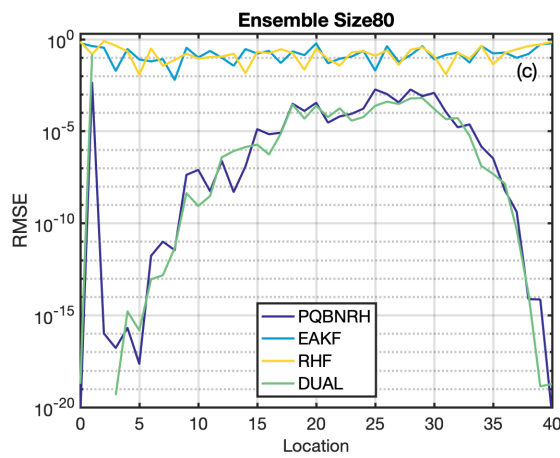
629 Figure 11 shows the RMSE for the source as a function of grid point for each of the four
 630 algorithms and all four ensemble sizes. The EAKF and RHF produce roughly comparable results
 631 that have a small dependence on ensemble size. The errors do not have a strong dependence
 632 on the grid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

640

641



642



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

646

647

648 6. Discussion and conclusions

649

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

656

657 The BNRH distribution has also been extended to handle duplicate ensemble members that are
658 expected to occur for variables with mixed distributions. The rank histogram diagnostic tool
659 was also extended to deal with duplicate ensemble members. An extension of the Lorenz-96
660 low-order dynamical system that includes an idealized advective tracer with local sources was
661 developed to test the new algorithms. This L96-T model should also provide challenging tests
662 for other data assimilation algorithms including variational methods and heuristically simplified
663 particle filters.

664
665 Results show that the BNRH works better than the EAKF or RHF for an OSSE with the L96-T
666 model. The RMSE is smaller for the bounded tracer concentration and source variables when
667 they are close to the bounds as might be expected. Results are also better when these variables
668 are not close to the bounds and for the unbounded standard state variables from L96. The RH
669 and PQBNRH algorithms use the BNRH distribution to compute observation increments.
670 However, the RH uses standard linear regression when updating state variables while the
671 PQBNRH includes the PPI transform using the BNRH distribution for the probability integral
672 transform. The RH results are similar to the EAKF results in this case, while the PQBNRH is
673 better for all variables and locations demonstrating that the PPI is a crucial part of the improved
674 performance. The DUAL case uses an EAKF for the L96 state which has no bounds and is
675 expected to be approximately normally distributed. There is some indication that the PQBNRH
676 is slightly better than the DUAL algorithm, but differences are not quantitatively significant. This
677 suggests a strategy of using the BNRH distribution for bounded variables but a normal
678 distribution for other variables may be useful for large model applications.

679
680 The BNRH as described allows duplicate ensemble members and the data assimilation process
681 can create additional duplicates; this happened for both concentration and source variable
682 ensembles in the OSSEs here. However, the assimilation process cannot eliminate duplicates. It
683 can change the value of ensemble members that are exactly at a bound in the prior ensemble.
684 This means that the model must eliminate duplicates if appropriate. That happens in
685 experiments here and is most clearly seen in figure 6b where all ensemble members are zero at

686 some times but not at subsequent times. Further investigation into methods that would allow
687 the assimilation to eliminate duplicates is warranted (see Appendix A).

688

689 All the OSSEs here were performed using the Data Assimilation Research Testbed (DART:
690 Anderson et al. 2009) which implements the QCEFF including the BNRH; the parallel algorithm
691 of Anderson and Collins (2007) was used. The results here only examined the use of normal or
692 BNRH distributions. DART software can support arbitrary distributions and currently supports
693 gamma, inverse gamma, log-normal, beta, and particle filter distributions. Previous work on
694 assimilation of bounded quantities has proposed using distributions like the log-normal,
695 gamma, and inverse gamma. However, the L96-T OSSE explored here presents specific
696 challenges for using these other distributions. The log-normal and inverse gamma distributions
697 do not have any probability at zero. This is clearly inappropriate for the mixed distributions in
698 the OSSE where much of the probability can be at 0 at some times. The gamma distribution can
699 have probability at zero. However, if the likelihood is a gamma distribution, the corresponding
700 observation error distribution is inverse gamma (Bishop 2016, A22). This means that
701 observations of the bounded quantities would not be able to have any probability at zero. This
702 is clearly problematic for the bounded quantities with realistic instruments. Further work on
703 explicitly using mixed distributions, for instance a combination of a log-normal distribution with
704 a discrete distribution, for applications like this is beyond the scope of this report.

705

706 The computational cost of the QCEFF algorithms including the BNRH is discussed in detail in
707 A23. There is almost no additional cost associated with allowing duplicate ensemble members.
708 A single additional if statement is required to test for a duplicate ensemble member in the
709 computation of the CDF. Similarly, a single if statement is required when computing the inverse
710 CDF to determine if the quantile is between the bounding quantiles of a discontinuous jump.
711 As noted in A23, the additional cost of a BNRH compared to an EAKF can be significant,
712 especially for low-order model applications. As discussed in Anderson (2019) and A23, much of
713 this cost is associated with the need to sort the ensemble members for each state variable.
714 However, the sorting order often changes little between assimilation steps. Caching the sort

715 order and then using sorts that are efficient for nearly sorted data can potentially result in large
716 computational cost reductions, but these methods have not yet been implemented in DART.

717

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

724

725 *Acknowledgements.* This material is based upon work supported by the NSF National Center for
726 Atmospheric Research, which is a major facility sponsored by the National Science Foundation
727 under Cooperative Agreement 1852977. Any opinions, findings, and conclusions or
728 recommendations expressed in this publication are those of the author and do not necessarily
729 reflect the views of the National Science Foundation. Thanks to three anonymous reviewers
730 who played a key role in improving the presentation. Thanks also to Ian Grooms, Moha
731 Gharamti, Joseph Chan, Hristo Chipilski, Ben Gaubert and the whole DAREs team for helpful
732 discussions about this material.

733

734 Data availability statement

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

736 https://github.com/NCAR/DART/releases/tag/MWR_QCEFF_Part3. All figures except figure 1
737 and the Appendix figures were produced by standard DART diagnostic tools that are
738 documented in this github repository.

739

740 APPENDIX A: Modified CDFs for discrete and mixed distributions

741

742 The CDF of a discrete distribution is defined by eqn. (4) in section 2. The QCEFF computes an
743 ensemble of quantiles by applying the CDF to each ensemble member. For the discrete CDF, the
744 average value of the quantiles will always be greater than 0.5 due to the less than or equals sign

745 in the second row of (4). This can lead to a bias in the analysis ensembles that result from
746 applying the QCEFF to discrete distributions.

747

748

749 Figure A1 shows the result of assimilating an observation on a 9-member prior ensemble from a
750 discrete distribution that is symmetrically distributed around 0. The prior CDF is shown in Fig.
751 A1a along with likelihoods for two observations with standard deviation of 1 and means of -1
752 and +1 respectively. The analysis ensemble that results from applying a QCEFF has some
753 duplicate ensemble members as indicated by the integers below the prior ensemble members.
754 Given the symmetric prior and the two likelihoods which are symmetric around 0, the analyses
755 for the two observations should be symmetric around 0. This is not the case; for instance, there
756 is 1 ensemble member with posterior value -1 for the -1 likelihood, but 2 members with value
757 +1 for the +1 likelihood. Table A1 compares the mean, standard deviation and skewness of the
758 two analysis ensembles. The mean of the +1 observation is considerably larger in magnitude
759 indicating that there is a bias toward more positive values in the algorithm. The +1 observation
760 leads to less spread; the values should be equal for an unbiased assimilation algorithm. Finally,
761 the magnitude of the skewness is less for the +1 observation.

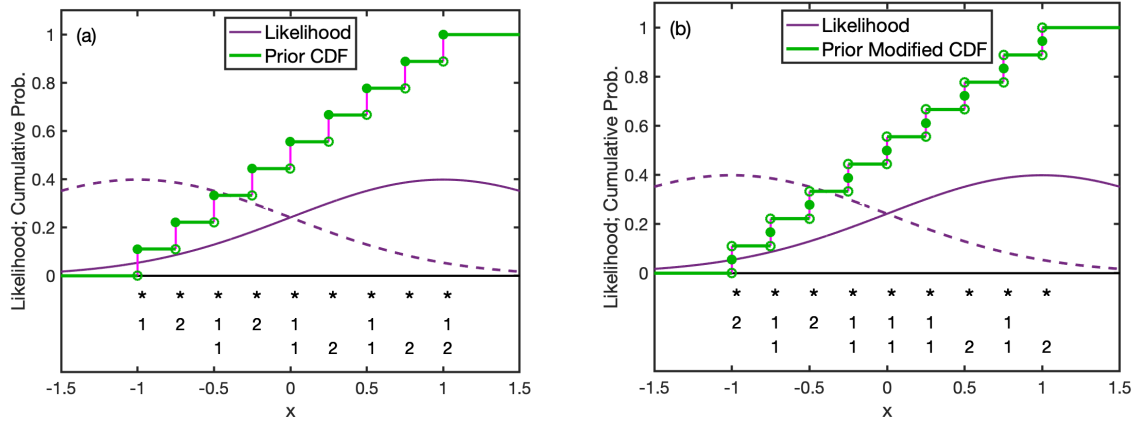
762

763 Suppose x_N^p is the largest prior ensemble value. The root of this problem is that it is almost
764 impossible to generate analysis ensembles that do not have at least one ensemble member
765 with the value x_N^p . The value of the CDF for x_N^p , $F(x_N^p)$, is 1 by definition. If the analysis
766 probability of x_N^p is even a tiny bit larger than 0, the CDF of the analysis will be less than one for
767 all values less than x_N^p . Hence, when inverting the CDF, the largest analysis member (with its
768 quantile of 1) will have value x_N^p . This problematic behavior is clearly seen in the analysis
769 ensemble for the -1 observation in Fig. A1a.

770

771 Figure A1b shows the results of the same prior ensemble and likelihoods but using the modified
772 CDF in eq. (5) to compute modified quantiles and its inverse (6) to invert them with the
773 modified CDF of the analysis distribution. Selecting the midpoint makes the mean value of $\tilde{F}(x)$

774 0.5 for any discrete distribution so that the modified quantiles are unbiased. In this case, the
 775 analysis ensemble members are symmetric around zero as also shown in Table 1. This
 776 assimilation algorithm is unbiased and does not have the problem related to the largest prior
 777 ensemble member.



778
 779 Figure A1: Results of idealized assimilation for a 9-member ensemble prior sampled from a
 780 discrete distribution. The prior ensemble is the same in both figures (asterisks). Two
 781 observation likelihoods are shown in purple with observed value -1 (dashed) and +1 (solid). The
 782 standard prior CDF is shown in green in (a) and the modified CDF in (b). The numbers below the
 783 prior ensemble indicate the number of analysis ensemble members with that value for the -1
 784 observation (higher row) and +1 observation (lower row). The magenta vertical segments show
 785 the value of the generalized inverse of the CDF for values in the discrete jumps.

786

787

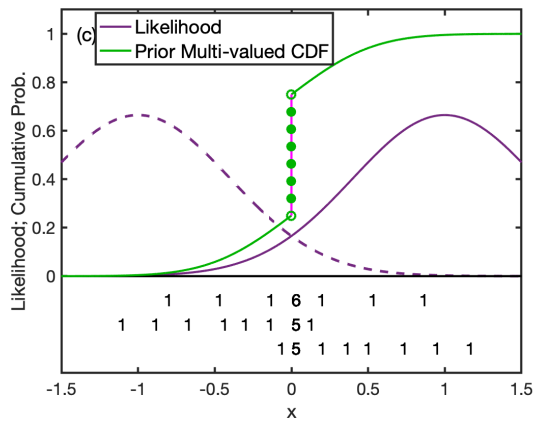
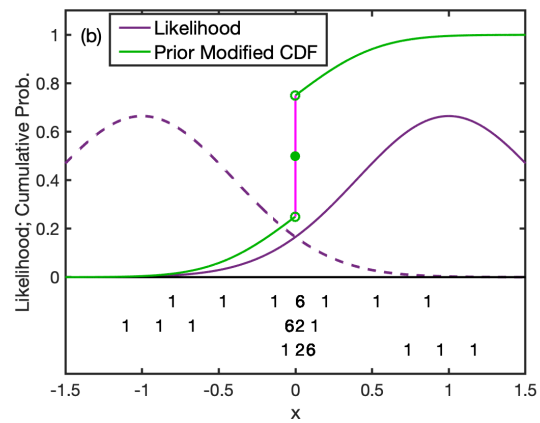
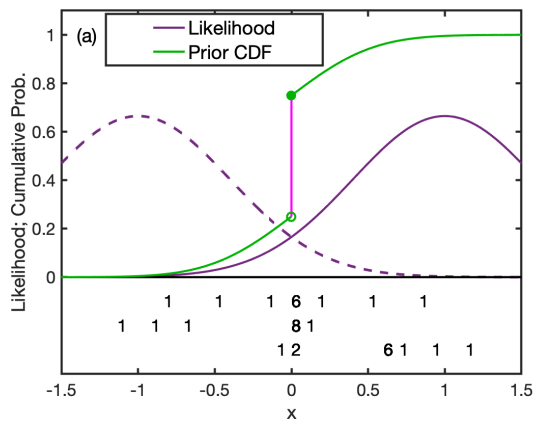
CDF Type	Observation	Ens. Mean	Ens. Std. Dev.	Ens. Skewness
Standard	-1	-0.2222	0.6428	0.7126
Standard	+1	0.4444	0.4965	-0.5900
Modified	-1	-0.3333	0.5863	0.5482
Modified	+1	0.3333	0.5863	-0.5482

788

789 Table A1: Analysis ensemble mean, standard deviation, and skewness for the idealized
 790 assimilation illustrated in Fig. A1 for the standard CDF, and the modified CDF, for observations
 791 with value -1 and +1.

792

793 Similar bias problems occur when applying the QCEFF to a mixed distribution. In this case, the
 794 standard CDF is defined in (9). Again, the quantiles computed with the CDF are biased towards
 795 positive values because of the definition at the discrete jumps. Figure A2 shows the results of
 796 assimilating a pair of observations with an ensemble that leads to a mixed distribution. The 12-
 797 member prior ensemble has 6 members with value exactly 0, and 6 other members that are
 798 distributed symmetrically around 0. It is assumed that the 6 duplicate members result from a
 799 discrete part of the distribution while the remaining members are part of a continuous normal
 800 distribution. The two observation likelihoods both have standard deviation 0.6 with means -1
 801 and +1 respectively. Figure A2a shows the CDF for this prior along with the prior and analysis
 802 ensembles. The results of assimilating the two observations are far from symmetric as also
 803 shown by the large differences in the magnitude of the mean, standard deviation and skewness
 804 in Table A2.



805
806

807

808 Figure A2: Results of idealized assimilation for a 12-member ensemble prior sampled from a
809 discrete distribution. The prior ensembles are the same in each panel, depicted by the top row
810 of integers below the axis with 6 members having the value 0. Two observation likelihoods are
811 shown in purple with observed value -1 (dashed) and +1 (solid). The standard prior CDF is
812 shown in green in (a), the modified CDF in (b), and the multiple value CDF in (c). The rows of
813 numbers below the prior ensemble depict the number and position of analysis ensemble
814 members for the -1 observation (middle row) and +1 observation (lower row). The magenta
815 vertical segments show the value of the generalized inverse of the CDF for values in the discrete
816 jumps and is the same in each figure.
817

818 A modified CDF can be defined as given in (10). The value of the modified CDF at the jump is
819 placed at the midpoint of the jump as shown in Fig. A2b. In this case, the analysis ensembles for
820 the two observations are symmetric around 0 as expected for an unbiased assimilation
821 algorithm and the first three moments in Table A2 have the same magnitude. As is the case for
822 the standard CDF, prior ensemble members with the same value (0 in this example) all have the
823 same quantile. The QCEFF treats them all identically so the analysis values of these six are
824 always identical.

825
826 Other alternative definitions for modified quantiles are also possible. Figure A2c shows a case
827 where a different quantile is assigned to each duplicate valued ensemble member. In this case,
828 referred to as the multi-valued CDF, the quantiles uniformly partition the jump in the standard
829 CDF. Note that this is not technically a function since it is multivalued at the jump. Figure A2c
830 shows that the results of assimilations with this multivalued modified quantile are unbiased.
831 However, they are qualitatively different from the results for the first modified quantile
832 method. In this case, prior ensemble members with the same value can have analysis ensemble
833 members with different values. The multi-valued CDF analysis retains a number of ensemble
834 members with the same value, zero, as the prior discrete value (Fig. A2c). The modified CDF
835 analysis (Fig. A2b) assigns new values to all of the zero-valued prior members, essentially
836 shifting the discrete value in the mixed distribution. Limited testing with several low-order
837 models suggested that the multi-valued CDF performed less well than the modified CDF. Similar
838 choices on posteriors have to be made for Gaussian anamorphosis applications (Moha
839 Gharamti, personal communication) and synergies with that work should be explored. Future

840 work will be needed to explore this in more detail for a variety of applications. Also note that it
 841 would be possible to define alternative CDFs that make different choices for boundary versus
 842 interior points with discrete probability. For instance, the values of 0 and 1 for sea ice
 843 concentration are physically special, while a discrete value in the interior might have different
 844 significance. Again, future work is needed to explore this issue.

845

CDF Type	Observation	Ens. Mean	Ens. Std. Dev.	Ens. Skewness
Standard	-1	-0.2214	0.4308	-1.2942
Standard	+1	0.5218	0.3727	-0.4031
Modified	-1	-0.2564	0.4123	-1.2692
Modified	+1	0.2564	0.4123	1.2692
Multi-valued	-1	-0.3023	0.4147	-0.8966
Multi-valued	+1	0.3023	0.4147	0.8966

846

847 Table A2: Analysis ensemble mean, standard deviation, and skewness for the idealized
 848 assimilation illustrated in Fig. A2 for the standard CDF, the modified CDF, and the multi-valued
 849 CDF for observations with value -1 and +1.

850

851

852 REFERENCES

853

854 Amezcua, J. and P. J. Van Leeuwen, 2014: Gaussian anamorphosis in the analysis step of the
 855 EnKF: a joint state-variable/observation approach. *Tellus A: Dynamic Meteorology and*
 856 *Oceanography*, 66. DOI: 10.3402/tellusa.v66.23493.

857 Anderson, J. L., 1996: A Method for Producing and Evaluating Probabilistic Forecasts from
 858 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).

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

862 Anderson, J. L., 2003: A local least squares framework for ensemble filtering. *Mon. Wea. Rev.*,
863 **131**, 634-642.

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

866 Anderson, J. L., 2019: A nonlinear rank regression method for ensemble Kalman filter data
867 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)
868 [0448.1](https://doi.org/10.1175/MWR-D-18-0448.1).

869 Anderson, J. L., 2022: A quantile-conserving ensemble filter framework. Part I: Updating an
870 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)
871 [21-0229.1](https://doi.org/10.1175/MWR-D-21-0229.1).

872 Anderson, J. L., 2023: A quantile-conserving ensemble filter Framework. Part 2: Updating an
873 observed variable. *Mon. Wea. Rev.*, **151**, 2759–2777, [https://doi.org/10.1175/MWR-D-23-](https://doi.org/10.1175/MWR-D-23-0065.1)
874 [0065.1](https://doi.org/10.1175/MWR-D-23-0065.1).

875 Anderson J. L., and S. L. Anderson, 1999: A Monte Carlo implementation of the nonlinear
876 filtering problem to produce ensemble assimilations and forecasts. *Mon. Wea. Rev.*, **127**,
877 2741–2758.

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

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

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

886 Beal, D., P. Brasseur, J. M. Brankart, Y. Ourmieres, and J. Verron, 2010: Characterization of
887 mixing errors in a coupled physical biogeochemical model of the North Atlantic:
888 Implications for nonlinear estimation using Gaussian anamorphosis. *Ocean Sci.*, **6**, 247–
889 262.

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

892 Bishop, C. H., 2016: The GIGG-EnKF: Ensemble Kalman filtering for highly skewed non-negative
893 uncertainty distributions. *Q. J. R. Meteorol. Soc.*, **142**, 1395-1412. doi:10.1002/qj.2742.

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

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

899 Chan, M., J. L. Anderson, and X. Chen, 2020: An efficient bi-Gaussian ensemble Kalman filter
900 for satellite infrared radiance data assimilation. *Mon. Wea. Rev.*, **148**, 5087–
901 5104, <https://doi.org/10.1175/MWR-D-20-0142.1>.

902 Doron, M., P. Brasseur, J. M. Brankart, S. N. Losa, and A. Melet, 2013: Stochastic estimation of
903 biogeochemical parameters from Globcolour ocean colour satellite data in a North
904 Atlantic 3D ocean coupled physical–biogeochemical model. *J. Marine Systems*, **117**, 81–
905 95.

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

908 Gharamti, M., 2018: Enhanced adaptive inflation algorithm for ensemble filters. *Mon. Wea.*
909 *Rev.*, **146**, 623–640, <https://doi.org/10.1175/MWR-D-17-0187.1>.

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

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

916 Hamill, T. M., 2001: Interpretation of rank histograms for verifying ensemble forecasts. *Mon.*
917 *Wea. Rev.*, **129**, 550–560, [https://doi.org/10.1175/1520-0493\(2001\)129<0550:IORHFV>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0550:IORHFV>2.0.CO;2).

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

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

923 Kurosawa, K., and J. Poterjoy, 2021: Data assimilation challenges posed by nonlinear
924 operators: A comparative study of ensemble and variational filters and smoothers. *Mon.*
925 *Wea. Rev.*, **149**, 2369–2389, <https://doi.org/10.1175/MWR-D-20-0368.1>.

926 Lien, G.Y., E. Kalnay, and T. Miyoshi, 2013: Effective assimilation of global precipitation:
927 Simulation experiments. *Tellus A* 65, DOI: 10.3402/tellusa.v65i0.19915

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

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

932 Riedel, C., 2023: In press

933 Simon, E., and L. Bertino, 2012: Gaussian anamorphosis extension of the DEnKF for combined
934 state parameter estimation: application to a 1D ocean ecosystem model. *J. Marine Syst.*,
935 **89**, 1–18.

936 Suhaila, J., K. Ching-Yee, Y. Fadhilah and F. Hui-Mean, 2011: Introducing the mixed distribution
937 in fitting rainfall data. *Open Journal of Modern Hydrology*, **1**, 11-22.
938 doi: 10.4236/ojmh.2011.12002.

939 Van Leeuwen, P. J., 2009: Particle filtering in geophysical systems. *Mon. Wea. Rev.*, **137**, 4089-
940 4114. doi.org/10.1175/2009MWR2835.1

941 Van Leeuwen, P. J., H. R. Künsch, L. Nerger, R. Potthast, and S. Reich, 2019: Particle filters for
942 high-dimensional geoscience applications: A review. *Quart. J. Roy. Met. Soc.*, **149**, 2335-
943 2365. doi.org/10.1002/qj.3551

944 Wilks, D. S., 2019: Indices of rank histogram flatness and their sampling properties. *Mon. Wea.*
945 *Rev.*, **147**, 763–769, <https://doi.org/10.1175/MWR-D-18-0369.1>.