

Document downloaded from:

<http://hdl.handle.net/10251/40393>

This paper must be cited as:

Xu, T.; Gómez-Hernández, JJ.; Li ., L.; Zhou ., H. (2013). Paralellized ensemble Kalman filter for hydraulic conductivity characterization. *Computers and Geosciences*. 52:42-49. doi:10.1016/j.cageo.2012.10.007.



The final publication is available at

<http://dx.doi.org/10.1016/j.cageo.2012.10.007>

Copyright Elsevier

Parallelized Ensemble Kalman Filter for Hydraulic Conductivity Characterization

Teng Xu^{a,*}, J. Jaime Gómez-Hernández^a, Liangping Li^a, Haiyan Zhou^a

^a*Group of Hydrogeology, Universitat Politècnica de València, Camino de Vera, s/n, 46022 Valencia, Spain*

Abstract

The ensemble Kalman filter (EnKF) is nowadays recognized as an excellent inverse method for hydraulic conductivity characterization using transient piezometric head data. Its implementation is well suited for a parallel computing environment. A parallel code has been designed that uses parallelization both in the forecast step and in the analysis step. In the forecast step, each member of the ensemble is sent to a different processor, while in the analysis step, the computations of the covariances are distributed between the different processors. An important aspect of the parallelization is to limit as much as possible the communication between the processors in order to maximize execution time reduction.

Four tests are carried out to evaluate the performance of the parallelization with different ensemble and model sizes. The results show the savings provided by the parallel EnKF, especially for a large number of ensemble realizations.

Keywords: Parallel EnKF; Cluster; Hydraulic conductivity; Parallel computing;

1. Motivation

The ensemble Kalman filter (EnKF) proposed by Evensen (Evensen, 2003) has proven to be an effective inverse method. It has been applied in many fields such as petroleum

*Corresponding author. Tel: +34 963879615 Fax: +34 963879492

Email addresses: tenxu@posgrado.upv.es (Teng Xu), jaime@dihma.upv.es (J. Jaime Gómez-Hernández), liali@upvnet.upv.es (Liangping Li), haizh@upvnet.upv.es (Haiyan Zhou)

4 engineering, oceanography, meteorology or hydrology (e.g., Aanonsen et al., 2009; Evensen,
5 2003; Dowell et al., 2010; Hendricks Franssen and Kinzelbach, 2008; Li et al., 2012b). While
6 it has proven to be more effective than alternative inverse methods (Hendricks Franssen and
7 Kinzelbach, 2009), it still has important computational needs.

8 The ensemble Kalman filter is based on the simultaneous analysis of a large number
9 of realizations. The first attempts to reduce CPU time usage were aimed at reducing the
10 number of realizations in the ensemble. The covariance localization is a modification of the
11 initial EnKF implementation that serves to reduce the ensemble size without compromising
12 much the quality of the outcome (e.g., Houtekamer and Mitchell, 1998; Anderson, 2007;
13 Devegowda et al., 2010; Zhang and Oliver, 2010). Other authors have proposed to reduce
14 the size of each realization, for instance, Li et al. (2012a) couple upscaling and the EnKF
15 for the inverse modeling of groundwater flow.

16 In any case, even for the most efficient implementation of the EnKF, the fact that it works
17 on each one of an ensemble of realizations makes it amenable to parallelization and thus take
18 advantage of multi-processor computers or of grid computing to reduce even further the time
19 needed for the algorithm to run. Although it will be reviewed in more detail later in the
20 paper, recall that the EnKF takes an ensemble of realizations, runs a forward model in each
21 realization, collects state data at observation locations and uses the difference between the
22 predicted and observed values to update each one of the realizations. The updating step
23 requires using all the ensemble realizations to compute the Kalman gain.

24 Keppenne (2000) proposed a parallel algorithm in which the forward model for each
25 ensemble member is run in a different processor, then, to compute the Kalman gain, a
26 domain decomposition is performed and each processor ends up with a small portion of all
27 ensemble realizations, finally, for the updating, each processor is responsible of the update of
28 a realization. (Keppenne and Rienecker, 2002a,b, 2003) also apply a domain decomposition
29 in the updating step and each processor is responsible of updating the subdomain used

30 for the Kalman gain computation in all realizations. This approach has the advantage of
31 avoiding the ensemble transpositions across processor that would be required if, after the
32 computation of the Kalman gain, each processor updates a full realization.

33 Recently, Tavakoli et al. (2011) have proposed a parallel EnKF algorithm applied to
34 petroleum engineering using a three-level parallelization. On the first level, each ensemble
35 member runs on a separate processor during the forecast step, on the second level uses a
36 parallel implementation of the reservoir simulator, and, on the third level, it handles the
37 matrix-vector multiplications involved in the Kalman gain computation and the updating of
38 the ensemble members.

39 In this paper we propose an alternative parallel EnKF algorithm and provide the com-
40 puter code to run it on a parallel environment using MODFLOW (McDonald and Harbaugh,
41 1988; Harbaugh et al., 2000) as the forward model. To the best of our knowledge this is the
42 first application of a parallel EnKF algorithm in the field of hydrogeology.

43 The paper proceeds with an overview of the EnKF, the strategy for parallelization, and
44 the evaluation of the algorithm on two examples.

45 **2. Overview of the EnKF**

46 The EnKF is the evolution of the Kalman filter to better address nonlinear state transfer
47 functions using a Monte-Carlo approach. The Kalman filter was proposed by Kalman et al.
48 (1960) as a data assimilation filter to improve the estimation of the state of a dynamic linear
49 system. Later, the extended Kalman filter (EKF) was proposed to address nonlinear systems;
50 the extension is based on a linearization of the nonlinear model, using a Taylor expansion,
51 for the computation of the time evolution of the covariance (McElhoe, 1966). The EKF has
52 been used in many fields, including hydrology (e.g., Hantush, 1997; Leng and Yeh, 2003;
53 Yeh and Huang, 2005), however, it has severe shortcomings in dealing with highly nonlinear
54 functions due to the accumulative error induced by the linearization process. Besides, the

55 algorithm itself is time and storage consuming, yielding it infeasible for large-scale system.
 56 To overcome these problems, the EnKF was proposed, the specifics of which are introduced
 57 next.

58 There are many uncertain factors in groundwater modeling: initial and boundary con-
 59 ditions, forcing terms, parameter values, ... (Hendricks Franssen and Kinzelbach, 2009).
 60 In this paper, we focus on the uncertainty about the parameter log-hydraulic conductivity.
 61 The state-transition equation is the standard three-dimensional groundwater flow equation,
 62 which is solved by MODFLOW, one of the most popular three-dimensional finite-difference
 63 groundwater flow simulators (e.g., McDonald and Harbaugh, 1988; Harbaugh et al., 2000).
 64 The filter has two steps. The first one is the prediction step, given by

$$s_i^f(t) = M(s_i^a(t - \Delta t) + w_i(t), w(t) \sim N(0, \sigma) \quad (1)$$

65 where $s_i^f(t)$ is the forecasted state of the system for a given set of parameters i , this state
 66 is function of the last estimate of the state of the system at $t - \Delta t$, $s_i^a(t - \Delta t)$ through a
 67 state-transition equation represented by M . Equation M is only an approximation of how
 68 the system behaves, therefore a model error w is added to the forecast, which is assumed to
 69 be Gaussian distributed, with zero mean and covariance σ .

70 The second step is the analysis step (Burgers et al., 1998), whereby the observed state
 71 measurements are used to update the forecast state to come up with a better estimate of
 72 the forecasted state.

$$s_i^a(t) = s_i^f(t) + k(t)[d_o(t) + e_i(t) - Hs_i^f(t)] \quad (2)$$

73

$$k(t) = p^f(t)H^T[Hp^f(t)H^T + R]^{-1} \quad (3)$$

74 where the forecasted state $s_i^f(t)$ is updated as a function of the difference between the
 75 predicted state $s_i^f(t)$ and the observations $(d_o(t) + e_i(t))$; matrix H is a measurement matrix

76 that serves to map the model predictions (generally over a fixed grid) onto the locations
77 where the observations are taken. Note that the observations are composed of two terms the
78 “real” state of the system $d_o(t)$ plus a measurement error $e_i(t)$ of mean zero and covariance
79 R . The amount by which the forecasted state should be modified is controlled by the
80 Kalman gain $k(t)$, which is a function of the state covariance $p^f(t)$. The ensemble Kalman
81 filter was developed precisely to avoid the shortcomings in computing the state covariance
82 by linearizing the state-transfer equation. For this purpose, an ensemble of realizations is
83 generated, and their state is forecasted and updated sequentially in time as data are collected.
84 At each time step, the covariance is numerically inferred from the ensemble of realizations,

$$p^f(t) = \frac{1}{N_e - 1} \sum_{i=1}^{N_e} (s^f(t) - \langle s^f(t) \rangle)(s^f(t) - \langle s^f(t) \rangle)^T \quad (4)$$

$$\langle s^f(t) \rangle = \frac{1}{N_e} \sum_{i=1}^{N_e} s^f(t) \quad (5)$$

85 where N_e is the number of realizations of the ensemble, $\langle \cdot \rangle$ denotes ensemble average, and
86 p is a matrix of dimensions $n \times n$, with n is the number of nodes at which the state of the
87 system is predicted by the numerical forecast model.

88 In inverse modeling applications, the state of the system is augmented so that it not
89 only includes the properly-speaking state of the system but also the parameters defining
90 the transfer function. In hydrogeology it is common to use such an augmented state, in
91 our case we use piezometric heads and log-hydraulic conductivities as state and parameters,
92 respectively. In our implementation we use an augmented vector, the state transfer function
93 in Eq. (1) leaves unchanged the logconductivity values, and updates the piezometric heads
94 according to the groundwater flow model; then, in the analysis step, we limit ourselves to
95 update the logconductivity values as explained next. For the sake of demonstration, we will
96 assume that observations are taken at groundwater model prediction locations; limiting the

97 component of the measurement matrix H to be 0's and 1's, and simplifying the expressions
 98 in Eqs. (2) and (3) as follows. We may partition the covariance matrix

$$p = \begin{bmatrix} p_{hh} & p_{Yh} \\ p_{hY} & p_{YY} \end{bmatrix} \quad (6)$$

99 where p_{hh} , p_{YY} are the covariance between hydraulic heads at measurement locations and
 100 the covariance between log-conductivities at all model gridnodes, respectively, and p_{Yh} and
 101 p_{hY} are the cross-covariances between a log-conductivity and a hydraulic head. And Eqs.
 102 (2) and (3) become

$$Y_i^a(t) = Y_i^f(t) + k(t)[d_o(t) + e_i(t) - h_i^f(t)] \quad (7)$$

103

$$k(t) = p_{Yh}[p_{hh} + R]^{-1} \quad (8)$$

104 where $Y_i^a(t)$ and $Y_i^f(t)$ are the elements of the augmented state vector corresponding to the
 105 logconductivities, and $h_i^f(t)$ is the predicted piezometric heads at measurement locations.

106 The output of the EnKF is an ensemble of realizations of hydraulic conductivity all
 107 of which are “conditioned” to the observation data. From this ensemble we can obtain
 108 average estimates and uncertainty estimates about the hydraulic conductivity, or we can
 109 post-process these fields to obtain optimal estimates, with their associated uncertainty, of
 110 response functions based on the ensemble, such as optimal pumping rates for the dewatering
 111 of a construction site. As more observational data are assimilated, the ensemble of hydraulic
 112 conductivities become more alike, and therefore the uncertainty associated is smaller, since
 113 there are less alternative realizations capable of reproducing the observed entire transient
 114 state information.

115 We end this introduction mentioned some of the main advantages and drawbacks of the
 116 EnKF. The main advantage of the EnKF is that is not an optimization approach but rather

117 a filtering approach. Therefore, there is no need for recursive evaluations of expensive cost
118 functions, just the need to compute the ensemble covariance and the mismatch between
119 predictions and observations, which will be used to update the conductivity fields. This
120 characteristic makes the EnKF very easy to implement and to apply in cases with many
121 different sources of information. The main drawback, leaving apart the Monte-Carlo aspect
122 of the method and the need to handle many realizations, is that the EnKF has been found
123 to collapse underestimating the final uncertainty (that is, as the fields keep been updated,
124 they tend to become more and more similar). There are two main reasons for this behavior,
125 one is the appearance of spurious correlations between distant points due to the numerical
126 nature of the covariance calculations, the other has to do with the number of realizations, if
127 it is not large enough, the empirical covariance based on a reduced number of realizations
128 tends to decrease as new updating steps are performed. The approaches to deal with these
129 problems are covariance localization and covariance inflation.

130 **3. Parallelization approach**

131 There are mainly three kinds of parallel computer architectures, the first kind is based
132 on shared memory, the second one is based on distributed memory, and the third one is a
133 hybrid in between both of them. The first type corresponds to multi-core computers, the
134 benefits are that the communication between the different ranks is very fast, and that it is
135 very easy to share data among them. However, this type of architecture is limited by the
136 total amount of memory available, and may be unable to address large models. The second
137 type corresponds to sharing resources among a grid of computers, the grid could easily be
138 enlarged attaching new computers to it, and therefore the limitation about the size of the
139 model they can handle disappears, on the contrary, the communication between the different
140 processors is much slower than in the shared-memory architecture. The best architecture
141 is the third one that uses a grid of multi-core computers, balancing the advantages and

142 disadvantages of the first and second architecture kinds.

143 In this paper, we employ a hybrid architecture. The cluster of computers, known as
144 “Pleiades”, consists of three HP Proliant DL 580, each with six-core four processors AMD
145 Opteron Model 8439 SE (six-core, 2.8GHz, 6MB L3, 105W), which amounts to a total of
146 24 cores per computer. All cores are 64-bit, and each computer has 256 GB of RAM.
147 The operating system is CentOS. Communication between processors is via message passing
148 interface (MPI). The cluster is networked via Ethernet TCP/IP.

149 As already explained previously, the EnKF consists of two steps: forecast and analysis
150 (or updating). We have analyzed the potential for parallelization of both steps.

151 Fig. 1 shows a flowchart of the proposed parallelization, which is explained next.

152 3.1. Forecast step

153 The most straightforward way to parallelize the forecast step in the EnKF is at the
154 realization level (Chen and Zhang, 2006). See box 1 of Fig. 1, let $s(n \cdot N_e)$ be the state
155 vector including all the ensemble members, where n represents state-vector size and N_e
156 is the number of realizations of the ensemble; if there are k processors, then the metric
157 can be decomposed into sub-metrics $s(n \cdot m_0), s(n \cdot m_1), \dots, s(n \cdot m_k)$, where m_0, m_1, \dots, m_k
158 ($m_0 + m_1 + \dots + m_k = N_e$) denotes the number of ensemble members that must be processed
159 by each processing element, PE_0, PE_1, \dots, PE_k respectively. Once the realizations are dis-
160 tributed among the processors, the MODFLOW simulator is called to generate the forecast
161 piezometric heads in all realizations. The distribution of the realizations to the processors
162 must be so that the load on all processors is as similar as possible. If N_e is a multiple of k ,
163 then each processor will receive N_e/k realizations, otherwise there will be $N_e \% k$ processors
164 receiving one extra realization.

165 Since the processors operate asynchronously, and the analysis step cannot be performed
166 until the forecast is performed in all realization, it is necessary to call the *MPI_Barrier*

167 command before starting the update step, to block each processor until all processors reach
 168 the barrier point in the code (Dewaraja et al., 2002).

169 3.2. Update step

170 After the forecast step, the state vector s^f , which contains all realizations, is distributed
 171 between the processors. Next, we have to compute the covariances p_{Yh} , and p_{hh} . This calcu-
 172 lation is distributed as follows, first accumulate the distributed state vector s^f in each proces-
 173 sor, $\sum s^{f0}$, $\sum s^{f1}$, ..., $\sum s^{fk}$, (recall that each processor is in charge of the forecast of a subset
 174 of all realizations) send these accumulated values to one of the processors PE_0 and compute
 175 the mean value of each component. Broadcast the mean values to all of the processors, and
 176 accumulate the products of the differences of the state vector with respect to their means in
 177 each one of the processors $\sum (s^{f0} - \langle s^f \rangle)(s^{f0} - \langle s^f \rangle)^T, \dots, \sum (s^{fk} - \langle s^f \rangle)(s^{fk} - \langle s^f \rangle)^T$.
 178 Then, send the accumulated products of differences to PE_0 and compute the covariances in
 179 this processor. Notice that the accumulated products of differences must be computed only
 180 for the individual entries in p_{Yh} , and p_{hh} , not for all possible entries of p^f in Eq. (6). With
 181 the covariances computed, the Kalman gain is obtained and broadcasted back to all proces-
 182 sors so that the updating Eq. (7) is applied to each member of the state vector distributed
 183 between the processors.

184 4. Application

185 The 3-D transient groundwater flow equation is:

$$S_s \frac{\partial h}{\partial t} - \nabla \cdot (k \nabla h) = W \quad (9)$$

186 where S_s is specific storage coefficient [L^{-1}]; h is the hydraulic head [L]; k is the hydraulic
 187 conductivity [LT^{-1}]; t is the time [T]; W denotes sinks and sources per unit volume [T^{-1}].

188 We use four test cases corresponding to two different model sizes. Cases 1, 2 and 3 use
189 a synthetic model with 50 by 50 by 5 cells, and the difference between them is the number
190 of realizations in the ensemble, which are 1200, 720 and 240. Case 4 uses a synthetic model
191 with 50 by 50 by 1 cell and 1200 realizations in the ensemble. All cells are 5 m by 5m by 2
192 m.

193 For the small model of size 50 by 50 by 1 cells, there are 75 observation wells in the
194 domain located as shown in Fig. 2. The left boundary has a specified head boundary equal
195 to 8 m, the right boundary has a specified flow of -0.0008 d^{-1} and the top and bottom
196 boundaries are no flow.

197 For the large model of size of 50 by 50 by 5 cells, there are 25 observation wells in the
198 domain, as shown in Fig. 3, that monitor the piezometric heads at the first, third and fifth
199 layers, for a total of 75 measurements. A verification well, located at row 30, column 20 and
200 layer 3 (see Fig. 3) is used to test the evolution of the piezometric head at an unsampled
201 location. The first layer of the left boundary is given a specified head of 8 m, the fifth layer
202 of the right boundary is given a prescribed flow of -0.008 d^{-1} . The rest of the boundary are
203 no flow boundaries.

204 For both models, the initial head is set 8 m throughout the domain. Specific storage S_s
205 is set as 0.0008 m^{-1} . The total simulation time is 500 days, discretized into 100 time steps
206 according to the following progression

$$\Delta t_k = \delta \Delta t_{k-1} \quad (10)$$

207 where δ is 1.05.

208 Log-conductivity $\ln(k)$ is assumed to be second-order stationary following a multi-Gaussian
209 distribution of zero $\ln \text{ m/d}$ mean, standard deviation $\sigma=1.5 \ln \text{ m/d}$, and an exponential co-
210 variance function.

$$r(h) = \sigma^2 [1 - \exp(-\frac{|h_x|}{\lambda_x} - \frac{|h_y|}{\lambda_y} - \frac{|h_z|}{\lambda_z})] \quad (11)$$

211 where the integral scales in the x , y , z directions are respectively $\lambda_x = 90$ m, $\lambda_y = 30$ m, $\lambda_z = 5$
 212 m.

213 The sequential Gaussian simulation module of the GSLIB software (Deutsch and Journel,
 214 1998) is used to generate the log-conductivity realizations. One of the realizations is chosen
 215 as the reference field.

216 4.1. Analysis

217 Fig. 4 shows the reference log-conductivity field, and Fig. 5a,5b,5c,5d,5e,5f shows the
 218 ensemble mean field for the tests 1, 2 and 3, at the beginning of the simulation (when no
 219 piezometric head data has been assimilated yet) and at the 50th time step. While comparing
 220 Fig.5b,5d,5f and Fig. 4, we find that the main features of the reference field are captured by
 221 the EnKF after the 50th data assimilation step. The larger the ensemble size, the smoother
 222 is the ensemble mean.

223 Fig. 6a, 6c and 6e shows the piezometric heads in the control well (not used for con-
 224 ditioning) computed on the initial ensemble of logconductivity realizations. They display a
 225 very large variability indicative of large prediction uncertainty. This variability is reduced
 226 when the piezometric heads are computed on the updated realizations, as seen in Fig. 6b,
 227 6d and 6f.

228 We can evaluate the goodness of the estimated field using the average deviation between
 229 the average of the ensemble members and the reference field using the root mean square
 230 error (*RMSE*).

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (s^{ref} - \langle s^a \rangle)^2} \quad (12)$$

231 Where N is the number of model gridblocks; s^{ref} is the value of the reference field; $\langle s^a \rangle$
 232 denotes the mean estimation of the ensemble fields.

233 Fig. 7 shows that the *RMSE* of test 3 has some fluctuations at early assimilation steps,
 234 yet after the 31st assimilation step, it begins to increase, becoming even larger than the
 235 starting value, which implies that for a small ensemble size the estimation of the covariance
 236 is poor (Hendricks Franssen and Kinzelbach, 2008). In the comparison of the *RMSE* between
 237 tests 1 and 2, the *RMSE* of test 1 is lower, plus, it is gradually smoother after the 27th
 238 assimilation step. So it can be concluded that the larger the ensemble, the better the
 239 estimation.

240 Similar results and conclusions can be reached for fourth test case, for which the model
 241 size is smaller while retaining the same number of measurements.

242 4.2. Parallelization analysis

243 Speedup and efficiency are usually used to evaluate the performance of parallel algorithms

$$\text{Speedup: } S_P = \frac{T_s}{T_P} \quad (13)$$

244

$$\text{Efficiency: } E_P = \frac{S_P}{P} \quad (14)$$

245 P is the number of the processors, T_s is the CPU time consumed under a serial implementa-
 246 tion of the algorithm, T_P is the CPU time consumed under a parallel implementation with
 247 P processors.

248 There is an obvious trade-off between the ensemble size and the CPU time, and also
 249 between the size of the numerical model and the CPU time. Table 1 and Fig. 8 show the
 250 performance of the parallel algorithm for test cases 1 (three-dimensional domain and 1200
 251 realizations), 3 (three-dimensional domain and 240 realizations) and 4 (two-dimensional
 252 domain and 1200 realizations).

253 From the table and figure we can see that the speedup is far from the ideal line for which
254 the improvement in CPU should follow the same proportion as the number of processors.
255 This is not a surprising result, there are two causes for this departure from the ideal per-
256 formance: there is a need for all processors to wait until they have finished certain tasks
257 for them to proceed on to the next task, and there is extra time needed for communication
258 between processors.

259 We notice some differences between the tests. Comparing tests 1 and 3, we notice that
260 the speedup is better for the case with the larger number of realizations, this is because for
261 test 3, each processor receives a smaller number of realizations and thus, proportionally, the
262 time spent in communication is larger for test 3 than for test 1. Apparently, for both tests,
263 the speedup could keep increasing if more processor were available. For test 4, however, the
264 speedup is similar to that of test 1 up to eight processors, then it appears to worsen, again,
265 the cause is found in the increase in the communication time among processors. It does not
266 seem that increasing the number of processors past 8 will increase the speedup (although
267 the overall CPU time will be still reduced) for test 4.

268 The efficiency curves have a similar behavior as the speedup curves. Efficiency worsens
269 past 8 processors for test 4, and for tests 1 and 3 we can conclude that the parallel algorithm
270 is more efficient the larger the ensemble size.

271 We can conclude that the parallel implementation of the EnKF runs with higher efficiency
272 for large size models and large ensembles than for small ones. In all cases, the final CPU
273 time is smaller than for the serial implementation, although the speedup is still far from
274 ideal.

275 In this parallel algorithm there are two inherent barriers to its processing capabilities, one
276 is data asynchrony, and the other is data-dependency. Data asynchrony makes all processors
277 be as slow as the slowest one, since they have to wait for all processors to finish a certain
278 task before they can proceed to the next one. Data dependency refers to the fact that the

279 distribution of realizations to the processors does not necessarily leaves all processors with
280 the same number of realizations. Data asynchrony mainly affects the forecast step, while
281 data dependency affects all steps and it is less noticed for large number of realizations.

282 **5. Summary and suggestion**

283 A parallel algorithm for the forecast and analysis steps of the EnKF has been presented
284 that reduces significantly the time to run the EnKF for large ensemble sizes. The efficiency
285 remains over 0.40 when using up to 12 processors for the two tests using 1200 realizations.

286 There are many measures that could be taken to reduce the communication time and
287 increase the efficiency such as decreasing the communication traffic, boosting communication
288 granularity or using a high-speed internet protocol to reduce the information transfer delay.
289 In addition, we must attempt to keep the load as balanced as possible among the processors,
290 that is, all processors should work on the same (or very similar) number of realizations.

291 Furthermore, the algorithm could be improved using a parallelized version of the ground-
292 water flow simulator, such as one employing domain decomposition, especially for large size
293 models.

294 **Acknowledgements**

295 The first author acknowledges the financial support from China Scholarship Council
296 (CSC). Financial support to carry out this work was also received from the Spanish
297 Ministry of Science and Innovation through project CGL2011-23295, and from the
298 Universitat Politècnica de València through project PERFORA.

299 Aanonsen, S., Nævdal, G., Oliver, D., Reynolds, A., Vallès, B., 2009. The ensemble Kalman
300 filter in reservoir engineering—a review. *SPE Journal* 14, 393–412.

301 Anderson, J., 2007. Exploring the need for localization in ensemble data assimilation using
302 a hierarchical ensemble filter. *Physica D: Nonlinear Phenomena* 230, 99–111.

303 Burgers, G., van Leeuwen, P., Evensen, G., Instituut, K.N.M., 1998. Analysis scheme in the
304 ensemble Kalman filter. *Monthly weather review* 126, 1719–1724.

305 Chen, Y., Zhang, D., 2006. Data assimilation for transient flow in geologic formations via
306 ensemble Kalman filter. *Advances in Water Resources* 29, 1107–1122.

307 Deutsch, C., Journel, A., 1998. *Gslib: Geostatistical Software Library and user’s Guide*.
308 New York 369.

309 Devegowda, D., Arroyo-Negrete, E., Datta-Gupta, A., 2010. Flow relevant covariance local-
310 ization during dynamic data assimilation using enkf. *Advances in Water Resources* 33,
311 129–145.

312 Dewaraja, Y., Ljungberg, M., Majumdar, A., Bose, A., Koral, K., 2002. A parallel Monte
313 Carlo code for planar and spect imaging: implementation, verification and applications in
314 131i spect. *Computer Methods and Programs in biomedicine* 67, 115–124.

315 Dowell, D., Zhang, F., Wicker, L., Snyder, C., Crook, N., 2010. Wind and temperature
316 retrievals in the 17 May 1981 arcadia, oklahoma, supercell: Ensemble Kalman filter ex-
317 periments .

318 Evensen, G., 2003. The ensemble Kalman filter: Theoretical formulation and practical
319 implementation. *Ocean Dynamics* 53, 343–367.

320 Hantush, M., 1997. Estimation of spatially variable aquifer hydraulic properties using
321 Kalman filtering. *Journal of Hydraulic Engineering* 123, 1027.

322 Harbaugh, A., et al., 2000. MODFLOW-2000, the US Geological Survey modular ground-
323 water Model: User Guide to Modularization Concepts and the Ground-water flow process.
324 US Geological Survey.

325 Hendricks Franssen, H., Kinzelbach, W., 2008. Real-time groundwater flow modeling with the
326 ensemble Kalman filter: Joint estimation of states and parameters and the filter inbreeding
327 problem. *Water Resources Research* 44, W09408.

328 Hendricks Franssen, H.J., Kinzelbach, W., 2009. Ensemble Kalman filtering versus sequential
329 self-calibration for inverse modelling of dynamic groundwater flow systems. *Journal of*
330 *Hydrology* 365, 261–274.

331 Houtekamer, P., Mitchell, H., 1998. Data assimilation using an ensemble Kalman filter
332 technique. *Monthly Weather Review* 126, 796–811.

333 Kalman, R., et al., 1960. A new approach to linear filtering and prediction problems. *Journal*
334 *of Basic Engineering* 82, 35–45.

335 Keppenne, C., 2000. Data assimilation into a primitive-equation model with a parallel
336 ensemble Kalman filter. *MONTHLY WEATHER REVIEW-USA* 128, 1971–1981.

337 Keppenne, C., Rienecker, M., 2002a. Development and initial testing of a parallel ensemble
338 Kalman filter for the poseidon isopycnal ocean general circulation model. *Monthly Weather*
339 *Review* 130, 2951–2965.

340 Keppenne, C., Rienecker, M., 2002b. Initial testing of a massively parallel ensemble Kalman
341 filter with the poseidon isopycnal ocean general circulation model. *Monthly Weather*
342 *Review(0027-0644)* 130, 2951–2965.

343 Keppenne, C., Rienecker, M., 2003. Assimilation of temperature into an isopycnal ocean
344 general circulation model using a parallel ensemble Kalman filter. *Journal of marine*
345 *systems* 40, 363–380.

346 Leng, C., Yeh, H., 2003. Aquifer parameter identification using the extended Kalman filter.
347 *Water resources research* 39, 1062.

- 348 Li, L., Zhou, H., Franssen, H., Gómez-Hernández, J., 2012a. Modeling transient groundwater
349 flow by coupling ensemble Kalman filtering and upscaling. *Water Resources Research* 48,
350 W01537.
- 351 Li, L., Zhou, H., Gómez-Hernández, J., Hendricks Franssen, H., 2012b. Jointly mapping hy-
352 draulic conductivity and porosity by assimilating concentration data via ensemble Kalman
353 filter. *Journal of Hydrology* 428-429, 152–169.
- 354 McDonald, M., Harbaugh, A., 1988. A modular three-dimensional finite-difference ground-
355 water flow model, techniques of water-resources investigations, book 6, chapter a1. US
356 Geological Survey 1, 988.
- 357 McElhoe, B.A., 1966. An assessment of the navigation and course corrections for a manned
358 flyby of Mars or Venus, in: *IEEE Transactions on Aerospace and Electronic Systems*,
359 Supplement, IEEE. pp. 613–623.
- 360 Tavakoli, R., Pencheva, G., Wheeler, M., 2011. Multi-level parallelization of ensemble
361 Kalman filter for reservoir history matching, in: *SPE Reservoir Simulation Symposium*.
- 362 Yeh, H., Huang, Y., 2005. Parameter estimation for leaky aquifers using the extended
363 Kalman filter, and considering model and data measurement uncertainties. *Journal of*
364 *Hydrology* 302, 28–45.
- 365 Zhang, Y., Oliver, D., 2010. Improving the ensemble estimate of the Kalman gain by boot-
366 strap sampling. *Mathematical Geosciences* 42, 327–345.

367 Figure 1: The parallel EnKF flow chart: box 1 is for the forecast step; box 2 is for the
368 updating step.

369 Figure 2: Flow domain and location of the 75 observation wells in the small model.

370 Figure 3: Flow domain and location of the 25 observation wells in the large model. The
371 filled triangles are observation wells measuring the piezometric heads at the first, third and
372 fifth layers, and the filled circle is a verification well.

373 Figure 4: Reference realization of $ln(k)$.

374 Figure 5: initial and updated ensemble means facies fields.(a), (b) show the ensemble
375 mean facies fields of the initial realizations and the 50th updated realizations of test 1; (c),
376 (d) show the ensemble mean facies fields of the initial realizations and the 50th updated
377 realizations for test 2; (e), (f) show the ensemble mean facies fields of the initial realizations
378 and the 50th updated realizations for test 3.

379 Figure 6: Piezometric head time evolution at the control well. The red curve is for the
380 reference field, the gray curves for each realization of the ensemble, and the green curve is
381 the ensemble mean. (a) and (b) show piezometric heads of the the initial realizations and
382 the 50th updated realizations for test 1; (c) and (d) show piezometric heads of the initial
383 realizations and the 50th updated realizations for test 2; (e) and (f) show piezometric heads
384 of the initial realizations and the 50th updated realizations for test 3. (For interpretation of
385 the references to color in this figure caption, the reader is referred to the web version of this
386 article.)

387 Figure 7: The root mean square error for three test.

388 Figure 8: Parallel speed-up and efficiency. (a) Speed up (b) efficiency.

389 Table 1: Parallel performance of three of the tests.

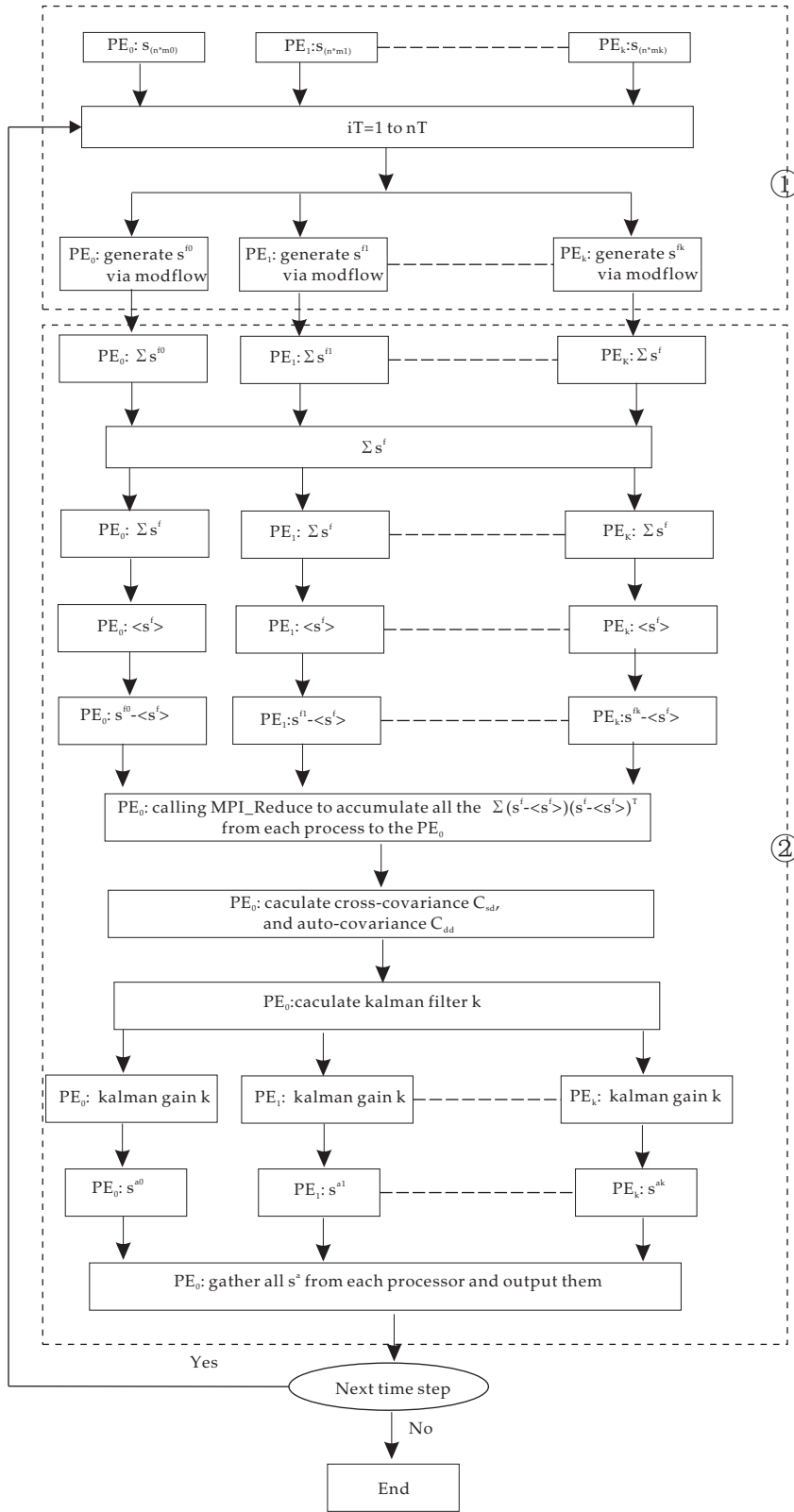


Figure 1
19

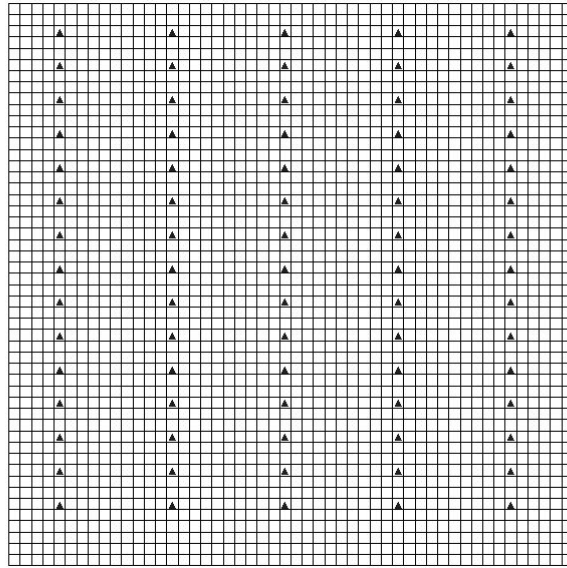


Figure 2

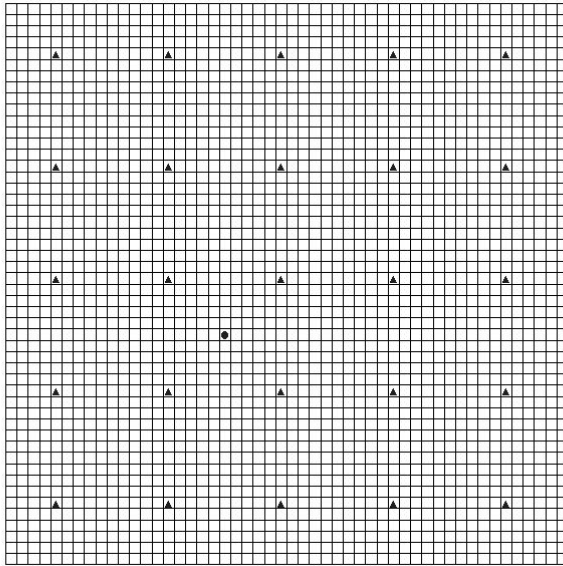


Figure 3

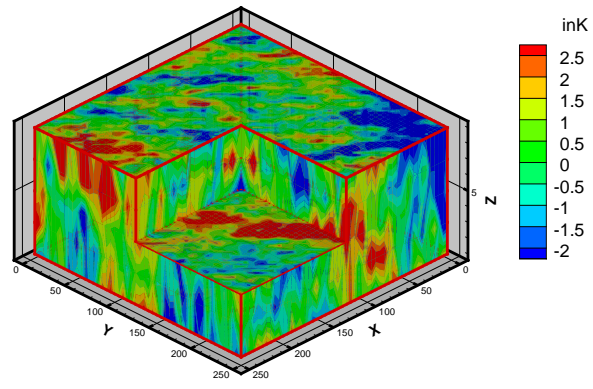


Figure 4

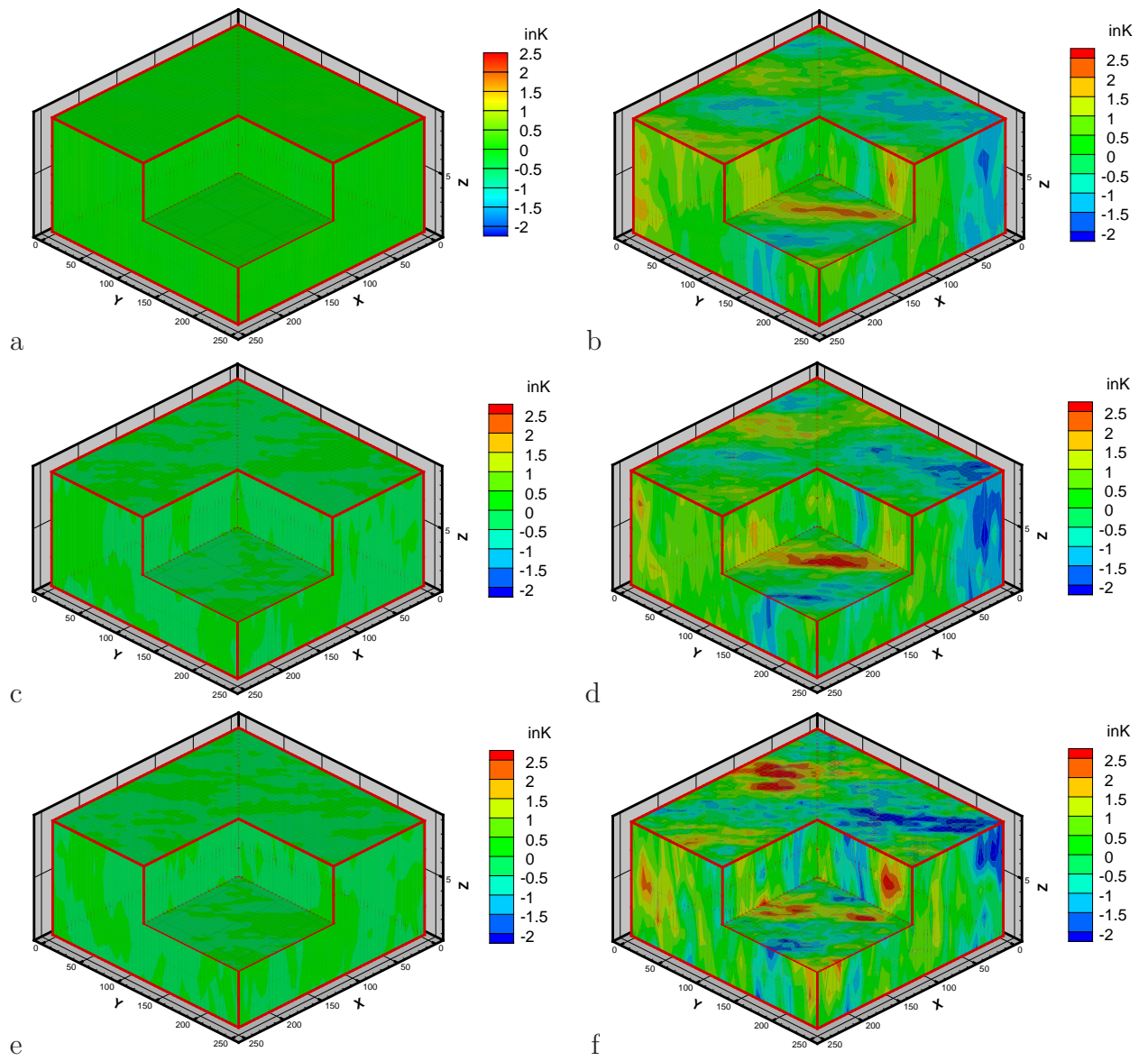


Figure 5

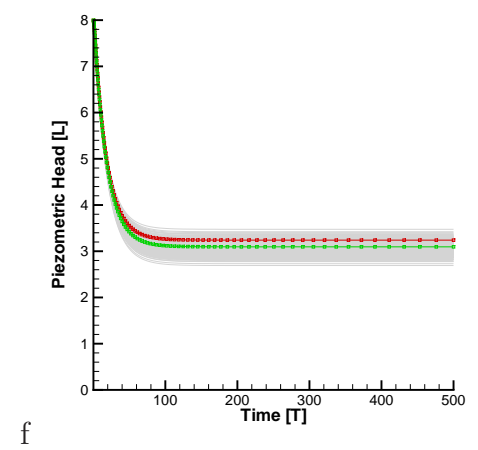
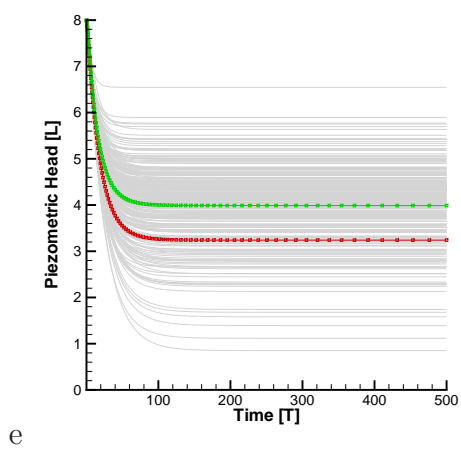
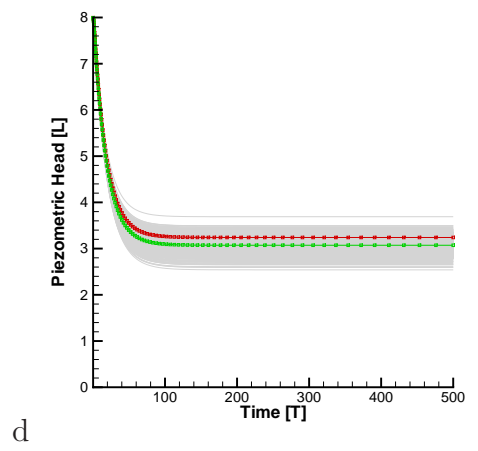
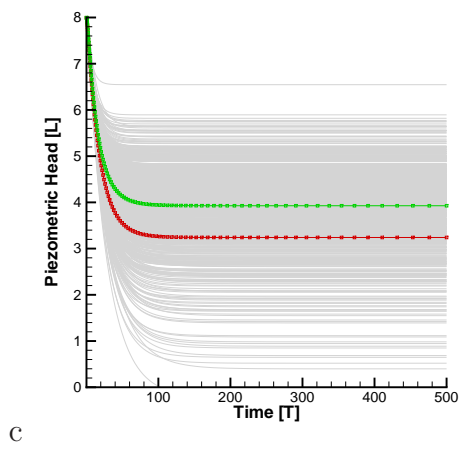
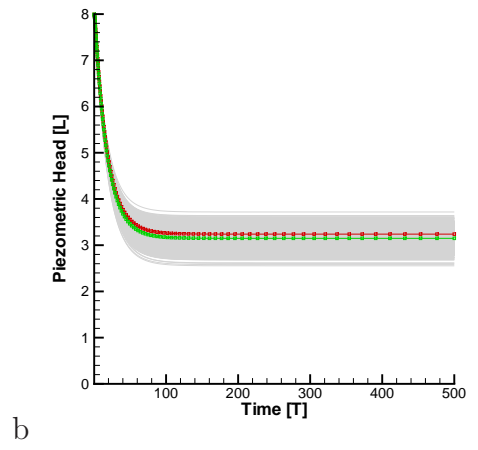
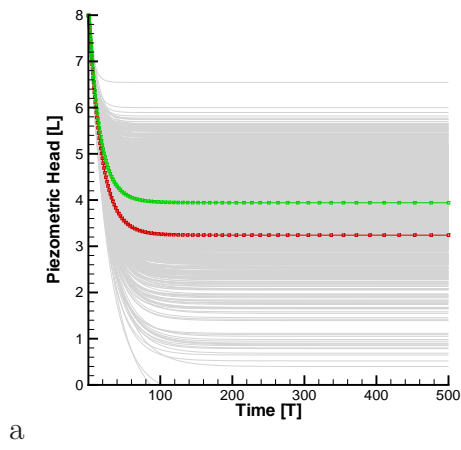


Figure 6

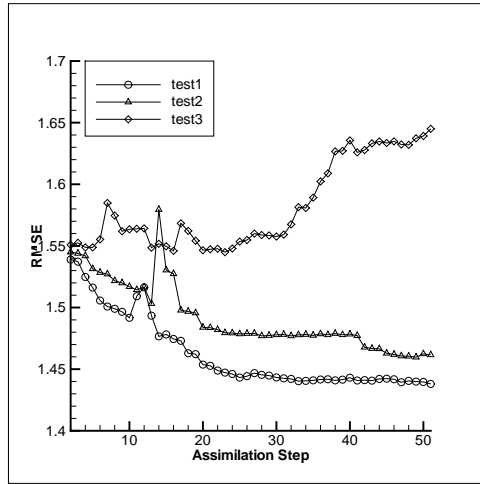


Figure 7

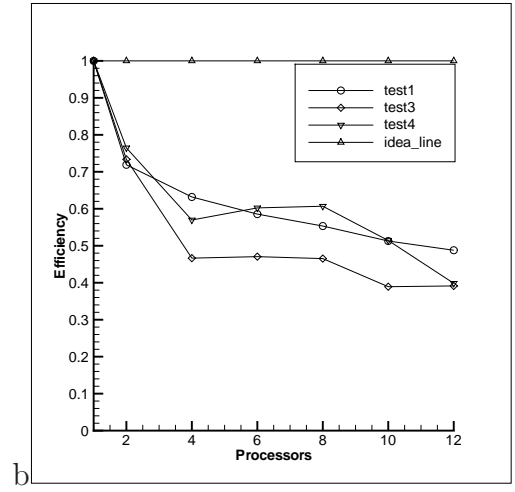
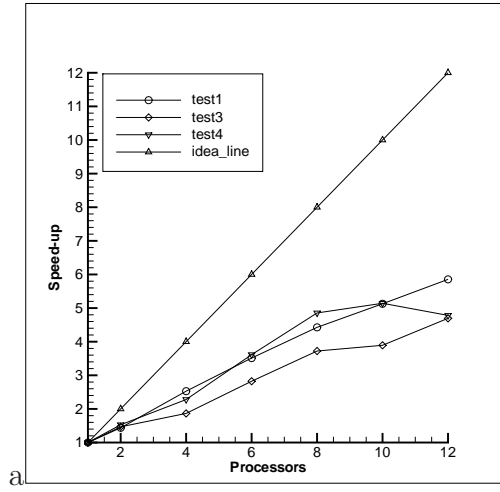


Figure 8

Table 1

Processors	Test1CPUs	Speedup	Efficiency	Test3CPUs	Speedup	Efficiency
1	20539.92	1.00	1.00	4866.50	1.00	1.00
2	14284.37	1.44	0.72	3316.35	1.47	0.73
4	8125.83	2.53	0.63	2607.83	1.87	0.47
6	5846.27	3.51	0.59	1723.75	2.82	0.47
8	4640.38	4.43	0.55	1307.34	3.72	0.47
10	4005.09	5.13	0.51	1250.23	3.89	0.39
12	3508.46	5.85	0.49	1036.38	4.70	0.39
Processors	Test4CPUs	Speedup	Efficiency			
1	10087.23	1.00	1.00			
2	6595.59	1.53	0.76			
4	4426.79	2.28	0.57			
6	2790.85	3.61	0.60			
8	2077.76	4.85	0.61			
10	1961.45	5.14	0.51			
12	2109.46	4.78	0.40			