

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

A Hybrid MPI/OpenMP Parallel Algorithm and Performance Analysis for an Ensemble Square Root Filter Designed for Multi-scale Observations

Yunheng Wang¹, Youngsun Jung¹, Timothy A. Supinie^{1,2} and Ming Xue^{1,2}

Center for Analysis and Prediction of Storms¹ and School of Meteorology²
University of Oklahoma, Norman Oklahoma 73072

August 2012

Submitted to J. Atmospheric and Oceanic Technology

Revised January 16, 2013

Corresponding author address:
Ming Xue
Center for Analysis and Prediction of Storms
University of Oklahoma,
120 David L. Boren Blvd, Norman OK 73072
mxue@ou.edu

31

Abstract

32

33

34

35

36

37

38

39

A hybrid parallel scheme for the ensemble square root filter (EnSRF) suitable for parallel assimilation of multi-scale observations including those from dense observational networks such as those of radar is developed based on the domain decomposition strategy. The scheme handles inter-node communication through message passing interface (MPI), and the communication within shared-memory nodes via Open Multi-Processing (OpenMP) threads; it also supports pure MPI and pure OpenMP modes. The parallel framework can accommodate high-volume remote-sensed radar (or satellite) observations as well as conventional observations that usually have larger covariance localization radii.

40

41

42

43

44

45

46

47

48

49

50

The performance of the parallel algorithm has been tested with simulated and real radar data. The parallel program shows good scalability in pure MPI and hybrid MPI/OpenMP modes, while pure OpenMP runs exhibit limited scalability on a symmetric shared-memory system. It is found that in MPI mode, better parallel performance is achieved with domain decomposition configurations in which the leading dimension of the state variable arrays is larger, because this configuration allows for more efficient memory access. Given a fixed amount of computing resources, the hybrid parallel mode is preferred to pure MPI mode on supercomputers with nodes containing shared-memory cores. The overall performance is also affected by factors such as the cache size, memory bandwidth, and the networking topology. Tests with a real data case with a large number of radars confirm that the parallel data assimilation can be done on a multi-core supercomputer with a significant speedup compared to the serial data assimilation algorithm.

51 **1. Introduction**

52 With significant advances in computing power in recent years, advanced data
53 assimilation (DA) techniques, such as the ensemble Kalman filter (EnKF) (Evensen 1994;
54 Evensen and Leeuwen 1996; Burgers et al. 1998; Houtekamer and Mitchell 1998; Anderson
55 2001; Bishop et al. 2001; Whitaker and Hamill 2002; Evensen 2003; Tippett et al. 2003) and
56 four-dimensional variational (4DVAR) (e.g., Le Dimet and Talagrand 1986; Courtier and
57 Talagrand 1987; Sun and Crook 1997; Gao et al. 1998; Wu et al. 2000; Caya et al. 2005), are
58 becoming more popular in both operational and research communities. However, they both incur
59 a high computational cost, one of the biggest constraints for their operational applications at very
60 high resolutions. Between EnKF and 4DVAR, the EnKF method appears to be more attractive
61 for convective scale numerical weather prediction (NWP), where nonlinear physical processes
62 have critical roles. EnKF can also provide a natural set of initial conditions for ensemble
63 forecasting. EnKF has been applied at scales ranging from global to convective and has produced
64 encouraging results (e.g., Snyder and Zhang 2003; Dowell et al. 2004; Tong and Xue 2005,
65 hereafter TX05; Xue et al. 2006; Jung et al. 2008; Buehner et al. 2010; Dowell et al. 2011;
66 Hamill et al. 2011; Snook et al. 2011; Jung et al. 2012).

67 Among variants of EnKF, the ensemble square-root Kalman filter (EnSRF) of Whitaker
68 and Hamill (2002) is widely used in convective-scale DA studies involving radar data. The
69 EnSRF, as well as the similar ensemble adjustment Kalman filter (EAKF, Anderson 2003) and
70 the classic perturbed-observation EnKF algorithm (Evensen 2003), is an observation-space-
71 based algorithm in which observations are assimilated one after another. Because of the
72 sequential nature of the EnSRF (and EAKF and classic EnKF), parallelization of the algorithm at
73 the observation level is not straightforward. It is possible to parallelize at the state variable level,

74 i.e., to perform the updating of the state variables in parallel because each observation updates
75 many state variables within the covariance localization radius of the EnSRF, and these operations
76 are independent. Such parallelization can be easily achieved on shared-memory platforms via
77 OpenMP directives, and is done with the Advanced Regional Prediction System (ARPS, Xue et
78 al. 2003) EnSRF system (e.g., Xue et al. 2006; Jung et al. 2008). A processing element (PE) on a
79 shared-memory or distributed-memory platform is an individual processor with single-core
80 processors or a processor core on multi-core CPUs. Each PE generally supports only a single
81 process or a single thread. The number of PEs available on shared-memory nodes (the term
82 “processing unit,” abbreviated PU, will be used to refer to a shared-memory node) usually limits
83 the scale of shared-memory parallelization (SMP) and the number of state variables that can be
84 updated simultaneously. Distributed-memory parallelization (DMP) via the Message Passing
85 Interface (MPI) library would allow the use of much larger computers, which are essential for
86 very-high-resolution DA and NWP over large domains (Xue et al. 2007).

87 Anderson and Collins (2007, hereafter AC07) proposed a modification to the standard
88 EAKF algorithm that is also applicable to EnSRF. In their algorithm, multiple observation priors
89 (background converted to observed quantities via observation operators) are first calculated in
90 parallel, and the observation priors corresponding to as yet unused observations are updated by
91 the filter together with the state vector, allowing easier parallelization at the state vector level
92 (for a given observation, multiple elements in the state vector are updated in parallel). However,
93 its state update procedure requires broadcasting the observation priors from one PU to the rest,
94 and more importantly, the processing of observations is still serial. Because of this, the algorithm
95 does not scale well when the number of PUs increases to the point where the cost of
96 communication starts to dominate or when the ratio of the number of observations to that of state

97 variables is large. Other parallel approaches have also been proposed by Keppenne and
98 Rienecker (2002) and Zhang et al. (2005). While both methods utilize domain decomposition,
99 they differ in whether communication among PUs is allowed. Because there is no cross-PU
100 communication in the algorithm of Zhang et al. (2005), the analysis near the PU boundaries is
101 not the same as that of scalar implementation, which is a potentially serious drawback of their
102 algorithm. Keppenne and Rienecker (2002), on the other hand, allow observations in other PUs
103 to update the states in the current PU, but their communication cost is potentially very high
104 because message passing is executed many times to properly exchange information among PUs.

105 In this paper, we develop a new parallelization algorithm for EnSRF (also suitable for
106 other similar serial ensemble filters) that is especially suitable for dense observations that
107 typically use relatively small horizontal covariance localization radii. Most NWP models,
108 including the ARPS and the Weather Research and Forecasting model (WRF), use horizontal
109 domain decomposition for effective parallelization (Sathye et al. 1997; Michalakes et al. 2004).
110 A domain-decomposition-based parallel DA strategy is attractive because it can share much of
111 the parallelization infrastructure with the prediction model. If the DA system and prediction
112 model use the same number and configuration of subdomains, transfer of model grids between
113 the two systems will be more straightforward either through disk or within computer memory.
114 Furthermore, with typical ensemble DA systems, the state arrays are usually moved between the
115 prediction model and DA system through disk I/O within the DA cycles; such I/O can take more
116 than half of the total wall clock time within each cycle (Szunyogh et al. 2008), making high-
117 frequency assimilation of observations on large, high-resolution, grids prohibitively expensive.
118 Our eventual goal is to achieve data exchange through message passing within computer
119 memory, bypassing disk I/O altogether; adopting a domain decomposition parallelization

120 strategy would simplify this process. Finally, the domain decomposition strategy makes grid-
121 based calculations within the DA system, such as spatial interpolation, easier.

122 The domain-decomposition-based strategy we propose takes advantage of the relatively
123 small localization radii typically used by very dense observations within ensemble algorithms,
124 because observations that do not influence state variables at the same grid points can be
125 processed in parallel. More sparse conventional observations tend to require larger localization
126 radii (Dong et al. 2011) and are therefore more difficult to process in parallel. In this case, a
127 strategy similar to that of AC07 is taken, in which observations are processed serially but still
128 using the same decomposed domains. Parallelization can be achieved at the state variable level in
129 the case; in other words, different parallelization strategies can be used in combination, taking
130 advantage of the serial nature of the ensemble algorithms. Note that this approach scales well
131 only for observations whose localization radius is large enough to impact most of the grid points
132 in the model domain, unless additional steps are taken to balance the load, as in AC07.

133 In addition to domain-decomposition-based parallelization, we also want to take
134 advantage of SMP capabilities of multi-core compute nodes that are available on essentially all
135 large parallel systems of today. SMP among cores on the same node eliminates explicit data
136 transport among the cores, thus reducing communication costs and contention for interconnect
137 ports. By performing domain decomposition for the nodes while parallelizing across the PEs
138 (e.g., cores) on the same PUs (e.g., nodes), the decomposed domains can be larger relative to the
139 localization radii, increasing the chance that observations on different decomposed domains can
140 be processed independently.

141 For the EnSRF algorithm, SMP is easily achieved at the state variable level, because each
142 observation will need to update all state variables within its localization radius, and these update

143 operations are independent. Thus, the state variable update can be parallelized using OpenMP
144 directives applied to the loops over the state variables. The combination of MPI and OpenMP
145 strategies gives hybrid parallelization. This paper describes a hybrid parallel scheme
146 implemented for the ARPS EnSRF system. In addition, observation data are organized into
147 batches to improve the load balance when assimilating data from a number of radars.

148 This paper is organized as follows. Section 2 reviews the EnSRF formulation and briefly
149 describes the ARPS model used in timing experiments. Section 3 introduces the parallel
150 algorithms for high-density radar data and conventional observations separately. It also describes
151 the OpenMP/MPI hybrid strategy as well as the observation organization. Validation of the
152 parallel implementation and its performance are examined in section 4. A summary and
153 conclusions are presented in section 5.

154 **2. The ARPS ensemble DA system**

155 The ARPS (Xue et al. 2000; Xue et al. 2001; Xue et al. 2003) model is a general-purpose,
156 multi-scale prediction system in the public domain. It has a non-hydrostatic, fully compressible
157 dynamic core formulated in generalized terrain-following coordinates. It employs the domain
158 decomposition strategy in the horizontal for massively parallel computers (Sathye et al. 1997;
159 Xue et al. 2007), and has been tested through real-time forecasts at convection-
160 permitting/allowing resolutions for many years (e.g., Xue et al. 1996), including forecasts in
161 continental US (CONUS-scale) domains at 4 and 1 km grid spacing (e.g., Xue et al. 2011),
162 assimilating data from all radars in the WSR-88D radar network using a 3DVAR method.

163 As mentioned earlier, the current ARPS EnKF DA system (Xue et al. 2006) is primarily
164 based on the EnSRF algorithm of Whitaker and Hamill (2002). In addition, an asynchronous
165 (Sakov et al. 2010) four-dimensional EnSRF (Wang et al. 2013) has also been implemented. The

166 system includes capabilities for parameter estimation (Tong and Xue 2008), dual-polarimetric
 167 radar data assimilation (Jung et al. 2008), simultaneous reflectivity attenuation correction (Xue et
 168 al. 2009), and the ability to handle a variety of data sources (Dong et al. 2011). Additionally, it
 169 has been coupled with a double-moment microphysics scheme (Xue et al. 2010; Jung et al. 2012).
 170 To be able to apply this system to large, convection-resolving domains such as those used by
 171 ARPS 3DVAR for continental scale applications (e.g., Xue et al. 2011) and be able to assimilate
 172 frequent, high-volume observations, efficient parallelization of the system is essential.

173 Briefly, in EnSRF, the ensemble mean and ensemble deviations are updated separately.
 174 The analysis equations for ensemble mean state vector $\bar{\mathbf{x}}$ and the ensemble deviations \mathbf{x}'_i are,
 175 respectively,

$$176 \quad \bar{\mathbf{x}}^a = \bar{\mathbf{x}}^b + \boldsymbol{\rho} \circ \mathbf{K} \left[\mathbf{y}^o - H(\bar{\mathbf{x}}^b) \right], \quad (1)$$

$$177 \quad \mathbf{x}'_i{}^a = \beta(\mathbf{I} - \alpha \boldsymbol{\rho} \circ \mathbf{K}H) \mathbf{x}'_i{}^b \quad (2)$$

178 where \mathbf{K} is the Kalman gain and \mathbf{y}^o the observation vector. Subscript i denotes the ensemble
 179 member and ranges from 1 to N with N being the ensemble size. H is the forward observation
 180 operator that projects state variables to observed quantities, which can be nonlinear. Symbol \circ in
 181 the equations represents the Schur (element-wise) product and $\boldsymbol{\rho}$ is the localization matrix,
 182 containing localization coefficients that are typically functions of the distance between the
 183 observation being processed and the state variable being updated. The analysis background $\bar{\mathbf{x}}^b$
 184 projected into observation space, i.e., $H(\bar{\mathbf{x}}^b)$, is called the observation prior. Superscripts a , b ,
 185 and o denote analysis, background, and observation, respectively. State vector \mathbf{x} includes in our
 186 case the grid point values of the three wind components (u , v , w), potential temperature (θ),
 187 pressure (p), the mixing ratios of water vapor (q_v), cloud water (q_c), rain water (q_r), cloud ice (q_i),

188 snow (q_s), and hail (q_h). When a two-moment microphysics parameterization scheme is used, the
 189 total number concentrations for the 5 water and ice species are also part of the state vector (Xue
 190 et al. 2010). Background state vectors $\bar{\mathbf{x}}^b$ and \mathbf{x}_i^b are either forecasts from the previous
 191 assimilation cycle or the states updated by observations processed prior to the current one. The
 192 parameter β is the covariance inflation factor. Variable α is a factor in the square root algorithm
 193 derived by Whitaker and Hamill (2002),

$$194 \quad \alpha = \left[1 + \sqrt{\mathbf{R}(\mathbf{H}\mathbf{P}^b\mathbf{H}^T + \mathbf{R})^{-1}} \right]^{-1}. \quad (3)$$

195 Here, \mathbf{R} is observation error covariance matrix, \mathbf{P}^b the background error covariance matrix, and
 196 \mathbf{H} the linearized observation operator. The Kalman gain matrix \mathbf{K} is given by

$$197 \quad \mathbf{K} = \mathbf{P}^b\mathbf{H}^T (\mathbf{H}\mathbf{P}^b\mathbf{H}^T + \mathbf{R})^{-1}. \quad (4)$$

198 In the above, matrices $\mathbf{P}^b\mathbf{H}^T$ and $\mathbf{H}\mathbf{P}^b\mathbf{H}^T$, representing the background error covariance
 199 between the state variables and observation priors, and that between observation priors,
 200 respectively, are estimated from the background ensemble, according to

$$201 \quad \mathbf{P}^b\mathbf{H}^T = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i^b - \bar{\mathbf{x}}^b) \left[H(\mathbf{x}_i^b) - \overline{H(\mathbf{x}^b)} \right]^T, \quad (5)$$

$$202 \quad \mathbf{H}\mathbf{P}^b\mathbf{H}^T = \frac{1}{N-1} \sum_{i=1}^N \left[H(\mathbf{x}_i^b) - \overline{H(\mathbf{x}^b)} \right] \left[H(\mathbf{x}_i^b) - \overline{H(\mathbf{x}^b)} \right]^T. \quad (6)$$

203 The overbars in Eqs. (5) and (6) denote the ensemble mean. When a single observation is
 204 analyzed, $\mathbf{P}^b\mathbf{H}^T$ becomes a vector having the length of the state vector \mathbf{x} . In practice, due to
 205 covariance localization, all elements in $\mathbf{P}^b\mathbf{H}^T$ are not calculated; those for grid points outside the
 206 localization radius of a given observation are assumed to be zero. In fact, it is this assumption
 207 that makes the design of our parallel algorithm practical; *observations whose domains of*
 208 *influence (as constrained by the covariance localization radii) do not overlap can be analyzed*

209 *simultaneously*. Another basic assumption with this algorithm (and most atmospheric DA
210 algorithms) is that observation errors are uncorrelated, so that observations can be analyzed
211 sequentially in any order. When the observations are processed serially, one at a time, the
212 observation error covariance matrix \mathbf{R} reduces to a scalar, as does matrix $\mathbf{HP}^b\mathbf{H}^T$. In this case,
213 $\mathbf{HP}^b\mathbf{H}^T$ is the background error variance at the observation point.

214 After an observation is analyzed based on Eqs. (1)-(6), the analyzed ensemble states \mathbf{x}_i^a
215 ($i=1\cdots N$), the sum of ensemble mean and deviations, become the new background states \mathbf{x}_i^b
216 for the next observation, and the analysis is repeated until all observations at a given time are
217 analyzed. An ensemble of forecasts then proceeds from the analysis ensemble until the time of
218 new observation(s); at that time the analysis cycle is repeated.

219 **3. The parallel algorithm for EnSRF**

220 For convective-scale weather, Doppler weather radar is one of the most important
221 observing platforms. The US National Weather Service (NWS) operates a network of over 150
222 Weather Surveillance Radar-1988 Doppler (WSR-88D) radars that continuously scan the
223 atmosphere, at a rate of one full volume scan every 5-10 minutes, producing radial velocity and
224 reflectivity data. One volume scan in precipitation mode typically contains 14 elevations with
225 approximately several million observations every 5 minutes.

226 The number of conventional observations, such as surface station measurements, upper
227 air soundings, and wind profiler winds, is small compared to radar observations; because they
228 typically represent weather phenomena of larger scales, their assimilation in EnKF typically uses
229 larger covariance localization radii, and therefore their influence reaches larger distances (Dong
230 et al. 2011). Because of the different characteristics of each data type, different parallel strategies
231 are employed for conventional and radar data.

232 *a. The parallel algorithm for high-density observations with small covariance localization radii*

233 The algorithm partitions the entire analysis domain into subdomains defined by the
234 number of participating MPI processes in the horizontal x and y directions. No decomposition is
235 performed in the vertical direction, and therefore, state variables are always complete in the
236 vertical columns. High-density radar observations (and other high-resolution observations
237 including those of satellite) are distributed to each subdomain according to their physical
238 locations. Fig. 1 illustrates an analysis domain that is partitioned into 4 physical subdomains
239 horizontally, to be handled by 4 PUs in the computing system. Each computational domain is
240 comprised of the physical subdomain (in darker gray for P1, separated with thick solid lines) and
241 extended boundary ‘halo’ zones surrounding the physical subdomain (in light gray for P1,
242 bounded by thin lines); the physical domain and the boundary halo zones combined together are
243 called computational subdomains. The width of the extended boundary halo zone for the DA
244 system is typically larger than the halo zone or ‘ghost cells’ needed for boundary condition
245 exchanges in parallel NWP models based on domain decomposition (e.g., Sathye et al. 1997).
246 The width of the halo zone in the ARPS model, for example, is only one grid interval on each
247 boundary.

248 The extended boundary zone on each side must be at least as wide as the maximum
249 localization radius (R) of observations handled by the algorithm in the subdomain. For radar
250 observations, R is usually equal to a few grid intervals. Each physical subdomain is further
251 divided into 4 patches that are separated by bold dashed lines in Fig. 1, and these patches are
252 labeled S1, S2, S3 and S4, respectively. The horizontal width of patch S2 and the vertical height
253 of patch S3 must be at least $2R$. The rest of the physical domain is assigned to patches S1 and S4
254 as in Fig. 1, and their horizontal width and height also must be at least $2R$. Thus, the width of the

255 physical subdomain must be larger than $4R$ for the algorithm to work. All other subdomains in
256 Fig. 1 are divided following the same patch pattern. Such a patch division assures that patches
257 with the same label in adjacent subdomains are at least $2R$ apart, so observations in any one
258 patch do not affect grid points in the same patch on other PUs and thus, they can be analyzed in
259 parallel. In other words, no two observations that are being analyzed in parallel will influence the
260 same grid point. In practice, we want to make patch S1 as large as possible, increasing the
261 chance that any two observations can be processed independently (see below). Thus, the width of
262 S2 and the height of S3 are assigned the minimum possible size of $2R$ (see Fig. 1), which leaves
263 the majority of the subdomain to patch S1.

264 The EnKF DA over the analysis domain is performed in 4 sequential steps for
265 observations within S1, S2, S3 and S4. In the first step, only observations within S1 on all PUs
266 are assimilated in parallel while observations on each S1 patch are assimilated sequentially. Let
267 P be the number of PUs. Then, there can be at most P observations being assimilated in parallel
268 at any time. After all observations located within S1 are assimilated, MPI communications are
269 required to properly update state variables at grid points within the extended boundary zones that
270 are shared with neighboring PUs. The same procedure is then repeated for observations within
271 S2, S3 and S4 in steps 2, 3, and 4.

272 The assimilation of observations within the same-labeled patches from all PUs can be
273 done in parallel because: 1) the grid points influenced by the observations analyzed in parallel
274 are separated far enough without overlap, and 2) the ensemble state arrays are extended beyond
275 the physical subdomain, so that the influence on state grids by observations within each
276 subdomain can be passed to its neighbor PUs with MPI communications. Best load balancing is
277 realized if the same-labeled patches contain the same number of observations so that all PUs can

278 complete each analysis step in approximately the same time. In practice, however, the number of
279 observations on each subdomain is usually different due to uneven spatial distribution of
280 observations (and of observation types). One way to improve parallelism is to make one patch
281 (S1 in our system) as large as possible, which increases the number of observations that can be
282 processed independently and improves the load balance. Assimilation of observations on S2, S3
283 and S4 may not be well balanced. However, because they tend to be smaller and contain fewer
284 observations, their effect on the assimilation time tends to be small.

285 Since high-density observations, such as radar data, usually assume relatively small
286 localization radii, the constraint that the width of the physical subdomain should be at least $4R$ in
287 each direction usually does not become a major problem, especially when the DA domain is
288 large. When a hybrid MPI-OpenMP parallelization strategy is used this problem can be further
289 alleviated (see later). While the proposed algorithm is valid for most meteorological observations
290 that can assume a small localization radius, certain ‘integral observations’ such as radar
291 reflectivity with path-integrated attenuation effect (e.g., Xue et al. 2009) and GPS slant-path
292 water vapor (e.g., Liu and Xue 2006) pose special challenge for the serial EnSRF algorithm in
293 general since their observation operators are non-local (Campbell et al. 2010).

294 *b. The parallel algorithm for conventional observations with large covariance localization radii*

295 Currently supported conventional observations in the ARPS EnKF system include surface
296 station, upper air sounding, wind profiler, and aircraft observations. Since the covariance
297 localization radii applied to these observations are usually large, the width of the extended
298 boundary zones described in section 3a would be impractical for these data, unless the
299 decomposed subdomains are much larger than the localization radii. This is usually only true
300 when a small number of subdomains is used. Therefore, we design and implement an alternative

301 algorithm for this type of observations. Because the number of conventional (or any other
302 coarse-resolution) observations is typically much smaller than the number of (dense) radar
303 observations, we can afford to process the observations serially while trying to achieve
304 parallelism at the state variable level, similar to the strategy taken by AC07.

305 In our current implementation, conventional observations within the entire analysis
306 domain are broadcast to all PUs and assimilated one by one. Only the PU containing the
307 observation to be analyzed computes the observation prior; it then broadcasts the observation
308 prior ensemble, $H(\mathbf{x}_i)$, to all other PUs. The state variables within the covariance localization
309 radius of this observation are updated simultaneously on each PU that carries the state variables
310 (Fig. 2). Since we do not need extra boundary zones, state variable updating occurs within the
311 computational subdomains of the original NWP model. However, a set of MPI communications
312 between PUs is still needed right after the analysis of each observation to update the state
313 variables within the halo zone to facilitate the spatial interpolation involved in observation
314 operators. These steps are repeated until all observations are assimilated.

315 Our current implementation does not pre-calculate $H(\mathbf{x})$ or update $H(\mathbf{x})$ as part of the
316 extended state vector as AC07 does, and we use a regular domain decomposition strategy to
317 distribute the state variables across the PUs. This implementation will have load balance issues
318 for conventional observations, especially when the covariance localization radii of these
319 observations are small relative to the size of the entire model domain. AC07 mitigates this
320 problem by distributing the state variables across PUs as heterogeneously as possible, i.e., by
321 distributing neighboring grid points across as many PUs as possible. Such an irregular
322 distribution of state variables makes it difficult to implement grid-point-based treatments within
323 the EnKF algorithms. The $H(\mathbf{x})$ pre-calculation and update strategy employed by AC07 allows

324 simultaneous calculation of observation priors. This can be an option in a future implementation;
325 in fact, the 4D EnSRF algorithm implemented by Wang et al. (2013) employs this strategy.

326 *c. Hybrid MPI-OpenMP parallelization*

327 All current supercomputers use compute nodes with multiple shared-memory cores. The
328 original ARPS EnSRF code supports OpenMP parallelization via explicit loop-level directives at
329 the state-variable-update level (Xue et al. 2006). Thus, it is straightforward to employ a hybrid
330 technique, using SMP among cores on the same node and DMP via MPI across nodes. Doing so
331 can reduce explicit data communication within nodes and allow for larger S1 patches within the
332 decomposed domains on each PU (see Fig. 1). Our hybrid implementation is designed such that
333 each MPI process spawns multiple threads. Since message passing calls are outside of the
334 OpenMP parallel sections, they are parallel thread safe, i.e., only the master thread in a process
335 makes calls to MPI routines. The final program is flexible enough to run in MPI only, OpenMP
336 only, or in MPI/OpenMP hybrid modes, on a single-node workstation or supercomputers made
337 up of multiple nodes.

338 *d. Parallel strategy for assimilating data from multiple radars*

339 In the ARPS EnKF system, full-resolution radar observations in the radar coordinates are
340 usually mapped horizontally to the model grid columns during preprocessing (Brewster et al.
341 2005). The original ARPS EnSRF implementation processes data from one radar at a time,
342 sequentially. This is convenient because the data are stored in arrays for individual radars on
343 elevation levels (Xue et al. 2006). For data from the same radar, only a few parameters are
344 needed to describe the radar characteristics. However, because each radar typically covers only a
345 portion of the model domain, this procedure severely limits the scalability of the analysis system

346 due to load imbalances (see Fig. 3). Figure 3a illustrates a domain that contains six radars labeled
347 A through F. If this domain is decomposed into four subdomains, all PUs, except P1, will be idle
348 when data from radar A are assimilated. The same is true for radars B through F. To mitigate this
349 problem, we develop a procedure that merges radar data into composite sets or batches so that
350 data from multiple radars can be processed at the same time.

351 In the analysis program, all vertical levels of radar observations at each horizontal grid
352 location are stored continuously as a vector column. The most general approach is to store all
353 columns of radar data in a single dynamically allocated storage array or data structure while
354 keeping track of the radar characteristics associated with each column. Each column may contain
355 different numbers of available radar elevations. When overlapping coverage exists, the grid
356 columns covered by multiple radars will have multiple columns of data (see Fig. 3a). To keep
357 track of data in reference to the analysis grid, it is convenient to define arrays that have the same
358 dimensions as the model grid in the horizontal directions, but such arrays will only be able to
359 store no more than one column of data at each grid location unless the last dimension is defined
360 dynamically or pre-defined to be large enough. While for optimally tuned EnKF, the order in
361 which observations are assimilated should not matter, in practice, because the ensemble spread
362 can be reduced too much by observations processed earlier before covariance inflation is applied,
363 the order of observation processing sometimes do matter somewhat. For this reason, we group
364 the radar data into several batches, the number of which is no bigger than the maximum number
365 of radars covering the same spot anywhere in the analysis domain. For a radar network that is
366 designed to maximize spatial coverage, such as the WSR-88D radar network, this maximum is
367 usually a single digit number; i.e., anywhere in the network, less than 10 radars observe the same
368 column.

369 Fig. 3 shows the spatial coverage of three batches of data that add up to all columns of
370 data available; those three batches of observations will be processed in sequence. Within regions
371 having multiple radar coverage, the radar from which data will be first picked can be chosen
372 randomly or based on the order the data were input to the program. Alternatively, the data
373 columns from the closest radar can be picked first. The last option is more desirable, as it
374 removes the randomness of the algorithm. Finally, because the radar data are no longer organized
375 according to radar, additional two-dimensional arrays are needed to store parameters for each
376 data column. When only a few elevations within a radar volume scan are analyzed using short
377 (e.g., 1 to 2 minutes) assimilation cycles, the vertical dimension of the arrays storing the
378 composite data sets need only to be a few.

379 With the above implementation, the load balance is significantly improved for the first
380 composite data set. It should be noted that we usually assimilate reflectivity data even in
381 precipitation-free regions, which has the benefit of suppressing spurious storms (Tong and Xue
382 2005). We note that load imbalance does still exist with radial velocity data in the first group
383 since they are usually only available in precipitation regions; however, their numbers are usually
384 much smaller than the total number of reflectivity data. In addition, load imbalances usually exist
385 with the second group of data and above but again the volume of data in these groups is small
386 since they only exist in overlapping regions, and these regions are usually spread over the
387 assimilation domain.

388 **4. Algorithm verification and performance analysis**

389 *a. Verification of the parallelized code*

390 The domain partition and batch processing inevitably change the sequence of
391 observations being assimilated in the EnKF system. Theoretically, the order in which the

392 observations are processed does not matter for observations with uncorrelated errors, to the
393 extent that sampling error does not impact the results. In practice, the analysis results may differ
394 significantly if the filter is not properly tuned, where the tuning typically includes covariance
395 inflation and localization.

396 A set of experiments has been performed to investigate the effect of domain
397 decomposition on the analysis of simulated radar observations in an observing system simulation
398 experiment (OSSE) framework. Convective storms are triggered by five 4-K ellipsoidal thermal
399 bubbles with a 60-km horizontal radius and 4-km vertical radius in an environment defined by 20
400 May 1977 Del City, Oklahoma, supercell sounding (Ray et al. 1981). The model domain is
401 $300 \times 200 \times 16 \text{ km}^3$ with horizontal and vertical grid spacings of 1 km and 500 m, respectively.
402 Forty ensemble members are initiated at 3000 seconds of model time. The full state vector has
403 1.4×10^9 elements. Simulated radar observations from three radars are produced, using the
404 standard WSR-88D VCP (Volume Coverage Pattern) 11, which contains 14 elevation levels.
405 The total number of observations is approximately 6.7×10^5 from three volume scans spanning 5
406 minutes each. Radar DA is first performed at 5-minute intervals from 3300 to 5700 seconds,
407 using the original serial ARPS EnSRF code to provide an ensemble for subsequent parallel
408 assimilation tests. The Milbrandt and Yau (2005) double-moment microphysics scheme is used
409 in both truth simulation and in DA. The environment and model configurations that are not
410 described here can be found in Xue et al. (2010).

411 Three parallel DA experiments are then performed at 6000 seconds, one running in pure
412 MPI mode, one in pure OpenMP mode, and one in pure OpenMP mode but processing
413 observations serially in a reversed order. These experiments are referred as MPI, OMP_F, and
414 OMP_B (F for forward and B for backward), respectively. For each experiment, average RMS

415 errors for the state variables are computed against the truth simulation at the grid points where
416 truth reflectivity is greater than 10 dBZ. The RMS errors of MPI and OMP_B are normalized by
417 the RMS errors of OMP_F and shown in Fig. 4 for individual state variables. Most of the
418 normalized errors are very close to 1, and all of them are between 0.95 and 1.05 for MPI. Among
419 the variables, the total number concentration for rain water shows the largest variability,
420 probably due to the high sensitivity of reflectivity to the rain drop size distribution. In fact, the
421 normalized error for rain water number concentration is an outlier for OMP_B, reaching close to
422 1.25, much larger than the normalized error of about 1.05 for MPI. These results suggest that the
423 effect of the domain partition on the analysis is small, and the differences are within the range of
424 sampling uncertainties of the ensemble system.

425 With respect to the parallel code implementation for conventional data analysis, domain
426 decomposition does not change the sequence of the observation processing (see section 3b).
427 Therefore, identical results from experiments OMP_F and MPI are guaranteed. The results from
428 the experiments when simulated surface observations are also included are not shown here.

429 *b. Performance evaluation with OSSE experiments*

430 The performance of our parallel EnKF system is evaluated with radar DA benchmark
431 experiments on a Cray XT5 system (called Kraken) at National Institute of Computational
432 Science (NICS) at the University of Tennessee, which has 9408 total compute nodes with 12
433 cores each (6 cores per processor, 2 processors per node), giving a peak performance of 1.17
434 petaFLOPS. With Kraken, users can set the number of MPI processes per node (1-12), the
435 number of MPI processes per processor (1-6), and the number of cores (OpenMP threads) per
436 MPI process (1-12). A number of experiments with combinations of different numbers of MPI
437 processes, OpenMP threads, cores per node, and cores per processor have been performed to

438 examine the timing performance of various configurations. The same case described in section
439 4a is used for benchmarking.

440 First, the scalability of the OpenMP implementation is investigated as a reference. Since
441 each Kraken node contains only 12 cores, the maximum number of threads that can be used for
442 an OpenMP job is 12. The OpenMP implementation shows scalability up to 8 cores (see Table 1),
443 beyond which the reduction in wall clock time becomes minimal. One very likely reason is the
444 contention accessing shared memory and cache by different cores of the Opteron processors used.

445 To evaluate the performance of our MPI implementation, we ran several OpenMP and
446 MPI experiments on a single compute node. Table 1 lists the wall clock times and relative
447 speedups for these experiments. The experiment names follow the convention o[*total cores used*
448] for OpenMP and m[*nproc_x*] × [*nproc_y*] for MPI experiments, where *nproc_x* and *nproc_y*
449 denote the number of PUs corresponding to the decomposed domains in the *x* and *y* directions,
450 respectively. Generally, the OpenMP jobs perform better than their MPI counterparts using the
451 same number of cores when running on a single node due to the communication overhead of
452 MPI processes and possibly better load balance with OpenMP. It is also noticed that the wall-
453 clock time is heavily influenced by the domain partitioning configuration in the *x* and *y*
454 directions. For example, m02×01 takes almost 1.4 times longer than m01×02, although both use
455 the same number of cores. Since FORTRAN arrays are stored contiguously in the column-major
456 order in the computer memory, a run that has a smaller partition number in the *x* direction than
457 the *y* direction (e.g., m01×02) is better at taking advantage of the spatial locality of the data in
458 memory. This can accelerate data loading from main memory into cache and improve cache
459 reuse. Conversely, an inefficient partition can degrade the performance even when more system
460 resources are used. For example, m03×02 using 6 cores has a much smaller speed improvement

461 over $m01 \times 01$ than experiments using 4 cores or even some experiments using 2 cores. These
462 results suggest that finding the optimal domain decomposition is important in achieving the best
463 performance with the given system resources.

464 Table 2 shows performance data collected from pure MPI runs, and from hybrid
465 MPI/OpenMP experiments that run on 4 Kraken nodes. All experiments are named as following:
466 $m(h)[nproc_x] \times [nproc_y]_o[number\ of\ processes\ per\ node]_o[number\ of\ threads\ per\ process]$,
467 where “m” denotes MPI runs and “h” denotes hybrid runs. For MPI runs, the number of threads
468 per process is always 1. Thus, “o[number of threads per process]” is omitted from the notations
469 for all MPI runs in Table 2. Since each Kraken node contains two processors, the processes on
470 each node are distributed to those processors as evenly as possible in order to obtain the best
471 possible performance.

472 It is found that the domain partitioning again plays an important role for the DA system
473 performance. For example, experiments that use 20 cores in total on 4 compute nodes show large
474 variability in the execution time. Among these experiments, $m02 \times 10_05$ has the best
475 performance, suggesting that $m02 \times 10_05$ utilizes the system cache most efficiently and/or has
476 the least message-passing overhead given 20 cores. Generally, the MPI experiments using more
477 nodes perform better than those experiments with the same domain partitioning but using fewer
478 nodes. For example, $m01 \times 04$ in Table 1 that uses one compute node takes 2660 seconds to finish
479 while $m01 \times 04_01$ in Table 2 running on 4 compute nodes takes only 2343 seconds. This is
480 consistent with the observation that performance is improved as available cache size increases.
481 Adding more processes improves the performance on 4 compute nodes. As an example,
482 $m06 \times 08_12$ takes less time than those experiments using 40 cores or less. This is because more
483 observations can be processed in parallel in the $m06 \times 08_12$ experiment than the others, even

484 though MPI communication costs are higher than in the other experiments. However, as
485 observed before with OpenMP experiments, access contention for the memory bandwidth and
486 the cache sharing as more cores are used may impede the performance at some point. It suggests
487 that there is a tradeoff between the number of processes and available computing resources and,
488 therefore, finding optimal configurations for MPI runs may not be straightforward because it
489 depends on a number of hardware factors.

490 For the hybrid runs, the wall-clock times of `m01×04_01` (i.e. `h01×04_01o1`),
491 `h01×04_01o2`, `h01×04_01o4`, `h01×04_01o6`, `h01×04_01o8` and `h01×04_01o12` decrease
492 monotonically, in that order. The decreasing trend of wall-clock time with increasing number of
493 threads is consistently found in other similar sets of experiments. It is also found that the hybrid
494 runs are as sensitive as the MPI runs to the domain partitioning, available cache, and other
495 hardware configuration factors. A hybrid experiment can outperform or underperform the
496 corresponding MPI experiments using the same resources (number of cores and number of nodes)
497 depending on the configuration (Table 2 and 3). For example, the minimum wall-clock time with
498 8 cores from 4 nodes in hybrid mode is 1471 seconds, which is smaller than the minimum time
499 required by a MPI run with 8 processes on 4 nodes (2169 seconds) in Table 3. On the other hand,
500 `h01×04_01o12` takes 733 seconds, more than the 606 seconds of `m06×08_12`, which uses the
501 same resources. It is also observed that a larger improvement is achieved by the hybrid jobs with
502 a fewer number of threads. This is because observations are processed one by one with OpenMP
503 processes. By using more MPI processes rather than more OpenMP threads, we can assimilate
504 more observations simultaneously and, hence, improve the parallel efficiency (see section 4c for
505 more details). In addition, cache availability and memory access contention with a large number
506 of threads in the hybrid experiments also affect program performance.

507 *c. Performance evaluation with a real data application*

508 The parallel ARPS EnKF system is applied to the 10 May 2010 Oklahoma-Kansas
509 tornado outbreak case. Over 60 tornadoes, with up to EF4 intensity, affected large parts of
510 Oklahoma and adjacent parts of southern Kansas, southwestern Missouri, and western Arkansas
511 on that day. This real data case is run on an SGI UV 1000 cc-NUMA shared-memory system at
512 the Pittsburgh Supercomputing Center (PSC). The system, called Blacklight, is comprised of 256
513 nodes containing 2 eight-core Intel Xeon processors each; its theoretical peak performance is 37
514 teraFLOPS. The cc-NUMA architecture allows for SMP across nodes. Up to 16 terabytes (TB)
515 of memory can be requested for a single shared memory job, while hybrid jobs can access the
516 full 32 TB of system memory.

517 The EnSRF analyses are performed on a grid with 4 km horizontal grid spacing, using 40
518 ensemble members. The domain consists of 443×483×53 grid points, and the model state
519 includes three velocity components, potential temperature, pressure, and mixing ratios of water
520 vapor, and 5 water and ice species. A single-moment microphysics scheme is used. The state
521 vector has 4.9×10^9 elements. Observations of radar reflectivity and radial velocity from 35
522 radars are analyzed from 1705 UTC to 1800 UTC at 5-minute intervals. Fig. 5 presents a
523 comparison between the radar observation mosaic at 1800 UTC on 10 May 2010 and the
524 corresponding analysis results by the parallel ARPS EnSRF system. Overall, the analyzed
525 reflectivity exhibits a good fit to the observations in shape, structure, and intensity. The
526 exceptions are several echoes in Texas, southeast Montana, and northwest Colorado, which are
527 due to the incomplete radar coverage over those areas. Several timing benchmark analyses at
528 1710 UTC are performed. There are about 1.3×10^6 observations from the 35 radars at this time
529 (see Fig. 6), more than any of the other times in the analysis window.

530 Our parallel benchmark experiments are run in pure OpenMP, pure MPI and hybrid
531 MPI/OpenMP modes. In all cases, all cores on the compute nodes were fully utilized, either by
532 individual MPI processes or by OMP threads. The experiment names and their configurations are
533 listed in Table 4. Guided by the timing results on Kraken, experiments are designed to use the
534 most optimal configurations, i.e., with a larger number of PUs in the y direction than in the x
535 direction. Each experiment in Table 4 was repeated 7 times. Because the timing results on
536 Blacklight show up to 185% variability due to system load variations, the best timing results for
537 each case are selected and presented here. Fig. 7 shows the best timing results of each case as a
538 function of the number of cores used. Very large variations in run time were found to be
539 attributable to disk I/O on a large shared file system; I/O times are therefore excluded in Fig. 7 to
540 allow us to focus on the time spent on the analyses. The times with and without including
541 message passing are shown.

542 Both MPI and hybrid runs show good scalability according to Fig. 7, and they outperform
543 pure OpenMP runs by a large margin except for the case of 16 cores. Because each physical node
544 of Blacklight has only 16 cores, when more than 16 cores are used by OpenMP, the memory
545 access will be across different physical nodes; this clearly leads to reduced parallelization
546 efficiency with the OpenMP runs. Also, with pure OpenMP, the parallelization is limited to the
547 state variable level, meaning all observations have to be processed serially (i.e., not
548 parallelization at the observation level).

549 Fig. 7 also shows that, when using the same amount of total resources, the hybrid runs
550 generally outperform pure MPI runs when both analysis and message passing times are included.
551 For the same number of cores used, pure MPI runs implies more PUs, i.e., more message passing
552 requests. Even though the pure MPI mode may be able to parallelize more at the observation

553 level, the message passing overhead can reduce the benefit. Not surprisingly, the hybrid
554 OpenMP/MPI runs are better in terms of total computational time. Among the hybrid groups,
555 jobs with fewer threads hence more MPI processes seem to give better performances, in terms of
556 the analysis time. This suggests that assimilating observations in parallel via MPI processes gives
557 a greater benefit before the increased message passing overhead becomes overwhelming.

558 We have noticed that I/O can easily take 60% to 80% of the total wall-clock time with
559 experiments in which all data I/O were handled by a single MPI process or the master OpenMP
560 thread. This I/O time can be reduced by distributing I/O loads among the MPI processes (but not
561 among OpenMP threads). Therefore, our solution is to let each MPI process read and write data
562 within its own subdomain, in the form of “split files”. This improves I/O parallelization and also
563 reduces time needed for communicating gridded information across PUs. With split files, only
564 data within the extended boundary zones need to be exchanged with neighboring PUs. Due to the
565 large variations in the I/O times collected on Blacklight, we ran another set of experiments on a
566 supercomputer with more consistent I/O performance between runs. It consists of 2.0 GHz quad-
567 core Pentium4 Xeon E5405 processors, with 2 processors on each node. Tests with split files on
568 this system, corresponding to h04×08_01o8 (see above naming conventions), reveal that the
569 times spent on I/O and message passing are reduced (the latter because of the reduced exchanges
570 of gridded information across MPI processes); the total wall-clock time for I/O and message
571 passing for one experiment was reduced from 1231 seconds to 188 seconds using split files.

572 **5. Summary and conclusions**

573 A parallel algorithm based on the domain decomposition strategy has been developed and
574 implemented within the ARPS EnKF framework. The algorithm takes advantage of the relatively
575 small spatial covariance localization radii typically used by high-resolution observations such as

576 those of radar. Assuming that the maximum horizontal covariance localization radius of the
577 observations to be analyzed in parallel is R , the horizontal area of a decomposed physical
578 subdomain should be at least $4R \times 4R$. An additional boundary zone of width R is added to each
579 side of the physical subdomains to create enlarged computational subdomains, which facilitate
580 information exchanges between neighboring subdomains. Each subdomain is assigned to one
581 processing unit (PU), within which no MPI message passing is required. The subdomains are
582 then further divided up into 4 sub-patches, denoted S1 through S4. The width and height of each
583 patch are required to be at least $2R$ to ensure any two observations that may be processed in
584 parallel are well separated. In practice, the size of S1 is made as large as possible within its
585 subdomain to increase the probability that observations from different subdomains can be
586 processed in parallel.

587 Observations within the 4 patches are processed sequentially, but data in the patches with
588 the same label in different subdomains are processed simultaneously. Distributed-memory
589 parallelization is therefore achieved at the observation level. The patch division ensures that most
590 of the analysis work is done in parallel when processing observations within patches S1 of all
591 PUs. To handle the load imbalance issue when assimilating observations from many radars, the
592 observation arrays are organized into batches. The maximum number of batches is limited by the
593 maximum number of radars covering the same location anywhere in the analysis domain. Such
594 an observation organization improves the uniformity of observation distribution within the first
595 observation batch and thereby improves load balance.

596 Conventional data that use larger covariance localization radii are still processed serially.
597 State variables influenced by a particular observation are updated synchronously on the PUs
598 carrying those state variables.

599 The algorithm supports three parallel modes: pure OpenMP, pure MPI, and
600 MPI/OpenMP hybrid. Within the PUs with multiple cores, shared-memory parallelization can be
601 achieved via OpenMP at the state variable update level. OpenMP parallelization reduces
602 message passing overhead and allows for larger decomposed domains, making the 4R
603 requirement easier to satisfy.

604 It was first confirmed via OSSEs that changing the sequence of observation processing
605 due to domain decomposition has little impact on the analysis. Parallel DA benchmark
606 experiments are performed on a Cray XT5 machine. The OpenMP implementation shows
607 scalability up to 8 threads, beyond which memory and cache access contention limit further
608 improvement. MPI and OpenMP runs on a single compute node show that OpenMP
609 parallelization runs faster because of the lower communication overhead. MPI jobs with a
610 smaller number of partitions in the x direction than in the y direction exhibit better performance.
611 The same also applies to most of the hybrid jobs, although all hybrid jobs do not outperform the
612 corresponding MPI jobs.

613 A real data case involving 35 radars is tested on an SGI UV 1000 cc-NUMA system
614 capable of shared-memory programming across physical nodes. Poor scalability with pure
615 OpenMP is observed when more than one node is used, but both MPI and hybrid runs show good
616 scalability on this system. Excluding message passing time, pure MPI runs exhibit best
617 performance. When message-passing time is included, the hybrid runs generally outperform pure
618 MPI runs. For this real data case, the EnKF analysis excluding I/O can be completed within 4.5
619 minutes using 160 cores of the SGI UV 1000.

620 Given a fixed amount of resources, the hybrid jobs improve more over pure MPI jobs
621 with fewer numbers of threads. Because MPI processes realize parallelization at the observation

622 level, they are more efficient than OpenMP threads. However, there is a tradeoff between a
623 performance improvement due to the parallel processing of observations and degradation due to
624 increased message passing overhead. On the other hand, a pure OpenMP strategy for EnKF
625 shows good scalability on symmetric shared-memory systems but is limited by the number of
626 cores available on the individual node and by the physical memory available on the node. With
627 pure OpenMP, data I/O can only be handled by a single process, reducing the overall scalability.

628 The MPI/OpenMP hybrid strategy combines the strengths of both methods. However,
629 care must be taken when partitioning the domain, because the configuration of MPI domain
630 partitioning has a significant impact on the system performance. Given the same resources, jobs
631 with smaller numbers of partitions in the x direction tend to run faster because FORTRAN arrays
632 are stored in the column-major order in memory. Timing experiments have also shown that
633 determining the optimal decomposition configuration on a specific computing system is not
634 straightforward because the performance depends on factors such as the subdomain size in the x
635 and y directions, the number of cores on each node, the cache sizes and memory bandwidth
636 available to each core, and the networking topology across the nodes.

637 In all configurations, data I/O constituted a large portion of the execution time.
638 Experiments on a small dedicated Linux cluster show that the time spent on I/O and message
639 passing are reduced significantly by distributing I/O loads among the MPI processes with
640 MPI/OpenMP hybrid or pure MPI runs.

641 Although a data batching strategy is developed to reduce the load imbalance issue, further
642 improvement could be obtained through dynamic load balancing. Another problem is the low
643 resource utilization during inter-node communications because all threads are idle except one:
644 the master thread. The development of runtime management algorithms, for example the

645 Scalable Adaptive Computational Toolkit (SACT) (Li and Parashar 2007; Parashar et al. 2008),
646 are expected to decrease runtime of the application automatically with reduced efforts from
647 developers. Finally we point out that our parallel algorithm can be easily applied to other serial
648 ensemble-based algorithm such as EAKF and the classic EnKF.

649

650 *Acknowledgements:* This work was primarily supported by NSF grant OCI-0905040, AGS-
651 0802888, and NOAA Warn-on-Forecast Project under NA080AR4320904. Partial support was
652 also provided by NSF AGS-0750790, AGS-0941491, AGS-1046171, and AGS-1046081. We
653 acknowledge David O’Neal of Pittsburgh Supercomputing Center (PSC) for assistance with the
654 use of Tuning and Analysis Utilities (TAU) on a PSC Altix cluster early in the work.
655 Computations were performed at the Pittsburgh Supercomputing Center (PSC), and the National
656 Institute for Computational Sciences (NICS), and the OU Supercomputing Center for Education
657 and Research (OSCER).

658

659 **References**

- 660 Anderson, J. L., 2001: An ensemble adjustment Kalman filter for data assimilation. *Mon. Wea.*
661 *Rev.*, **129**, 2884-2903.
- 662 Anderson, J. L., 2003: A local least square framework for ensemble filtering. *Mon. Wea. Rev.*,
663 **131**, 634-642.
- 664 Anderson, J. L. and N. Collins, 2007: Scalable implementations of ensemble filter algorithms for
665 data assimilation. *J. Atmos. Ocean. Technol.* , **24**, 1452-1463.
- 666 Bishop, C. H., B. J. Etherton, and S. J. Majumdar, 2001: Adaptive sampling with the ensemble
667 transform Kalman filter. Part I: Theoretical aspects. *Mon. Wea. Rev.*, **129**, 420.
- 668 Brewster, K., M. Hu, M. Xue, and J. Gao, 2005: Efficient assimilation of radar data at high
669 resolution for short-range numerical weather prediction. *WWRP Int. Symp. Nowcasting*
670 *Very Short Range Forecasting*, CDROM 3.06.
- 671 Buehner, M., P. L. Houtekamer, C. Charette, H. L. Mitchell, and B. He, 2010: Intercomparison
672 of variational data assimilation and the ensemble Kalman filter for global deterministic
673 NWP. Part II: One-month experiments with real observations. *Mon. Wea. Rev.*, **138**,
674 1567-1586.
- 675 Burgers, G., P. J. v. Leeuwen, and G. Evensen, 1998: Analysis scheme in the ensemble Kalman
676 filter. *Mon. Wea. Rev.*, **126**, 1719-1724.
- 677 Campbell, W. F., C. H. Bishop, and D. Hodyss, 2010: Vertical covariance localization for
678 satellite radiances in ensemble Kalman filters. *Mon. Wea. Rev.*, **138**, 282-290.
- 679 Caya, A., J. Sun, and C. Snyder, 2005: A comparison between the 4D-VAR and the ensemble
680 Kalman filter techniques for radar data assimilation. *Mon. Wea. Rev.*, **133**, 3081–3094.

681 Courtier, P. and O. Talagrand, 1987: Variational assimilation of meteorological observations
682 with the adjoint equation. Part II: Numerical results. *Quart. J. Roy. Meteor. Soc.*, **113**,
683 1329-1347.

684 Dong, J., M. Xue, and K. K. Droegemeier, 2011: The analysis and impact of simulated high-
685 resolution surface observations in addition to radar data for convective storms with an
686 ensemble Kalman filter. *Meteor. Atmos. Phys.*, **112**, 41-61.

687 Dowell, D., F. Zhang, L. J. Wicker, C. Snyder, and N. A. Crook, 2004: Wind and temperature
688 retrievals in the 17 May 1981 Arcadia, Oklahoma supercell: Ensemble Kalman filter
689 experiments. *Mon. Wea. Rev.*, **132**, 1982-2005.

690 Dowell, D. C., L. J. Wicker, and C. Snyder, 2011: Ensemble Kalman filter assimilation of radar
691 observations of the 8 May 2003 Oklahoma City supercell: Influence of reflectivity
692 observations on storm-scale analysis. *Mon. Wea. Rev.*, **138**, 1152-1171.

693 Evensen, G., 1994: Sequential data assimilation with a nonlinear quasi-geostrophic model using
694 Monte Carlo methods to forecast error statistics. *J. Geophys. Res.*, **99**, 10143-10162.

695 Evensen, G., 2003: The ensemble Kalman filter: Theoretical formulation and practical
696 implementation. *Ocean Dynamics*, **53**, 343-367.

697 Evensen, G. and P. J. v. Leeuwen, 1996: Assimilation of geosat altimeter data for the Agulhas
698 current using the ensemble Kalman filter with a quasigeostrophic model. *Mon. Wea. Rev.*,
699 **124**, 85-96.

700 Gao, J., M. Xue, Z. Wang, and K. K. Droegemeier, 1998: The initial condition and explicit
701 prediction of convection using ARPS adjoint and other retrievals methods with WSR-
702 88D data. *12th Conf. Num. Wea. Pred.*, Phoenix AZ, Amer. Meteor. Soc., 176-178.

703 Hamill, T. M., J. S. Whitaker, M. Fiorino, and S. G. Benjamin, 2011: Global ensemble
704 predictions of 2009's tropical cyclones initialized with an ensemble Kalman filter. *Mon.*
705 *Wea. Rev.*, **139**, 668-688.

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

708 Jung, Y., M. Xue, and M. Tong, 2012: Ensemble Kalman filter analyses of the 29-30 May 2004
709 Oklahoma tornadic thunderstorm using one- and two-moment bulk microphysics
710 schemes, with verification against polarimetric data. *Mon. Wea. Rev.*, **140**, 1457-1475.

711 Jung, Y., M. Xue, G. Zhang, and J. Straka, 2008: Assimilation of simulated polarimetric radar
712 data for a convective storm using ensemble Kalman filter. Part II: Impact of polarimetric
713 data on storm analysis. *Mon. Wea. Rev.*, **136**, 2246-2260.

714 Keppenne, C. L. and M. M. Rienecker, 2002: Initial testing of a massively parallel ensemble
715 Kalman filter with the Poseidon isopycnal ocean general circulation model. *Mon. Wea.*
716 *Rev.*, **130**, 2951-2965.

717 Le Dimet, F. X. and O. Talagrand, 1986: Variational algorithms for analysis and assimilation of
718 meteorological observations: Theoretical aspects. *Tellus*, **38A**, 97-110.

719 Li, X. and M. Parashar, 2007: Hybrid Runtime Management of Space-Time Heterogeneity for
720 Parallel Structured Adaptive Applications. *IEEE Transactions on Parallel and*
721 *Distributed Systems (TPDS)*, **18**, 1202-1214.

722 Liu, H. and M. Xue, 2006: Retrieval of moisture from slant-path water vapor observations of a
723 hypothetical GPS network using a three-dimensional variational scheme with anisotropic
724 background error. *Mon. Wea. Rev.*, **134**, 933-949.

725 Michalakes, J., J. Dudhia, D. Gill, T. Henderson, J. Klemp, W. Skamarock, and W. Wang, 2004:
726 The Weather Research and Forecast Model: Software Architecture and Performance. .
727 *Proceedings, 11th ECMWF Workshop on the Use of High Performance Computing In*
728 *Meteorology*, Reading U.K.

729 Milbrandt, J. A. and M. K. Yau, 2005: A multi-moment bulk microphysics parameterization. Part
730 I: Analysis of the role of the spectral shape parameter. *J. Atmos. Sci.*, **62**, 3051-3064.

731 Parashar, M., X. Li, and S. Chandra, 2008: *Advanced Computational Infrastructure for Parallel*
732 *and Distributed Adaptive Applications*. Vol. To appear, John Wiley & Sons.

733 Ray, P. S., B. Johnson, K. W. Johnson, J. S. Bradberry, J. J. Stephens, K. K. Wagner, R. B.
734 Wilhelmson, and J. B. Klemp, 1981: The morphology of severe tornadic storms on 20
735 May 1977. *J. Atmos. Sci.*, **38**, 1643-1663.

736 Sakov, P., G. Evensen, and L. Bertino, 2010: Asynchronous data assimilation with the EnKF.
737 *Tellus*, **62**, 24-29.

738 Sathye, A., M. Xue, G. Bassett, and K. K. Droegemeier, 1997: Parallel weather modeling with
739 the Advanced Regional Prediction System. *Parallel Computing*, **23**, 2243-2256.

740 Snook, N., M. Xue, and J. Jung, 2011: Analysis of a tornadic mesoscale convective vortex based
741 on ensemble Kalman filter assimilation of CASA X-band and WSR-88D radar data. *Mon.*
742 *Wea. Rev.*, **139**, 3446-3468.

743 Snyder, C. and F. Zhang, 2003: Assimilation of simulated Doppler radar observations with an
744 ensemble Kalman filter. *Mon. Wea. Rev.*, **131**, 1663-1677.

745 Sun, J. and N. A. Crook, 1997: Dynamical and microphysical retrieval from Doppler radar
746 observations using a cloud model and its adjoint. Part I: Model development and
747 simulated data experiments. *J. Atmos. Sci.*, **54**, 1642-1661.

748 Szunyogh, I., E. J. Kostelich, G. Gyarmati, E. Kalnay, B. R. Hunt, E. Ott, and E. Satterfield,
749 2008: A local ensemble transform Kalman filter data assimilation system for the NCEP
750 global model. *Tellus*, **60A**, 113-130.

751 Tippett, M. K., J. L. Anderson, C. H. Bishop, T. M. Hamill, and J. S. Whitaker, 2003: Ensemble
752 square root filters. *Mon. Wea. Rev.*, **131**, 1485-1490.

753 Tong, M. and M. Xue, 2005: Ensemble Kalman filter assimilation of Doppler radar data with a
754 compressible nonhydrostatic model: OSS Experiments. *Mon. Wea. Rev.*, **133**, 1789-1807.

755 Tong, M. and M. Xue, 2008: Simultaneous estimation of microphysical parameters and
756 atmospheric state with radar data and ensemble square-root Kalman filter. Part II:
757 Parameter estimation experiments. *Mon. Wea. Rev.*, **136**, 1649–1668.

758 Wang, S., M. Xue, and J. Min, 2013: A four-dimensional asynchronous ensemble square-root
759 filter (4DEnSRF) and tests with simulated radar data. *Tellus*, In press.

760 Whitaker, J. S. and T. M. Hamill, 2002: Ensemble data assimilation without perturbed
761 observations. *Mon. Wea. Rev.*, **130**, 1913-1924.

762 Wu, B., J. Verlinde, and J. Sun, 2000: Dynamical and microphysical retrievals from Doppler
763 radar observations of a deep convective cloud. *J. Atmos. Sci.*, **57**, 262-283.

764 Xue, M., K. K. Droegemeier, and V. Wong, 2000: The Advanced Regional Prediction System
765 (ARPS) - A multiscale nonhydrostatic atmospheric simulation and prediction tool. Part I:
766 Model dynamics and verification. *Meteor. Atmos. Phy.*, **75**, 161-193.

767 Xue, M., M. Tong, and K. K. Droegemeier, 2006: An OSSE framework based on the ensemble
768 square-root Kalman filter for evaluating impact of data from radar networks on
769 thunderstorm analysis and forecast. *J. Atmos. Ocean Tech.*, **23**, 46–66.

770 Xue, M., K. K. Droegemeier, and D. Weber, 2007: Numerical prediction of high-impact local
771 weather: A driver for petascale computing. *Petascale Computing: Algorithms and*
772 *Applications*, Taylor & Francis Group, LLC, 103-124.

773 Xue, M., M. Tong, and G. Zhang, 2009: Simultaneous state estimation and attenuation correction
774 for thunderstorms with radar data using an ensemble Kalman filter: Tests with simulated
775 data. *Q. J. Roy. Meteor. Soc.*, **135**, 1409-1423.

776 Xue, M., Y. Jung, and G. Zhang, 2010: State estimation of convective storms with a two-moment
777 microphysics scheme and an ensemble Kalman filter: Experiments with simulated radar
778 data *Q. J. Roy. Meteor. Soc.*, **136**, 685-700.

779 Xue, M., D.-H. Wang, J.-D. Gao, K. Brewster, and K. K. Droegemeier, 2003: The Advanced
780 Regional Prediction System (ARPS), storm-scale numerical weather prediction and data
781 assimilation. *Meteor. Atmos. Phys.*, **82**, 139-170.

782 Xue, M., K. K. Droegemeier, V. Wong, A. Shapiro, K. Brewster, F. Carr, D. Weber, Y. Liu, and
783 D. Wang, 2001: The Advanced Regional Prediction System (ARPS) - A multi-scale
784 nonhydrostatic atmospheric simulation and prediction tool. Part II: Model physics and
785 applications. *Meteor. Atmos. Phys.*, **76**, 143-165.

786 Xue, M., K. Brewster, K. K. Droegemeier, V. Wong, D. H. Wang, F. Carr, A. Shapiro, L. M.
787 Zhao, S. Weygandt, D. Andra, and P. Janish, 1996: The 1996 CAPS spring operational
788 forecasting period: Realtime storm-scale NWP, Part II: Operational summary and
789 examples. *Preprint, 11th Conf. Num. Wea. Pred.*, Norfolk, VA, Amer. Meteor. Soc., 297-
790 300.

791 Xue, M., F. Kong, K. W. Thomas, Y. Wang, K. Brewster, J. Gao, X. Wang, S. J. Weiss, A. J.
792 Clark, J. S. Kain, M. C. Coniglio, J. Du, T. L. Jensen, and Y. H. Kuo, 2011: CAPS

793 Realtime Storm Scale Ensemble and High Resolution Forecasts for the NOAA
794 Hazardous Weather Testbed 2010 Spring Experiment. *24th Conf. Wea. Forecasting/20th*
795 *Conf. Num. Wea. Pred.*, Amer. Meteor. Soc., Paper 9A.2.

796 Zhang, S., M. J. Harrison, A. T. Wittenberg, A. Rosati, J. L. Anderson, and V. Balaji, 2005:
797 Initialization of an ENSO forecast system using a parallelized ensemble filter. *Mon. Wea.*
798 *Rev.*, **133**, 3176-3201.

799
800

801 **List of Figures**

802 Fig. 1. A schematic of the domain decomposition strategy for the analysis of high-density
803 observations, illustrated with 4 processing units (PUs, denoted by P1 through P4). Letters $i-l$
804 denote observations that are assumed to be equally spaced and letters $a-h$ indicate the
805 influence limits (as determined by the covariance localization radii of EnKF) of those
806 observations. In this example, observations i and l are far enough apart that they will not
807 influence any of the same state variables; they are among the observations that are analyzed
808 simultaneously in the first step of the procedure. Observations j and k are analyzed in the
809 second step, but they must be analyzed sequentially. Note that in practice, there will be
810 many more observations within patches S1 and S2 of subdomains P1 to P4 than shown in
811 the figure.

812 Fig. 2. A schematic for analyzing conventional data. Three steps are involved when analyzing
813 one observation whose location is denoted by a black dot in the figure: 1) PU14 computes
814 $H(x_i)$ (where i is the ensemble index); 2) $H(x_i)$ are broadcasted to all PUs; 3) state variables
815 x_i within the influence range of this observation (within the large circle) are updated in
816 parallel by the PUs that carry the state variables.

817 Fig. 3. Composite radar data batches organized such that within each batch, no more than one
818 column of data exists for each grid column. (a) Observations from six radars (A-F) with
819 their coverage indicated by the maximum range circles are remapped onto the model grid.
820 (b) Observations of the first batch, (c) observations of the second batch, and (d)
821 observations of the third batch. If there are more observations unaccounted for, additional
822 data batch(es) will be formed.

823 Fig. 4. RMS errors averaged over the grid points where truth reflectivity is greater than 10 dBZ
824 and normalized by the errors of experiment OMP_F. The state variables are the 16 ARPS
825 prognostic variables: three velocity components (u , v , and w), potential temperature (pt),
826 pressure (p), mixing ratios of water vapor (q_v), cloud water (q_c), rain water (q_r), cloud ice
827 (q_i), snow aggregate (q_s), hail (q_h) and their respective number concentrations (N_{lc} , N_{lr} , N_{li} ,
828 N_{ls} , and N_{lh} , associated with a two-moment microphysics scheme used).

829 Fig. 5. (a) The observed radar reflectivity mosaic and (b) the reflectivity field analyzed by the
830 parallel EnKF algorithm, at model grid level 20 at 1800 UTC 10 May 2010.

831 Fig. 6. The model domain and coverage of 35 WSR-88D radars with 230 km range rings for the
832 10 May 2010 real data test case.

833 Fig. 7. Wall clock times of the EnKF analyses as a function of the total number of compute cores
834 used, for the 10 May 2010 real data case in the analysis domain shown in Fig. 6, obtained
835 on the PSC Blacklight (an SGI UV 1000). OMP denotes pure OpenMP runs, MPI denotes
836 pure MPI runs, and H_o4, H_o8, and H_o16 denote hybrid runs with 4, 8, and 16 OpenMP
837 threads within each MPI process, respectively. In all cases, all cores on the compute nodes
838 were fully utilized, either by individual MPI processes or by OMP threads. Solid lines
839 denote the total time excluding message passing, and the dashed lines show the total times
840 including message passing. Data I/O times are excluded from all statistics.

841

842

843 Table 1. Timing comparisons of OpenMP experiments with MPI experiments on
844 one compute node. Speedup for OpenMP and MPI experiments are computed
845 relative to o1 and m01×01, respectively.
846

Experiment	Total number of cores used	Wall Clock Time (s)	Speedup
o1	1	6310	1.00
o2	2	3617	1.75
o4	4	2597	2.43
o6	6	1919	3.29
o8	8	1597	3.95
o12	12	1607	3.93
m01×01	1	6815	1.00
m01×02	2	3994	1.71
m02×01	2	5698	1.20
m01×04	4	2660	2.56
m02×02	4	3690	1.85
m04×01	4	2896	2.35
m03×02	6	4177	1.63
m02×04	8	2100	3.25
m04×02	8	2413	2.82

847

848

849

850 Table 2. Timing comparisons of pure MPI experiments with hybrid MPI/OpenMP experiments
 851 on 4 compute nodes. Speedup is computed relative to experiment m01 × 01 (6815 seconds in
 852 Table 1).
 853

Experiment	Total number of cores used	Wall Clock Time (s)	Speedup
m01×04_01	4	2343	2.91
m02×02_01		3577	1.91
m04×01_01		2750	2.48
m02×04_02	8	2169	3.14
m04×02_02		2330	2.92
m03×04_03	12	1575	4.33
m06×02_03		1699	4.01
m04×04_04	16	1327	5.14
m02×10_05	20	915	7.45
m10×02_05		1357	5.02
m04×05_05		1082	6.30
m05×04_05		1085	6.28
m03×08_06	24	880	7.74
m06×04_06		1049	6.50
m04×10_10	40	637	10.70
m10×04_10		720	9.47
m06×08_12	48	606	11.25
h01×04_01o2	8	1471	4.63
h01×04_01o4	16	1129	6.04
h01×04_01o6	24	831	8.20
h01×04_01o8	32	772	8.83
h01×04_01o12	48	733	9.30
h02×04_02o2	16	1200	5.68
h02×04_02o4	32	908	7.51
h02×04_02o6	48	709	9.61

854

855

856 Table 3. Comparison of the minimum time taken in hybrid mode with that in MPI mode using
857 the same number of cores on 4 compute nodes
858

Number of cores	Hybrid case	Minimum time (s)	MPI case	Minimum time (s)	Difference (s)
8	h01×04_01o2	1471	m02×04_02	2169	698
16	h01×04_01o4	1129	m04×04_04	1327	198
24	h01×04_01o6	831	m03×08_06	880	49
40	h02×10_05o2	635	m04×10_10	637	2
48	h03×08_06o2	604	m06×08_12	606	2

859

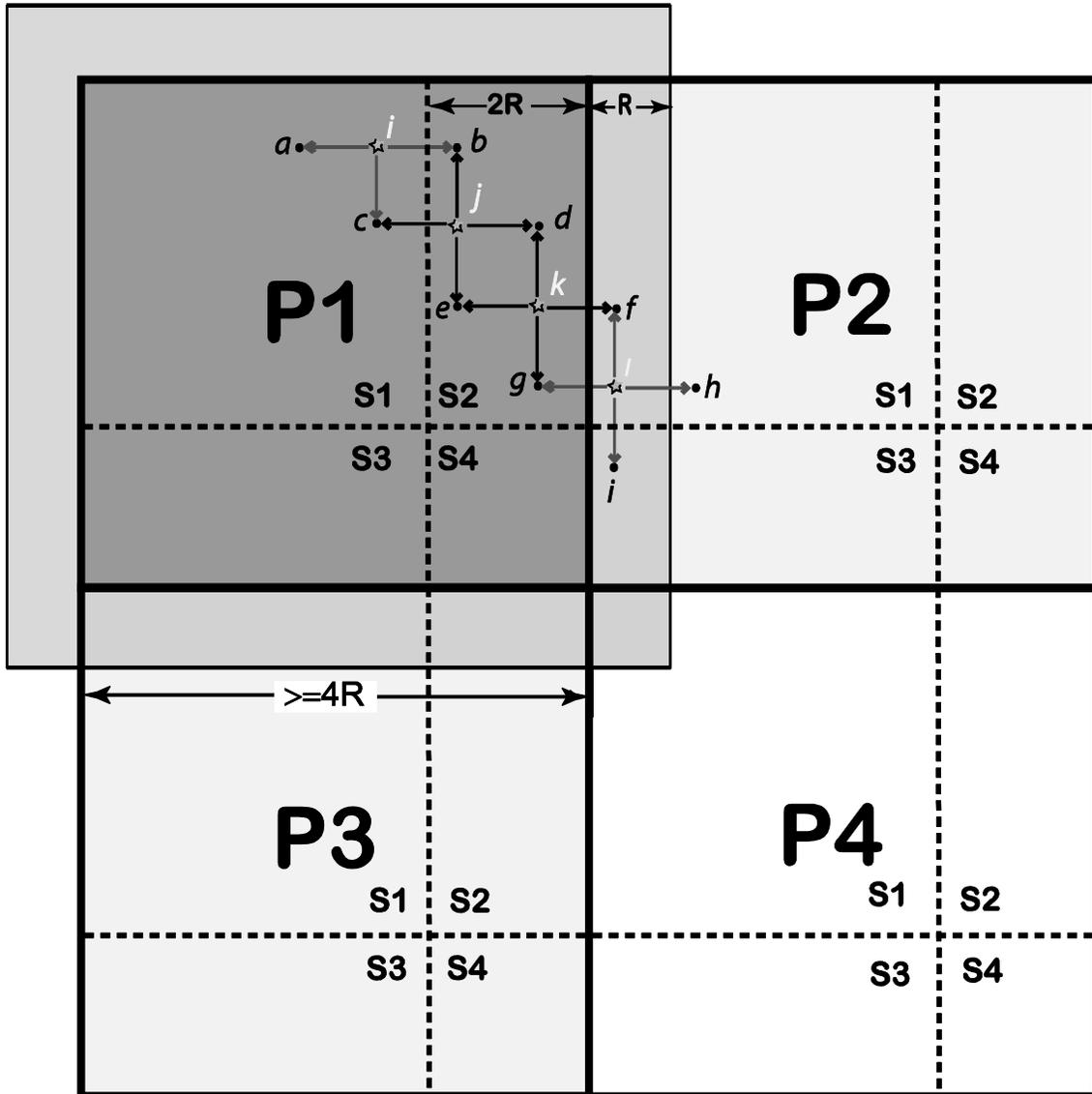
860

861
862

Table 4. The names and configurations of real data experiments.

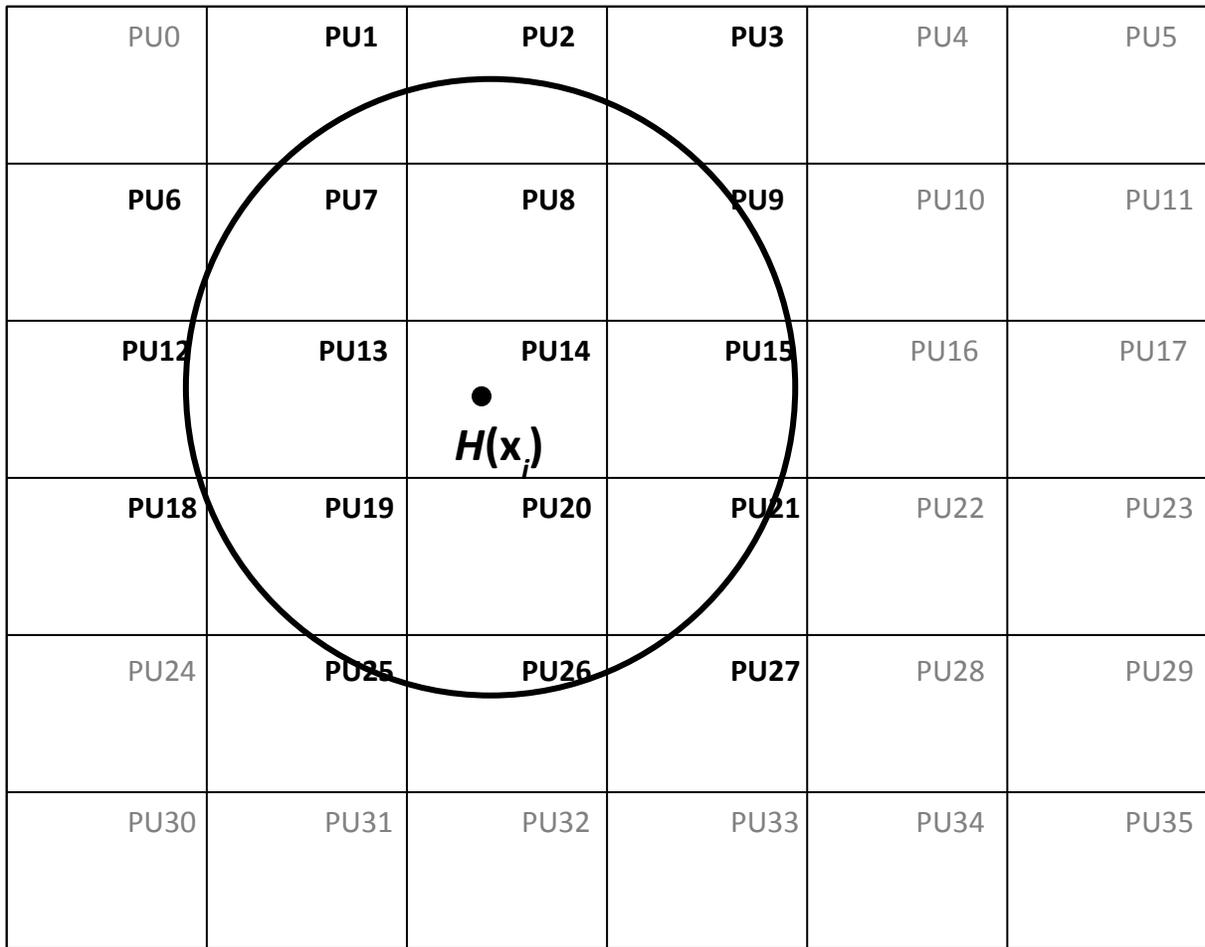
Experiment		Number of PUs in <i>x</i> direction	Number of PUs in <i>y</i> direction	Number of threads per PU	Total number of Cores used
OpenMP	o16			16	16
	o32			32	32
	o64			64	64
	o80			80	80
	o160			160	160
MPI	m16	1	16		16
	m32	2	16		32
	m64	3	16		64
	m80	5	16		80
	m160	10	16		160
Hybrid Group 1	h4o4	1	4	4	16
	h8o4	1	8		32
	h16o4	2	8		64
	h20o4	2	10		80
	h40o4	4	10		160
Hybrid Group 2	h2o8	1	2	8	16
	h4o8	1	4		32
	h8o8	1	8		64
	h10o8	2	5		80
	h20o8	4	5		160
Hybrid Group 3	h2o16	1	2	16	32
	h4o16	1	4		64
	h5o16	1	5		80
	h10o16	2	5		160

863
864



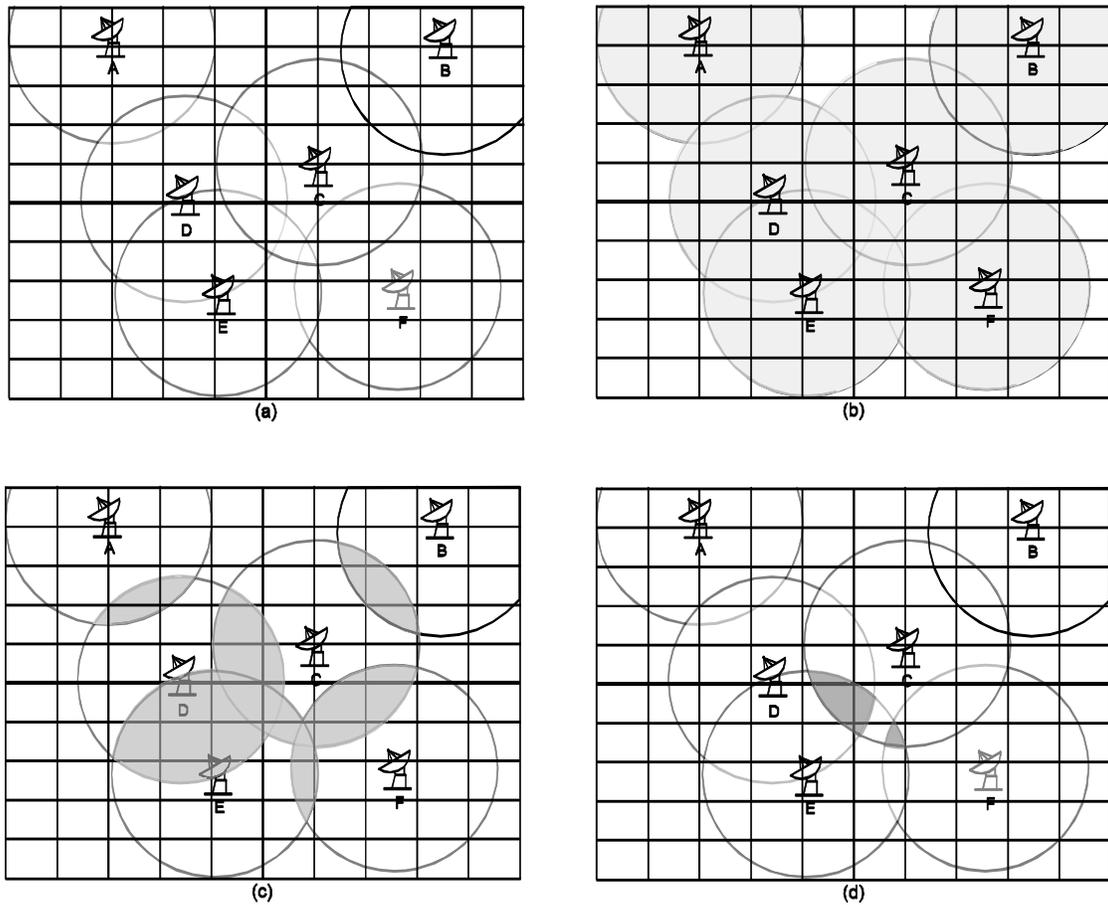
865

866 Fig. 1. A schematic of the domain decomposition strategy for the analysis of high-
 867 density observations, illustrated with 4 processing units (PUs, denoted by P1
 868 through P4). Letters *i-l* denote observations that are assumed to be equally spaced
 869 and letters *a-h* indicate the influence limits (as determined by the covariance
 870 localization radii of EnKF) of those observations. In this example, observations *i*
 871 and *l* are far enough apart that they will not influence any of the same state
 872 variables; they are among the observations that are analyzed simultaneously in the
 873 first step of the procedure. Observations *j* and *k* are analyzed in the second step,
 874 but they must be analyzed sequentially. Note that in practice, there will be many
 875 more observations within patches S1 and S2 of subdomains P1 to P4 than shown
 876 in the figure.



877

878 Fig. 2. A schematic for analyzing conventional data. Three steps are involved
 879 when analyzing one observation whose location is denoted by a black dot in the
 880 figure: 1) PU14 computes $H(x_i)$ (where i is the ensemble index); 2) $H(x_i)$ are
 881 broadcasted to all PUs; 3) state variables x_i within the influence range of this
 882 observation (within the large circle) are updated in parallel by the PUs that carry
 883 the state variables.

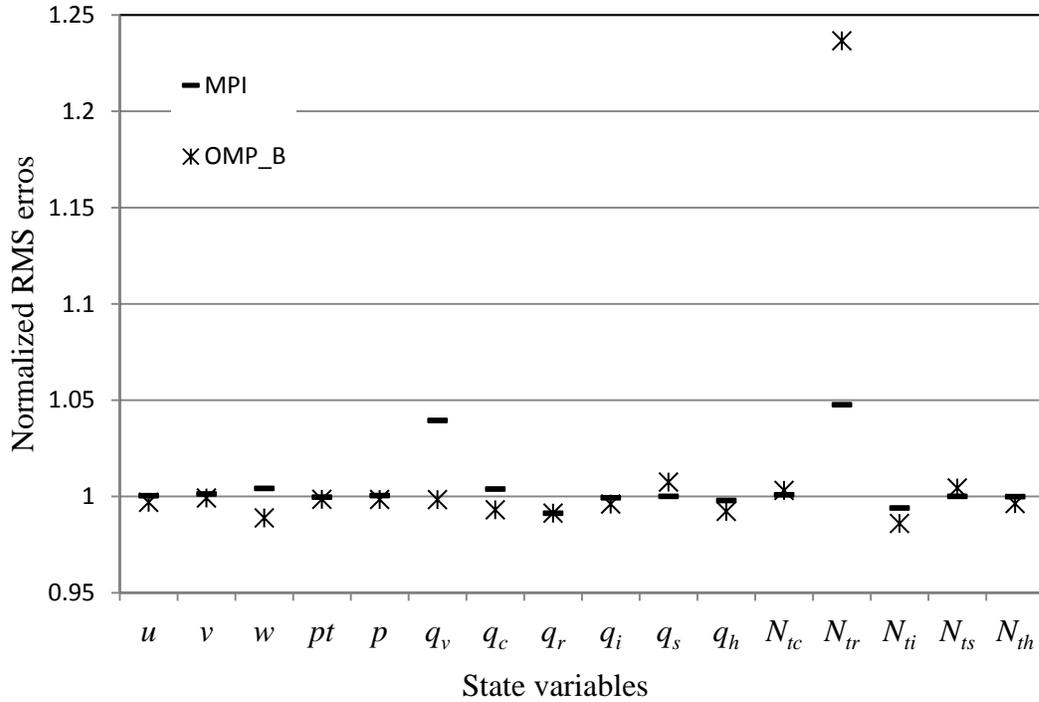


884

885 Fig. 3. Composite radar data batches organized such that within each batch, no
 886 more than one column of data exists for each grid column. (a) Observations from
 887 six radars (A-F) with their coverage indicated by the maximum range circles are
 888 remapped onto the model grid. (b) Observations of the first batch, (c) observations
 889 of the second batch, and (d) observations of the third batch. If there are more
 890 observations unaccounted for, additional data batch(es) will be formed.

891

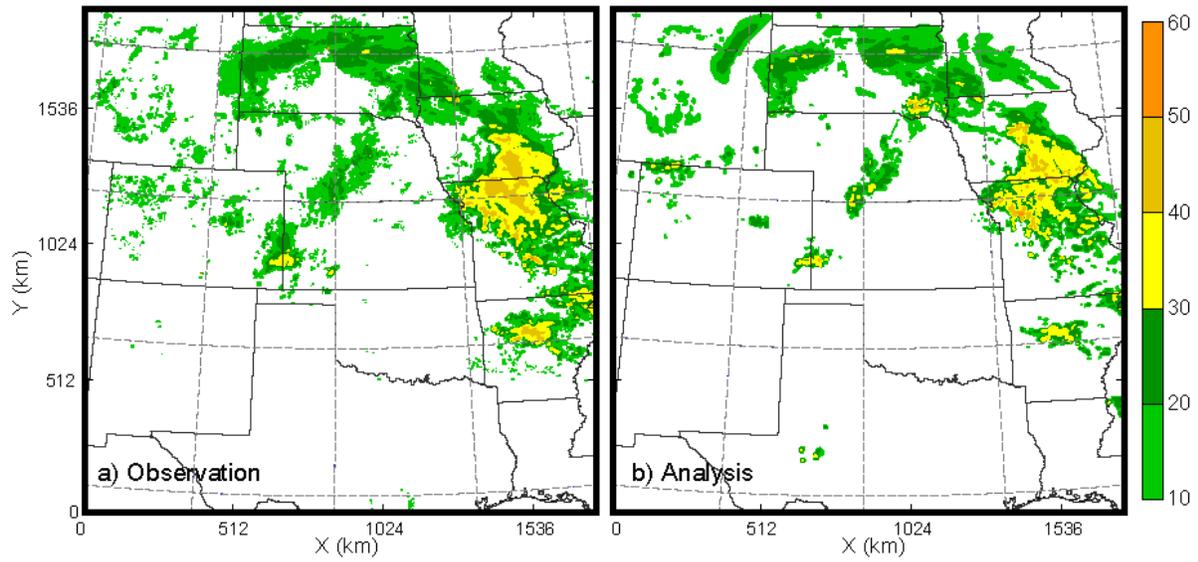
892



893

894 Fig. 4. RMS errors averaged over the grid points where truth reflectivity is greater
 895 than 10 dBZ and normalized by the errors of experiment OMP_F. The state
 896 variables are the 16 ARPS prognostic variables: three velocity components (u , v ,
 897 and w), potential temperature (pt), pressure (p), mixing ratios of water vapor (q_v),
 898 cloud water (q_c), rain water (q_r), cloud ice (q_i), snow aggregate (q_s), hail (q_h) and
 899 their respective number concentrations (N_{tc} , N_{tr} , N_{ti} , N_{ts} , and N_{th} , associated with a
 900 two-moment microphysics scheme used).
 901

902

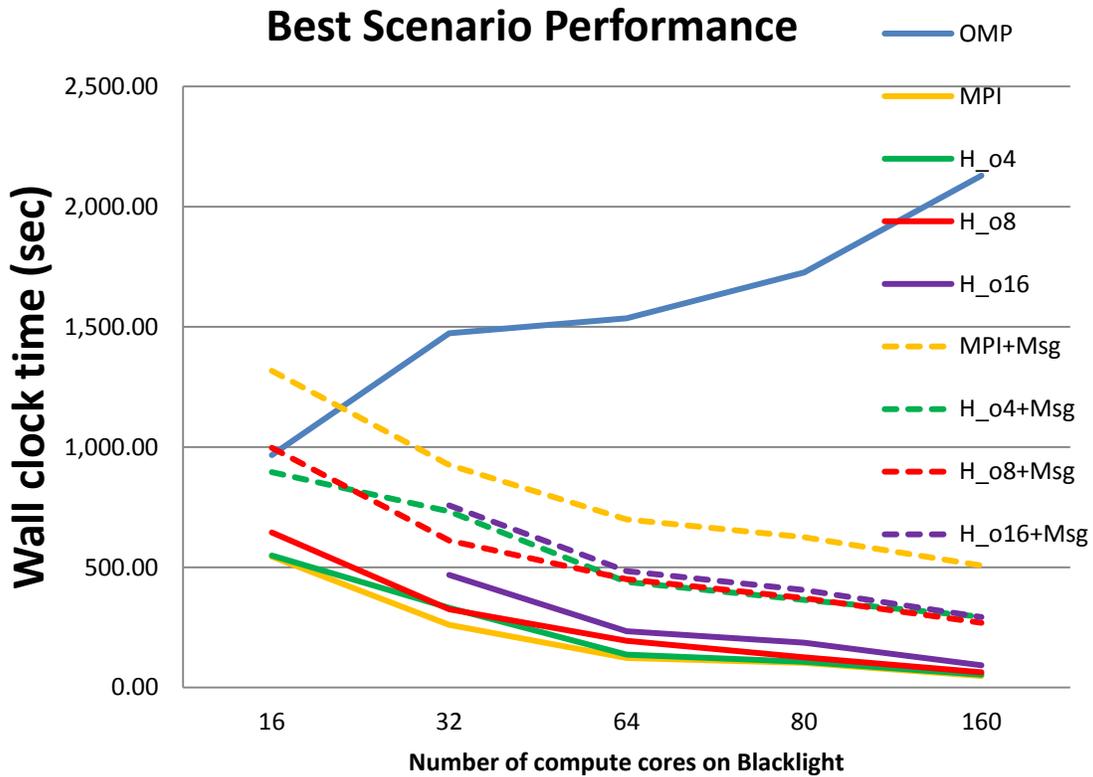


903

904 Fig. 5. (a) The observed radar reflectivity mosaic and (b) the reflectivity field
905 analyzed by the parallel EnKF algorithm, at model grid level 20 at 1800 UTC 10
906 May 2010.

907

908



913

914 Fig. 7. Wall clock times of the EnKF analyses as a function of the total number of
 915 compute cores used, for the 10 May 2010 real data case in the analysis domain
 916 shown in Fig. 6, obtained on the PSC Blacklight (an SGI UV 1000). OMP denotes
 917 pure OpenMP runs, MPI denotes pure MPI runs, and H_o4, H_o8, and H_o16
 918 denote hybrid runs with 4, 8, and 16 OpenMP threads within each MPI process,
 919 respectively. In all cases, all cores on the compute nodes were fully utilized, either
 920 by individual MPI processes or by OMP threads. Solid lines denote the total time
 921 excluding message passing, and the dashed lines show the total times including
 922 message passing. Data I/O times are excluded from all statistics.