

AN ADAPTIVE COVARIANCE PARAMETERIZATION TECHNIQUE FOR THE ENSEMBLE GAUSSIAN MIXTURE FILTER*

ANDREY A POPOV[†] AND RENATO ZANETTI[‡]

Abstract. The ensemble Gaussian mixture filter combines the simplicity and power of Gaussian mixture models with the provable convergence and power of particle filters. The quality of the ensemble Gaussian mixture filter heavily depends on the choice of covariance matrix in each Gaussian mixture. This work extends the ensemble Gaussian mixture filter to an adaptive choice of covariance based on the parameterized estimates of the sample covariance matrix. Through the use of the expectation maximization algorithm, optimal choices of the covariance matrix parameters are computed in an online fashion. Numerical experiments on the Lorenz '63 equations show that the proposed methodology converges to classical results known in particle filtering. Further numerical results with more advanced choices of covariance parameterization and the medium-size Lorenz '96 equations show that the proposed approach can perform significantly better than the standard EnGMF, and other classical data assimilation algorithms.

Key words. data assimilation, Gaussian mixture model, particle filtering, expectation maximization

MSC codes. 60G25, 62L12, 62M20, 93E11

1. Introduction. Sequential data assimilation [5, 37] aims to perform Bayesian inference on the state of some natural process from an inaccurate computational model and sparse and noisy observations. Traditionally, particle filter methods have been viewed as theoretically nice, but practically useless for inference of high dimensional systems. Recent advances in particle filters for high dimensions [43] have challenged this view.

Classical particle filters such as the bootstrap particle filter [37] make an empirical measure assumption on the prior distribution. Conversely filters such as the ensemble Kalman filter [17, 10] make an assumption that the first two moments of the empirical distribution are the only ones relevant to performing the inference, similar to a Gaussian assumption on the distribution. Gaussian mixture models (GMM) can extend the idea of an empirical measure approximation of the prior to a larger set of possible prior distributions, that combines the best of both worlds: it is able to represent non-Gaussian distributions while still assuming that the only two moments that matter to each mixture mode are its mean and covariance. The Gaussian sum filter [41] takes advantage of nice properties of GMMs, suffers from needing to propagate the covariance of each mode, and from requiring the need of many heuristics to ensure that the modes do not degenerate [36].

The ensemble Gaussian mixture filter (EnGMF) [4, 25, 44] and its related cousin the adaptive Gaussian mixture filter (AGMF) [42, 43] are sequential data assimilation algorithms that make use of the GMM approximation to the prior through the use of kernel density estimation (KDE) techniques. The EnGMF is based on the observation that Gaussian mixture models, under linear observation assumptions, are closed under multiplication [3]. The quality of the inference produced by the EnGMF is directly

*Submitted to the editors of SIAM SISC January 13, 2024.

Funding: This work was sponsored in part by DARPA (Defense Advanced Research Projects Agency) under STTR contract number W31P4Q-21-C-0032.

[†]Oden Institute for Computational Engineering and Sciences, University of Texas at Austin (andrey.a.popov@utexas.edu).

[‡]Department of Aerospace Engineering & Engineering Mechanics, University of Texas at Austin (renato@utexas.edu).

42 related to the accuracy of the GMM assumption about the prior distribution. This
 43 prior distribution is typically determined in whole by Monte Carlo samples through
 44 choices of the means and covariances of the GMM. While the choice of means is
 45 readily apparent as that of the Monte Carlo samples, the choice of covariance is less
 46 straightforward, and is the subject of much of the research surrounding KDE [40, 9].

47 The key innovations of this work are as follows: We first provide theoretical
 48 results that show the convergence of the EnGMF for a certain class of probabilistic
 49 assumption on the bandwidth parameter in the classical EnGMF algorithm. We next
 50 generalize the EnGMF by introducing the parameterization of statistical covariance
 51 matrix estimates from other ensemble-based filters to the EnGMF. We finally show
 52 how the EnGMF machinery could be used to choose the value of these parameters
 53 in an optimal adaptive fashion by utilizing the expectation maximization algorithm.
 54 Thus the sum total of these results is the adaptive Gaussian mixture filter (AEnGMF)
 55 which utilizes all this machinery for inference.

56 This work is organized as follows: we first introduce the data assimilation prob-
 57 lem and the EnGMF in section 2. We next present the adaptive ensemble Gaussian
 58 mixture filter in section 3 along with the expectation maximization algorithm in sub-
 59 section 3.1. Numerical experiments are provided in section 4, and concluding remarks
 60 in section 5.

61 **2. Background.** Assume that we are given a model that evolves a natural
 62 process of interest from time index $i - 1$ to time index i ,

63 (2.1)
$$x_i^t = \mathcal{M}(x_{i-1}^t) + \xi_i,$$

64 with model error ξ_i . For the remainder of this paper, we assume that the model error
 65 ξ_i is always zero, and thus the model (2.1) is exact.

66 The goal is to estimate the true state x^t of said process given some non-linear
 67 observation,

68 (2.2)
$$y_i = \mathcal{H}(x_i^t) + \eta_i,$$

69 with observation operator \mathcal{H} and an additive error term η_i . Denote with Y_i all the
 70 observations up to and including time index i ,

71 (2.3)
$$Y_i = \{y_1, y_2, \dots, y_i\}.$$

72 Given a prior at time index i , namely $x_i^b = x_i^t | Y_{i-1}$, we aim to perform Bayesian
 73 inference on these two sources of information,

74 (2.4)
$$p(x_i^b | y_i) \propto p(x_i^b) p(y_i | x_i^b)$$

75 resulting in the ‘analysis’, $x_i^a = x_i^b | y_i = x_i^t | Y_i$.

76 *Remark 2.1* (Model error). All of the derivations and algorithms presented in this
 77 work do not require the model error in (2.1) to be zero. This assumption is merely
 78 made for convenience in this work.

79 We next describe how a solution to (2.4) can be achieved using Monte Carlo
 80 sampling and the EnGMF.

81 **2.1. The Ensemble Gaussian Mixture Filter.** Assume that we have a col-
 82 lection of N particles at time index i that is represented as $\mathbf{X}_i^b = [x_{i,1}^b, x_{i,2}^b, \dots, x_{i,N}^b]$
 83 and is composed of weighted samples from the prior distribution $p(x_i^b)$ with weights

84 $\{u_{i,j}\}_{j=1}^N$. Given non-linear observations of the truth (2.2), our aim is to find a col-
 85 lection of N particles $\mathbf{X}_i^a = [x_{i,1}^a, x_{i,2}^a, \dots, x_{i,N}^a]$ containing samples from the posterior
 86 distribution, such that the posterior is the prior conditioned by the observations,

$$87 \quad (2.5) \quad x_i^a = x_i^b | y_i,$$

88 solving the Bayesian inference problem (2.4).

89 As the prior distribution of the particles is unknown, our prior knowledge can be
 90 used to construct an approximation thereto. From kernel density estimation theory,
 91 the ensemble Gaussian mixture filter (EnGMF) assumes that the distribution of the
 92 prior state at time index i , x_i^b , is given by the Gaussian mixture,

$$93 \quad (2.6) \quad x_i^b \sim \sum_{j=1}^N u_{i,j} \mathcal{N}(\bar{x}_{i,j}^b, \mathbf{B}_{i,j}^b),$$

94 where each mean exactly corresponds to one of the particles in the ensemble,

$$95 \quad (2.7) \quad \bar{x}_{i,j}^b := x_{i,j}^b.$$

96 The observation distribution at time index i is given by the Gaussian mixture,

$$97 \quad (2.8) \quad y_i | x_i^t \sim \sum_{k=1}^M v_{i,k} \mathcal{N}(\bar{y}_{i,k}, \mathbf{R}_{i,k}),$$

98 which is a generalization of the typical Gaussian assumptions on the observation error
 99 made in data assimilation literature.

100 The posterior distribution at time index i is defined [4] by the Gaussian mixture,

$$101 \quad (2.9) \quad x_i^a \sim \sum_{j=1}^N \sum_{k=1}^M w_{i,j,k} \mathcal{N}(\bar{x}_{i,j,k}^a, \mathbf{B}_{i,j,k}^a),$$

102 with the following set of definitions,

$$103 \quad (2.10) \quad \begin{aligned} \bar{x}_{i,j,k}^a &= \bar{x}_{i,j}^b - \mathbf{G}_{i,j,k} (\mathcal{H}(\bar{x}_{i,j}^b) - \bar{y}_{i,k}), \\ \mathbf{B}_{i,j,k}^a &= (\mathbf{I} - \mathbf{G}_{i,j,k} \mathbf{H}_{i,j}^T) \mathbf{B}_{i,j}^b, \\ \mathbf{G}_{i,j,k} &= \mathbf{B}_{i,j}^b \mathbf{H}_{i,j}^T (\mathbf{H}_{i,j} \mathbf{B}_{i,j}^b \mathbf{H}_{i,j}^T + \mathbf{R}_{i,k})^{-1}, \\ w_{i,j,k} &\propto u_{i,j} v_{i,k} \mathcal{N}(\bar{y}_{i,k} | \mathcal{H}(\bar{x}_{i,j}^b), \mathbf{H}_{i,j} \mathbf{B}_{i,j}^b \mathbf{H}_{i,j}^T + \mathbf{R}_{i,k}), \\ \mathbf{H}_{i,j} &= \left. \frac{d\mathcal{H}}{dx} \right|_{x=\bar{x}_{i,j}^b}. \end{aligned}$$

104 where $\bar{x}_{i,j,k}^a$, are the analysis Gaussian mixture means, $\mathbf{B}_{i,j,k}^a$ are the analysis Gauss-
 105 ian mixture covariances, $\mathbf{G}_{i,j,k}$ is similar to a gain matrix, $w_{i,j,k}$ are the Gaussian
 106 mixture weights, and $\mathbf{H}_{i,j}$ is the linearization of the observation operator around $\bar{x}_{i,j}^b$.
 107 When the observation operator is linear, $\mathcal{H}(x) = \mathbf{H}x$, the posterior GMM (2.9) is ex-
 108 actly the posterior corresponding to the assumed prior (2.6) and the observation (2.8)
 109 distributions.

110 Each Gaussian distribution in (2.6) has the following probability density function,

$$111 \quad (2.11) \quad \mathcal{N}(x | \bar{x}_{i,j}^b, \mathbf{B}_{i,j}^b) = \left| 2\pi \mathbf{B}_{i,j}^b \right|^{-\frac{1}{2}} e^{-\frac{1}{2} (x - \bar{x}_{i,j}^b)^T \mathbf{B}_{i,j}^{b,-1} (x - \bar{x}_{i,j}^b)},$$

112 with the other Gaussian distributions in (2.8)–(2.10) having a similar form.

113 *Remark 2.2* (Normalization Factors). Note that when either the observation
 114 operator \mathcal{H} is non-linear, or the kernel covariance matrices $\mathbf{B}_{i,j}$ are not identical, extra
 115 care must be taken when computing the weights $w_{i,j,k}$ in (2.10), as the covariances
 116 $\mathbf{H}_{i,j}\mathbf{B}_{i,j}^b\mathbf{H}_{i,j}^T + \mathbf{R}_{i,k}$ are not necessarily equal. This means that, the normalization
 117 factors in each Gaussian term (similar to (2.11)),

$$118 \quad (2.12) \quad \left| 2\pi \left(\mathbf{H}_{i,j}\mathbf{B}_{i,j}^b\mathbf{H}_{i,j}^T + \mathbf{R}_{i,k} \right) \right|^{-\frac{1}{2}},$$

119 are required to be computed. This can be performed in a computationally efficient
 120 manner through the use of the Cholesky decomposition and the log-sum-exp trick [7].

121 While the transformation of the distribution in (2.9) results in an estimate of the
 122 posterior distribution, the means $\bar{x}_{i,j}^a$ of this distribution are not actually samples
 123 from this distribution, thus it is not the case that the posterior samples are equivalent
 124 to these means,

$$125 \quad (2.13) \quad x_{i,j}^a \neq \bar{x}_{i,j}^a.$$

126 A resampling procedure is therefore required in order to obtain independently and
 127 identically distributed (iid) samples from (2.9). What follows is one such procedure.

128 *Procedure 2.3* (EnGMF resampling). Given the final posterior Gaussian mixture
 129 distribution in (2.9), it is possible to resample S samples from the posterior GMM
 130 through the following procedure:

- 131 1. for $s = 1, \dots, S$, sample the random variable ℓ from the discrete distribution
- 132 defined by the weights $\{w_{i,j,k}\}_{j,k}$,
- 133 2. sample $\mathbf{X}_{i,s}^a$ from the Gaussian $\mathcal{N}\left(x|\bar{x}_{i,\ell}^a, \mathbf{B}_{i,\ell}^a\right)$,
- 134 enabling samples to be generated from the posterior.

135 *Remark 2.4* (Arbitrary Sampling of the Posterior). Note, that using *Proce-*
 136 *dure 2.3* we are able to arbitrarily sample from the posterior distribution. This means
 137 that the number of posterior samples S could be significantly larger or significantly
 138 smaller than the original number of samples N used to generate said posterior.

139 *Remark 2.5* (Independent and Identically Distributed Samples). Note that while
 140 we make the convenient assumption that the samples generated by *Procedure 2.3* are
 141 iid, this is not actually the case. The parameters of each mode of the GMM are
 142 actually functions of the prior samples, and are themselves random variables, and
 143 thus introduce a dependence if two samples come from the same mode. This hidden
 144 dependence of the particles means that they are merely conditionally independent,
 145 making them exchangeable, but not independent in general.

146 *Remark 2.6* (Prior Uniform Weights). If *Procedure 2.3* is utilized to re-sample
 147 the particles at every step of the assimilation, then the prior distribution weights in
 148 (2.6) are all uniform $u_{i,j} = \frac{1}{N}$ under the assumption of a uniform transition density.

149 *Remark 2.7* (Differences Between the EnGMF and the AGMF). Unlike the En-
 150 GMF, in the AGMF (see [42, 43]), instead of resampling like in *Procedure 2.3*, the
 151 weights are scaled by a defensive factor towards uniformity,

$$152 \quad (2.14) \quad w_{i,j,k} = \alpha_i w_{i,j,k} + (1 - \alpha_i) \frac{1}{N},$$

153 with a ‘defensive factor’, α_i , which ensures that the weights do not degenerate. The
 154 new particles are taken to be the means of the candidate posterior distribution (2.9).
 155 This idea is not explored in this work.

156 The posterior GMM (2.9) is only a good representation of the exact posterior if
 157 the prior GMM assumption (2.6) is a good approximation of the distribution that
 158 originated the particles \mathbf{X}_i^b . For a finite ensemble N it is possible that the EnGMF
 159 analysis is a poor representation of the truth. The prior GMM assumption (2.6), in
 160 the particle limit $N \rightarrow \infty$, converges to the exact prior under certain assumption
 161 on the covariances, $\mathbf{B}_{i,j}$, which is a known result from kernel density estimation
 162 literature [40]. We now show that under some assumptions on the covariances, $\mathbf{B}_{i,j}$,
 163 the posterior GMM (2.9) produced by the EnGMF procedure also converges, in the
 164 particle limit, $N \rightarrow \infty$, to a distribution obtained from performing exact Bayesian
 165 inference.

166 **THEOREM 2.8 (EnGMF convergence).** *Assuming the observation distribution*
 167 *is exact (2.8), if the means $\bar{x}_{i,j}^b$ in the estimated prior distribution GMM (2.6) are*
 168 *samples from the underlying exact distribution with weights $u_{i,j}$, and the prior Kernel*
 169 *covariance matrices tend to zero in the limit of ensemble size, $\lim_{N \rightarrow \infty} \mathbf{B}_{i,j}^b = 0$, then*
 170 *the EnGMF with the resampling procedure Procedure 2.3 converges to a filter in the*
 171 *class of sequential importance resampling (SIR) filters.*

172 *Proof.* Given the assumptions above, in the limit of ensemble size, $N \rightarrow \infty$, the
 173 prior distribution GMM converges to the empirical distribution,

$$174 \quad (2.15) \quad p(x_i^b) = \frac{1}{N} \sum_{j=1}^N u_{i,j} \delta_{x_i^b - \bar{x}_{i,j}^b},$$

175 which converges weakly to the underlying prior distribution. Then as the prior Kernel
 176 covariance tends towards zero, the posterior GMM estimate defined by (2.9) and (2.10)
 177 converges to the empirical measure,

$$178 \quad (2.16) \quad p(x_i^a) = \frac{1}{N} \sum_{j=1}^N \left(\sum_{k=1}^M w_{i,j,k} \right) \delta_{x_i^a - \bar{x}_{i,j}^b},$$

179 which converges weakly to the exact posterior distribution. Then, the EnGMF re-
 180 sampling in Procedure 2.3 makes the EnGMF converge to an SIR filter. \square

181 **2.2. EnGMF rate of convergence for scalar parameterization.** We now
 182 motivate the importance of choosing a good parameterization of the covariance matrix
 183 is more cost effective than simply increasing the number of particles N . In the follow-
 184 ing, we show that the rate of convergence of the prior GMM estimate of the EnGMF
 185 is sub-linear under a scalar parameterization. This makes the choice of parameter the
 186 dominant factor that determines the goodness-of-fit of the distribution.

187 The most common parameterization of the prior GMM is a scalar parameterization
 188 that modifies the scaling of the covariance in a way that is guaranteed to degenerate
 189 in the limit of ensemble size. In the case of this scalar bandwidth parameterization,
 190 the prior covariance estimates become

$$191 \quad (2.17) \quad \mathbf{B}_{i,j}^b(\beta_{i,N}^2) = \beta_{i,N}^2 \mathbf{P}_{i,N}^b, \quad 1 \leq j \leq N,$$

192 where \mathbf{P}_i^b is the known (or approximated) covariance of the prior distribution, and β_i^2
 193 is a scaling factor yet to be determined. Following the derivations in [40], we provide
 194 formulations of the error and optimal density for the covariance parameterization
 195 in (2.17).

196 Given the unknown true prior distribution $p_{x_i^b}$, and its GMM approximation $\tilde{p}_{x_i^b}$,
 197 the most common metric for determining the convergence of the latter to the former
 198 is the mean integral squared error (MISE) at time index i ,

$$199 \quad (2.18) \quad \text{MISE}_i(p_{x_i^b}, \tilde{p}_{x_i^b}) = \mathbb{E}_{\mathbf{X}_i^b} \left[\int_{\Omega_x} \left(p_{x_i^b}(x) - \tilde{p}_{x_i^b}(x) \right)^2 dx \right],$$

200 Observe that the convergence of the prior GMM estimate is not the same as
 201 convergence of the posterior GMM estimate. Nevertheless convergence of the prior
 202 estimate implies convergence of the posterior estimate. It is therefore the case that the
 203 rate of convergence of the prior estimate is directly related to the rate of convergence
 204 of the prior estimate.

205 The MISE is typically approximated using the dominant terms of its expansion
 206 into the approximated mean integrated square error (AMISE) given by,

$$207 \quad (2.19) \quad \text{AMISE}(p_{x_i^b}, \tilde{p}_{x_i^b}) = \frac{1}{4} \beta_i^4 \alpha^2 \gamma_i + N^{-1} \beta_i^{-n} \delta,$$

208 where for the GMM approximation of a distribution α and δ ,

$$209 \quad (2.20) \quad \alpha = 1, \quad \delta = (2\sqrt{\pi})^{-n},$$

210 are known constants that depend on the dimension n , and γ_i ,

$$211 \quad (2.21) \quad \gamma_i = \int_{\Omega_x} \text{tr}^2 \left(\nabla_x^2 p_{x_i^b}(x) \right) dx,$$

212 is dependent on the true prior distribution at time index i . As the true prior is
 213 unknown, a known reference distribution ϕ can be used to compute an approximation
 214 for (2.21),

$$215 \quad (2.22) \quad \tilde{\gamma}_i = \int_{\Omega_x} \text{tr}^2 \left(\nabla_x^2 \phi(x) \right) dx,$$

216 where ϕ is often taken to be the standard Gaussian distribution.

217 If our assumed estimate of the parameter $\tilde{\gamma}_i$ is correct, the optimal bandwidth
 218 that minimizes the AMISE in (2.19) is,

$$219 \quad (2.23) \quad \beta_i^2 = \left(\frac{\delta n}{\alpha^2 \tilde{\gamma}_i N} \right)^{\frac{2}{n+4}},$$

220 by satisfying the first order optimality conditions of (2.19).

221 Plugging (2.23) back into (2.19), the error can now be written as,

$$222 \quad (2.24) \quad \text{AMISE}(p_{x_i^b}, \tilde{p}_{x_i^b}) = \underbrace{\frac{\delta}{4} \left(\frac{n\delta}{\alpha^2} \right)^{-\frac{n}{n+4}}}_{\text{const.}} \underbrace{\left(4\tilde{\gamma}_i^{\frac{n}{n+4}} + n\gamma_i \tilde{\gamma}_i^{-\frac{4}{n+4}} \right)}_{\text{reference mismatch}} \underbrace{N^{-\frac{4}{n+4}}}_{\text{conv. rate}}$$

223 where the first term is a constant and can be ignored, the second term is dependent on
 224 the mismatch between the true (2.21) and approximated (2.22) γ_i terms, and the third
 225 term determines the rate of convergence in N . Note that rate of coverage is sublinear
 226 but close to linear for small state-space dimensions n , and is purely sub-linear for even
 227 a modestly small n .

228 As the rate of convergence in (2.24) is sub-linear, as determined by the third term,
 229 the multiplicative terms in front of the rate of convergence play a very dominant role.
 230 The first term is constant, and thus can be ignored. The second term,

$$231 \quad (2.25) \quad 4\tilde{\gamma}_i^{\frac{n}{n+4}} + n\gamma_i\tilde{\gamma}_i^{-\frac{4}{n+4}},$$

232 therefore largely determines the error. It is trivial to see that (2.25) is minimized
 233 when $\tilde{\gamma}_i = \gamma_i$, therefore the error is only minimized when the reference distribution ϕ
 234 in (2.22) matches the true distribution in (2.21).

235 Thus, when there is a large discrepancy between the reference distribution ϕ and
 236 the true distribution, an adaptive choice of the bandwidth parameter (2.23) could
 237 produce a much more significant decrease in error than simply increasing the ensemble
 238 size N .

239 **3. Adaptive ensemble Gaussian Mixture Filter.** Following the observa-
 240 tions provided by Theorem 2.8 and by the discussion in subsection 2.2, we want
 241 to choose covariance matrices $\mathbf{B}_{i,j}$ found in the prior GMM assumption (2.6) in an
 242 intelligent and adaptive manner such that the convergence properties are satisfied.
 243 We additionally attempt to fulfill a desire useful to the practitioner: that practical
 244 convergence is achieved with as small as possible number of particles.

245 To that end, in this work we explore arbitrary parameterized covariance matrices
 246 in the prior GMM (2.6),

$$247 \quad (3.1) \quad p(x_i|\theta_i) = \sum_{j=1}^N u_{i,j} \mathcal{N}\left(x_i|x_{i,j}^b, \mathbf{B}_{i,j}^b(\theta_i)\right),$$

248 where each covariance $\mathbf{B}_{i,j}^b(\theta_i)$ is a matrix function of some (small number of) param-
 249 eters θ_i .

250 The aim of the parameterization in (3.1) is to find a set of parameters θ_i that
 251 can both be chosen adaptively at each step, and can ensure that the EnGMF does
 252 not violate the assumptions of Theorem 2.8 and, additionally, possibly attempts to
 253 minimize the error presented in subsection 2.2.

254 We now provide a way by which we can solve for the optimal parameters θ_i in (3.1)
 255 through the expectation maximization algorithm.

256 **3.1. Expectation Maximization.** The expectation maximization (EM) algo-
 257 rithm [8, 6] finds the set of the parameters θ_i that maximize $p(\theta_i|y_i)$ which is the
 258 conditional distribution of the parameters given the observations, at time index i .

259 Given some initial set of parameters $\theta_i^{(0)}$, the expectation maximization algorithm
 260 proceeds in an iterative fashion in two steps:

261 The *expectation step*,

$$262 \quad (3.2) \quad \mathbb{E}_{x_i^b|y_i, \theta_i^{(m)}} \log p(x_i^b, y_i, \theta_i)$$

263 constructs the function representing the expectation of the joint distribution of the
 264 prior state, the observations, and the parameters. The joint distribution in (3.2) can
 265 be written in terms of the prior (2.6), observation (2.8), and parameter distributions
 266 as,

$$267 \quad (3.3) \quad p(x_i^b, y_i, \theta_i) = p(y_i|x_i^b, \theta_i)p(x_i^b|\theta_i)p(\theta_i),$$

268 and, as the observation GMM (2.8) is not dependent on the parameters, (3.3) can be
 269 simplified to,

$$270 \quad (3.4) \quad p(x_i^b, y_i, \theta_i) = p(y_i|x_i^b)p(x_i^b|\theta_i)p(\theta_i),$$

271 where the prior distribution (3.1) is parameterized in terms of its covariance (3.1),
 272 and the parameter distribution

$$273 \quad (3.5) \quad p(\theta_i),$$

274 is determined on a case-by-case basis.

275 The *maximization step* aims to find the value of the parameters θ_i , that maximize
 276 the log joint distribution (3.4),

$$277 \quad (3.6) \quad \theta_i^{(m+1)} = \arg \max_{\theta_i} \mathbb{E}_{x_i^b|y_i, \theta_i^{(m)}} [\log p(x_i^b|\theta_i) + \log p(\theta_i)],$$

278 where $\theta_i^{(m)}$ are the parameters from the previous step, $x_i^b|y_i, \theta_i^{(m)}$ are samples from
 279 the posterior distribution (2.9) given the previous set of parameters $\theta_i^{(m)}$, and the
 280 term $\log p(y_i|x_i^b)$ is constant and thus can be safely ignored due to the fact that it
 281 does not influence the optimization problem. Recall Remark 2.4 that in the EnGMF,
 282 it is possible to generate an unlimited number of i.i.d. samples from the posterior
 283 distribution, thus the maximization step (3.6) can be computed to an arbitrary level
 284 of accuracy, given some reasonable assumptions on the distribution of the sample
 285 mean.

286 *Remark 3.1 (Invertible Covariances).* Note that the prior covariance $p(x_i^b|\theta_i)$
 287 in (3.6) requires that the covariance matrices $\mathbf{B}_{i,j}^b(\theta_i)$ in (3.1) are invertible, which is
 288 not necessarily required by the standard EnGMF.

289 *Remark 3.2 (Meaningful Representation of the Prior).* Note that the effect of
 290 the expectation maximization algorithm is to pick a parameterization of the prior
 291 estimate that best matches the posterior. For arbitrary parameterizations this would
 292 simply produce another copy of the posterior. The parameterization in (3.1) does
 293 not allow this to happen, as only the covariance is modified, and the mixture weights
 294 and means are not. This ensures that for (almost) all choices of the parameters θ ,
 295 the prior estimate is still a useful representation of the prior. This means that the
 296 EM algorithm merely chooses the prior estimate that is most useful in subsequently
 297 representing the posterior.

298 **3.1.1. Stochastic Optimization.** The maximization step (3.6) requires the so-
 299 lution of a stochastic optimization problem. Much of the recent literature on stochastic
 300 optimization has been focused on machine learning applications [1]. As the number
 301 of parameters in θ_i is small, it is possible to take advantage of methods that are
 302 built for the small parameter size case and that differ from typical machine learning
 303 optimization methods. Thus, in this work we utilize a variant of Newton’s method.

304 We can write the loss in the maximization step (3.6) as,

$$305 \quad (3.7) \quad \mathcal{L}(x_i, \theta_i) = \mathbb{E}_{x_i^b|y_i, \theta_i^{(m)}} [\log p(x_i|\theta_i) + \log p(\theta_i)],$$

306 where the posterior can be written as the following,

$$307 \quad (3.8) \quad x_i^{a,(m)} = x_i^b|y_i, \theta_i^{(m)},$$

308 which is useful shorthand for the following derivations.

309 One algorithm for finding the maximum of the loss function (3.7) is Newton’s
 310 method,

$$311 \quad (3.9) \quad \theta_i^{(m+1,p+1)} = \theta^{(m+1,p)} + \alpha_m \left(\mathbb{E}_{x_i^{a,(m)}} [\nabla_{\theta}^2 \mathcal{L}(x, \theta_i^{(m+1,p)})] \right)^{-1} \mathbb{E}_{x_i^{a,(m)}} [\nabla_{\theta} \mathcal{L}(x, \theta_i^{(m+1,p)})],$$

312 where α_m is the step-size (also known as the learning rate in the machine learning
 313 community), and the initial parameter value for the algorithm is $\theta_i^{(m+1,1)} := \theta_i^{(m)}$,
 314 which is the parameter from the previous maximization step (3.6). As it is challenging
 315 to compute the expected values in (3.9) analytically, some sort of approximation
 316 procedure is required.

317 As this is a stochastic optimization procedure, the two expected values in (3.9)
 318 cannot be calculated with the same random samples, as that would introduce unin-
 319 tended bias and variance into the update. In this work we utilize the sub-sampled
 320 version of Newton’s method [39] built specifically to handle this scenario. In sub-
 321 sampled Newton’s method, independent samples of $x_i^{a,(m)}$ are used to approximate
 322 the Hessian $\mathbb{E}_{x_i^{a,(m)}} [\nabla_{\theta}^2 \mathcal{L}(x, \theta_i^{(m+1,p)})]$ and the gradient $\mathbb{E}_{x_i^{a,(m)}} [\nabla_{\theta} \mathcal{L}(x, \theta_i^{(m+1,p)})]$. If
 323 the number of samples used is identical, then the number of samples required is
 324 double that of the stochastic gradient descent (SGD) algorithm which only requires
 325 the computation of $\mathbb{E}_{x_i^{a,(m)}} [\nabla_{\theta} \mathcal{L}(x, \theta^{(m+1,p)})]$. As Newton’s method achieves faster
 326 convergence than SGD, it is the authors’ belief that for this particular scenario the
 327 benefits of this approach outweigh the additional costs.

328 *Remark 3.3* (Quasi-Newton Methods). Instead of computing an estimate to the
 329 Hessian $\mathbb{E}_{x_i^{a,(m)}} [\nabla_{\theta}^2 \mathcal{L}(x, \theta_i^{(m+1,p)})]$ at every step, it is possible to only compute the
 330 Hessian at the initial step $\mathbb{E}_{x_i^{a,(m)}} [\nabla_{\theta}^2 \mathcal{L}(x, \theta_i^{(m+1,1)})]$ and use this approximation for
 331 all subsequent steps. This type of computationally efficient computation is a type of
 332 Quasi-Newton method [30] that is often used in practical applications.

333 *Remark 3.4* (Alternative Optimization Algorithms). Alternative stochastic opti-
 334 mization algorithms could also be utilized. The classic stochastic gradient descent
 335 algorithm [38] is an alternative which would require a smaller step-size α_m . Another
 336 alternative is ADAM [23] which would require to keep track of separate momentum
 337 and velocity terms.

338 *Remark 3.5* (Incremental Expectation Maximization). If the expectation maxi-
 339 mization algorithm is performed online in sequential data assimilation, it is not nec-
 340 essary to perform many steps of either the expectation maximization algorithm, or
 341 sub-sampled Newton’s method (3.9). In this work we initialize the parameters ex-
 342 pectation maximization algorithm subsection 3.1 from the previous time step of the
 343 data assimilation algorithm. This can be weakly justified as a type of incremental
 344 expectation maximization [28], and in the authors’ experience significantly increases
 345 the utility of the proposed approach.

346 We now discuss several different strategies for parameterizing the kernel covari-
 347 ance (3.1).

348 **3.2. Bandwidth-based covariance.** In kernel density estimation, choosing the
 349 optimal covariance matrices has had considerable research interest [40]. And, as dis-
 350 cussed in section 2 has a considerable impact on the efficacy of the EnGMF algorithm.

351 From subsection 2.2, recall the covariance parameterization,

352 (3.10)
$$\mathbf{B}_{i,j}^b(\beta_{i,N}^2) = \beta_{i,N}^2 \mathbf{P}_{i,N}^b, \quad 1 \leq j \leq N,$$

353 where $\beta_{i,N}^2$ is known as the bandwidth parameter [40]. This is one particular case
 354 of the covariance in the prior GMM (3.1) that is a focus of this work. The samples
 355 covariance in (3.10)

356 (3.11)
$$\mathbf{P}_{i,N}^b = \frac{1}{N-1} \mathbf{X}_i^b \left(I_N - \frac{1}{N} \mathbf{1}\mathbf{1}^T \right) \mathbf{X}_i^{b,T} \approx \mathbb{E} \left[(x_i^b - \mathbb{E}[(x_i^b)]) (x_i^b - \mathbb{E}[(x_i^b)])^T \right],$$

357 is known as the empirical covariance, and is an estimate of the covariance matrix of
 358 the prior state x_i^b . In (3.10), the only parameter is the bandwidth estimate, $\theta_i = \beta_{i,N}^2$.

359 *Remark 3.6* (Stochastic Newton’s for the Bandwidth Parameter). When the sub-
 360 sampled version of the stochastic Newton’s method (3.9) is applied to the covariance
 361 parameterized by the bandwidth parameter (3.10), then both the stochastic estimate
 362 of the gradient and the stochastic estimate of the Hessian are scalars. This enables
 363 the computation of the maximization step (3.6) to be performed with minimal linear
 364 system solves.

365 The prior kernel covariance estimate in (3.10) takes advantage of the underlying
 366 covariance of the data, and is thus a type of online estimate however, the resulting
 367 accuracy of the density estimate is still highly dependent on the bandwidth parameter
 368 $\beta_{i,N}^2$.

369 It is known from [40] that if the underlying exact prior distribution of x^b in (2.6) is
 370 Gaussian, that the optimal choice of bandwidth parameter β^2 in (3.10) that minimizes
 371 the mean integrated square error is,

372 (3.12)
$$\beta_{i,N,\text{Gaussian}}^2 = \left(\frac{4}{N(n+2)} \right)^{\frac{2}{n+4}},$$

373 which is also known as Silverman’s rule of thumb.

374 In practice most probability distributions of interest are not Gaussian, and (3.12)
 375 can result in a very poor approximation of the underlying density [40], thus a more
 376 refined choice of the bandwidth parameter is required.

377 *Theorem 2.8* showed that a sufficient condition for the convergence of the EnGMF
 378 is that the covariance estimate tends towards zero as $N \rightarrow \infty$. We now show a
 379 condition on the bandwidth parameter that is sufficient for the EnGMF to converge.

380 *LEMMA 3.7.* *Given the sequence of parameters $\{\beta_{i,N}^2\}_{N=1}^\infty$ parameterized by the*
 381 *particle amount N , a sufficient condition for the covariance estimate (3.10),*

382 (3.13)
$$\mathbf{B}_{i,N} = \beta_{i,N}^2 \mathbf{P}_{i,N}^b,$$

383 *to degenerate in the limit of particle number,*

384 (3.14)
$$\mathbf{B}_{i,N} \xrightarrow{D} \delta_{\mathbf{0}},$$

385 *is that the bandwidth parameter tends towards zero,*

386 (3.15)
$$\lim_{N \rightarrow \infty} \beta_{i,N}^2 = 0,$$

387 *in the limit of particle number N .*

388 *Proof.* Observe that,

$$389 \quad (3.16) \quad \begin{aligned} \mathbb{E}[\mathbf{B}_{i,N}] &= \beta_{i,N}^2 \mathbb{E}[\mathbf{P}_{i,N}^b], \\ \text{Cov}[\mathbf{B}_{i,N}] &= \beta_{i,N}^4 \text{Cov}[\mathbf{P}_{i,N}^b], \end{aligned}$$

390 meaning that both the mean and covariance of the random variable tend towards zero
391 as $N \rightarrow \infty$, as required. \square

392 **COROLLARY 3.8.** *The sequence $\{\beta_{i,N,\text{Gaussian}}^2\}_{N=1}^\infty$ of bandwidth parameters de-*
393 *fined by Silverman's rule of thumb (3.12) satisfies the conditions of Lemma 3.7.*

394 **Lemma 3.7** showed that when $\beta_{i,N}^2$ is a constant, and converges to zero in the limit
395 of particle number $N \rightarrow \infty$, the EnGMF converges. However, as we have uncertainty
396 about the bandwidth parameter, it is natural to think about it as a random variable
397 with some distribution. Thus an important choice is that of the distribution of the
398 bandwidth parameter. Care must be taken to ensure that this choice is sufficient to
399 make the resulting algorithm converge.

400 We provide a sufficient condition on the distribution of the bandwidth parameter
401 $\beta_{i,N}^2$ from (3.10), as an extension of **Theorem 2.8**. We therefore extend **Lemma 3.7** to
402 bandwidth parameters that are random variables with some prior distribution in the
403 expectation maximization algorithm.

404 **THEOREM 3.9.** *Given the sequence of random variables $\{\beta_{i,N}^2\}_{N=1}^\infty$ with a se-*
405 *quence of distributions $\{p(\beta_{i,N}^2)\}_{N=1}^\infty$ parameterized by the particle amount N , a suf-*
406 *ficient condition for the covariance estimate (3.10),*

$$407 \quad (3.17) \quad \mathbf{B}_N = \beta_{i,N}^2 \mathbf{P}_N^b,$$

408 *to tend towards zero in distribution in the limit of particle number,*

$$409 \quad (3.18) \quad \mathbf{B}_N \xrightarrow{D} \mathbf{0},$$

410 *is that the distribution of the bandwidth parameter tends towards the delta distribution*
411 *around zero,*

$$412 \quad (3.19) \quad \lim_{N \rightarrow \infty} p(\beta_{i,N}^2) = \delta_0,$$

413 *ensuring that $\beta_{i,N}^2$ almost surely becomes 0.*

414 *Proof.* If the distribution of $\beta_{i,N}^2$ converges to δ_0 , then the solution to the max-
415 imization step in the EM algorithm (3.6) almost surely becomes a constant, namely
416 that $\beta_{i,N}^2 \xrightarrow{\text{a.s.}} 0$, as required. \square

417 **3.2.1. Choosing the bandwidth distribution.** One way in which the condi-
418 tions of **Theorem 3.9** could be satisfied is through an intelligent choice of the proba-
419 bility distribution of the bandwidth parameter $p(\beta_{i,N}^2)$.

420 A common choice in the literature, the principal of maximum entropy (PME) [22]
421 could be used to find a good candidate for this distribution. If we assume that the
422 expected value of the bandwidth, $\beta_{i,N}^2$, is Silverman's rule of thumb (3.12), and we
423 have no other information available, then the distribution that satisfies the PME is
424 the exponential distribution,

$$425 \quad (3.20) \quad p(\beta_{i,N}^2) = \beta_{i,N,\text{Gaussian}}^{-2} e^{-\beta_{i,N,\text{Gaussian}}^{-2} \beta_{i,N}^2},$$

426 this distribution, however, always has a single mode at zero, thus, from the authors’
 427 experience, is ill-suited for use in expectation maximization.

428 It is possible to perform some slight-of-hand in order to make this assumption
 429 more tractable. It is more efficient to look at $\beta_{i,N}$, the square-root of the bandwidth
 430 parameter. If (3.20) is the distribution of $\beta_{i,N}^2$, then $\beta_{i,N}$ is distributed according to

$$431 \quad (3.21) \quad p(\beta_{i,N}) = 2\beta_{i,N,\text{Gaussian}}^{-2} \beta_{i,N} e^{-\beta_{i,N,\text{Gaussian}}^{-2} \beta_{i,N}^2}$$

432 which is the Rayleigh distribution [31] with known mode $2^{-1/2}\beta_{i,N,\text{Gaussian}}$. We as-
 433 sume this Rayleigh distribution (3.21) on the bandwidth parameter for the remainder
 434 of this paper.

435 It is also possible to assume a more general distribution around $\beta_{i,N}^2$, such as a
 436 gamma distribution, though this choice would introduce another free parameter into
 437 the algorithm; an undesirable outcome.

438 While the parameterized covariance in (3.10) is well-studied, it has a few limita-
 439 tions that prevent it from being used in high-dimensional inference, chief among those
 440 being the fact that the covariance estimate in (3.11) can potentially be low-rank, and
 441 thus generate a covariance that is not invertible Remark 3.1, thus we can introduce
 442 covariance matrix estimates that have extra parameters in order to mitigate this issue.

443 *Remark 3.10.* It is important to note that the optimal bandwidth is deterministic,
 444 but unknown. The uncertainty that transforms our knowledge about the bandwidth
 445 into a random variable is purely from the point of view of the agent performing the
 446 inference. For a more in-depth discussion about choosing distributions for parameters
 447 see [22].

448 **3.3. Covariance Shrinkage Estimates.** Covariance shrinkage [13, 12, 14, 11,
 449 24] aims to use extra prior information about the covariance of x_i^b in (2.6) in order
 450 to have a more accurate covariance estimate in the case when the number of samples
 451 is smaller than the dimension of the dynamical system $N < n$. Covariance shrinkage
 452 methods have previously been employed for ensemble data assimilation [29, 34] and
 453 for regularization in particle filtering [35].

454 Assume that we have prior information about the covariance structure of x_i^b in
 455 the form of a ‘target’ covariance matrix \mathbf{T}_i . The covariance shrinkage estimate of the
 456 covariance, scaled by the bandwidth, is given by,

$$457 \quad (3.22) \quad \mathbf{B}_{i,j}^b = \beta_{i,N}^2 \left[\gamma_i \mu_i \mathbf{T}_i + (1 - \gamma_i) \mathbf{P}_i^b \right]$$

458 where,

$$459 \quad (3.23) \quad \mu_i = n^{-1} \text{tr } \mathbf{C}_i, \quad \mathbf{C}_i = \mathbf{T}_i^{-\frac{1}{2}} \mathbf{P}_i^b \mathbf{T}_i^{-\frac{1}{2}}$$

460 is a rescaling factor, and γ_i is the shrinkage factor, which we treat as a parameter.

461 Under Gaussian assumptions on the samples, $x_{i,j}^b$, a good known shrinkage factor
 462 is,

$$463 \quad (3.24) \quad \gamma_{i,\text{RBLW}} = \min \left[\frac{N-2}{N(N+2)} + \frac{(n+1)N-2}{N(N+2)(n-1)\hat{U}_i}, 1 \right],$$

$$\hat{U}_i = \frac{1}{n-1} \left(\frac{n \text{tr } \mathbf{C}_i^2}{\text{tr}^2 \mathbf{C}_i} - 1 \right)$$

464 called the Rao-Blackwell Ledoit-Wolf estimator [13].

465 One possible choice of the target matrix \mathbf{T}_i that does not require any prior knowl-
 466 edge is the diagonal of the empirical covariance (3.11),

$$467 \quad (3.25) \quad \mathbf{T}_i = \boxtimes \mathbf{P}_i^b,$$

468 where the notation of \boxtimes is introduced to signify the matrix consisting of only the
 469 diagonal of the subsequent term. Observe that,

$$470 \quad (3.26) \quad \text{tr} \left[\left(\boxtimes \mathbf{P}_i^b \right)^{-\frac{1}{2}} \mathbf{P}_i^b \left(\boxtimes \mathbf{P}_i^b \right)^{-\frac{1}{2}} \right] = n,$$

471 meaning that when the target matrix is defined by (3.25), the scaling factor $\mu_i = 1$.
 472 This means that direct calculation the matrix \mathbf{C}_i , from (3.23), is only required for the
 473 calculation of $\gamma_{i,\text{RBLW}}$ in (3.24).

474 *Remark 3.11* (On $p(\gamma_i)$). A commonly made assumption is that parameters are
 475 independently distributed, therefore the distribution of $p(\beta_{i,N}^2)$ can be chosen inde-
 476 pendently of the distribution $p(\gamma_i)$. As there are no requirements that the optimal γ_i
 477 is dependent on ensemble size, it is a natural choice to assume a uniform likelihood,

$$478 \quad (3.27) \quad p(\gamma_i) \propto 1,$$

479 which is a typical assumption in parameter estimation [8].

480 **3.4. Covariance Localization.** In the geosciences, states usually have some
 481 sort of innate spatial structure. State variables that are spatially far apart are gener-
 482 ally more weakly correlated than states that are closer together. Taking advantage of
 483 this fact, covariance localization [5] is a matrix tapering technique which aims to re-
 484 duce the impact of spurious correlations caused by undersampled ($N \ll n$) covariance
 485 matrix estimates.

486 In this work we focus on what is known as the B-localization methodology, and
 487 combine it with the bandwidth scaling (3.10) in the following manner,

$$488 \quad (3.28) \quad \mathbf{B}_{i,j}^b = \beta_i^2 \left(\rho(r_i) \circ \mathbf{P}_i^b \right),$$

489 where the matrix $\rho(r_i)$ contains a set of decorrelation variables parameterized by the
 490 localization radius r_i , and \circ is the element-wise Schur product.

491 A common choice for ρ is known as Gaussian localization,

$$492 \quad (3.29) \quad \rho(r_i)_{\ell,q} = e^{-\frac{1}{2} \frac{d(\ell,q)^2}{r_i^2}},$$

493 where $d(\ell, q)$ represents the spatial distance between the variables at index ℓ and index
 494 q .

495 *Remark 3.12* (Choice of Localization Radius r_i). The choice of localization ra-
 496 dius r_i in (3.28) can be informed by the temporal covariance of the model of interest
 497 if the model of interest is Ergodic, however in practice, the best localization radius is
 498 almost always determined empirically.

499 Adaptive-in-time choices for r_i have been explored for the ensemble Kalman filter
 500 in [33].

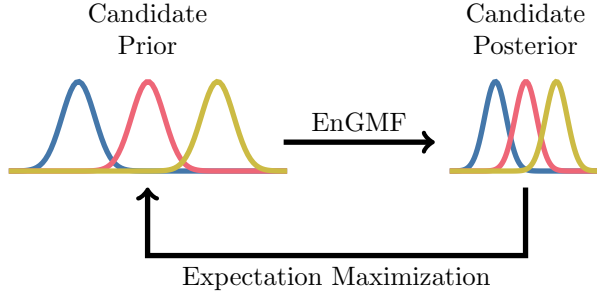


FIG. 1. An illustration of the outer loop of the AEnGMF algorithm. In this example the candidate prior is made up of three distinct Gaussian modes which are plotted separately. The candidate prior is transformed into the candidate posterior through the standard EnGMF update (2.10). Next the expectation maximization algorithm (3.2) and (3.6) is performed, and a new candidate prior is obtained. This procedure repeats until some desired level of convergence is achieved.

501 **3.5. Practical Implementation of the AEnGMF.** We are now able to combine
 502 all the elements presented in this section to fully describe the inner-workings of
 503 the adaptive ensemble Gaussian mixture filter (AEnGMF). The AEnGMF operates
 504 as follows. First a choice of parameterized covariance is made by the user. This
 505 choice determines the parameters that are optimized for. Next a choice of parameter
 506 distribution is required. In this work the bandwidth parameter is assumed to be distributed
 507 according to the Rayleigh distribution subsection 3.2.1, and the rest of the
 508 parameters are assumed to be proportional to one, thus of no additional consequence.

509 At each step of the algorithm, the previous choice of covariance parameters is
 510 carried over, $\theta_i^{(1,1)} := \theta_{i-1}$. This choice from Remark 3.5 is motivated by incremental
 511 approaches to expectation maximization, and lends itself particularly well to parameterized
 512 covariances that do not depend on their parameters changing a lot from step
 513 to step.

514 Next, M iterations of the expectation-maximization algorithm are performed.
 515 The expectation maximization algorithm can be treated as the ‘outer-loop’ [30]
 516 in this optimization procedure.

517 The cost function is solved using P loops of sub-sampled Newton’s method (3.9)
 518 with a constant learning-rate α making this the ‘inner-loop’ algorithm. As it is possible
 519 to sample from the posterior arbitrarily Remark 2.4, the gradient and Hessian
 520 calculations can be performed using a different number of samples, S , than that of
 521 the number of particles N . Specifically, the gradient is computed using S samples
 522 from the candidate posterior, and the Hessian is computed using S separate samples
 523 from the candidate posterior, for a total of $2S$ samples.

524 As the AEnGMF is a particle filter, resampling of N particles is performed at the
 525 end of the algorithm with the EnGMF resampling procedure Procedure 2.3.

526 The outer loop of the algorithm is illustrated in Figure 1, and a detailed step-by-
 527 step look at the algorithm can be seen in Algorithm 3.1.

528 *Remark 3.13* (Choosing M , P , α , and S). In the authors’ experience, it is much
 529 more advantageous to perform multiple iterations of the expectation maximization
 530 algorithm than that of sub-sampled Newton’s method, thus it is advantageous to
 531 take $M \geq P$. It is also advantageous to oversample the gradient and Hessian, thus
 532 $S \geq N$. By far the hardest choice to make is that of the learning-rate α . A learning
 533 rate that is too large ($\alpha \approx 1$) could cause the algorithm to become unstable, thus

Algorithm 3.1 The Adaptive Ensemble Gaussian Mixture Filter

Input Initial ensemble \mathbf{X}_0^a , initial estimate for the parameters θ_0 , outer loop iteration count M , inner loop iteration count P , learning rate α , and number of internal samples S .

```
for  $i = 1, \dots$  do
  % Propagate the ensemble forward in time through the model
   $\mathbf{X}_i^b = \mathcal{M}(\mathbf{X}_{i-1}^a)$ 
  % Initialize the parameters  $\theta$  to the parameters from the previous step
   $\theta_i^{(1,1)} := \theta_{i-1}$ 
  % Perform the expectation maximization loop  $M$  times
  for  $m = 1, \dots, M$  do
    % Construct loss function
     $\mathcal{L}(x, \theta) := \mathbb{E}_{x^b|y, \theta_i^{(m)}} [\log p(x|\theta) + \log p(\theta)]$ 
    % Initialize the inner loop  $\theta$  parameter
     $\theta_i^{(m+1,1)} := \theta_i^{(m,P+1)}$ 
    % Perform  $P$  steps of subsampled Newton's method.
    for  $p = 1, \dots, P$  do
      % Sample  $S$  particles from the candidate posterior.
       $\mathbf{X}^a \sim_{(S)} \pi(x|y, \theta_i^{(m)})$ 
      % Compute the loss gradient
       $g := \mathbb{E}_{\mathbf{X}^a} [\nabla_{\theta} \mathcal{L}(X, \theta_i^{(m+1,p)})]$ 
      % Similarly, compute sample Hessian using different samples
       $\mathbf{X}^a \sim_{(S)} \pi(x|y, \theta_i^{(m)})$ 
       $H := \mathbb{E}_{\mathbf{X}^a} [\nabla_{\theta}^2 \mathcal{L}(X, \theta_i^{(m+1,p)})]$ 
      % Compute new estimate of the parameters
       $\theta_i^{(m+1,p+1)} := \theta_i + \alpha H^{-1} g$ 
    end for
    % Set the current  $\theta$  parameter
     $\theta_i^{(m+1)} := \theta_i^{(m+1,P+1)}$ 
  end for
  % Set the  $\theta$  parameter for the current time index
   $\theta_i := \theta_i^{(M+1)}$ 
  % Sample a new ensemble of  $N$  particles with new parameters  $\theta_i$ 
   $\mathbf{X}_i^a \sim_{(N)} \pi(x|y, \theta_i)$ 
end for
```

534 increasing the cost instead of decreasing it, and thus make it choose parameters θ
535 that are worse than the original choices. A learning rate that is too small ($\alpha \rightarrow 0$)
536 could lead to parameters that react poorly to the changing conditions of the states.
537 From the practitioner's point of view, this is by far the most important parameter
538 to choose correctly. Ideally, the parameter α can be chosen through some type of
539 line search technique [30] that ensures that steps are always taken in a direction that
540 decreases the error, though this is not explored in this work.

541 *Remark 3.14* (Considerations for the High-dimensional Setting). There are many
542 considerations to be made for getting the AEnGMF to work in the high-dimensional
543 setting. First is that the computation of the covariance cannot be made explicitly.

544 This can be resolved by utilizing a covariance estimate that does not have to be ex-
 545 plicitly computed like that of the shrinkage estimate in (3.22). Matrix inverse vector
 546 products of (3.22) can also be computed without explicit computation of the entire
 547 matrix. The normalizing factor issue in Remark 2.2 can also be mitigated by this
 548 covariance matrix estimate. Another problem is the resampling procedure in Pro-
 549 cedure 2.3, which requires the computation of matrix square root vector products.
 550 Methods such as those proposed in [2, 19, 15] could be utilized to solve this issue,
 551 though this is still an open problem. A final consideration is that the intermediate re-
 552 sults or approximations of the densities could represent non-physical states, thus spe-
 553 cial consideration must be given to those problems, especially in the high-dimensional
 554 setting.

555 *Remark 3.15* (Computational complexity of the AEnGMF). Take n to be the
 556 dimension of the full state, m to be the dimension of the observations, q to be the
 557 dimensions of the parameters θ , and N to be the ensemble size. The dominant terms
 558 of the EnGMF update involve constructing the covariance, computing the gain matrix
 559 and updating the covariance, which in the worst case has computational complexity
 560 $\mathcal{O}(m^3N + m^2nN + n^2N)$. The resampling procedure Procedure 2.3 has complexity
 561 $\mathcal{O}(n^3N)$. For the AEnGMF, the cost of computing the Hessian in the worst case,
 562 where the gradients need to be computed by repeat evaluation of the cost function,
 563 the complexity becomes $\mathcal{O}(q^2n^3N)$. Thus, the total complexity of the AEnGMF is

$$564 \quad (3.30) \quad \mathcal{O}(MP[C_{\text{EnGMF}} + q^2n^3N] + C_{\text{EnGMF}}),$$

565 where C_{EnGMF} is the complexity of the EnGMF. Thus, the cost of the algorithm has
 566 to be weighted against the cost of propagating more particles through the forward
 567 model dynamics.

568 **4. Numerical Experiments.** The aim of the numerical experiments is first to
 569 demonstrate the viability and convergence of the proposed AEnGMF on small-scale
 570 problem, and secondly to demonstrate the more complicated covariance parameteri-
 571 zation approaches on a larger-scale problem.

572 **4.1. Lorenz '63.** With the first set of experiments we aim to look at a highly
 573 non-linear system with a non-linear observation operator. We focus on the stan-
 574 dard EnGMF case of the Kernel covariance parameterized by the bandwidth param-
 575 eter (3.10).

576 We take the the 3-variable Lorenz '63 equations [26],

$$577 \quad (4.1) \quad \begin{aligned} x' &= \sigma(y - x), \\ y' &= x(\rho - z), \\ z' &= xy - \beta z, \end{aligned}$$

578 with canonical parameters $\sigma = 10$, $\rho = 28$, and $\beta = \frac{8}{3}$. The time between assimilations
 579 is taken to be $\Delta t = 0.5$, which allows the system enough time to evolve in a highly
 580 non-linear manner. The non-linear dynamics are propagated through time with an
 581 adaptive Runge-Kutta method [16] with absolute and relative tolerances of 10^{-11} in
 582 order to simulate a costly forward model calculation.

583 It is known [20] that the system (4.1) has three critical points, one at the origin,

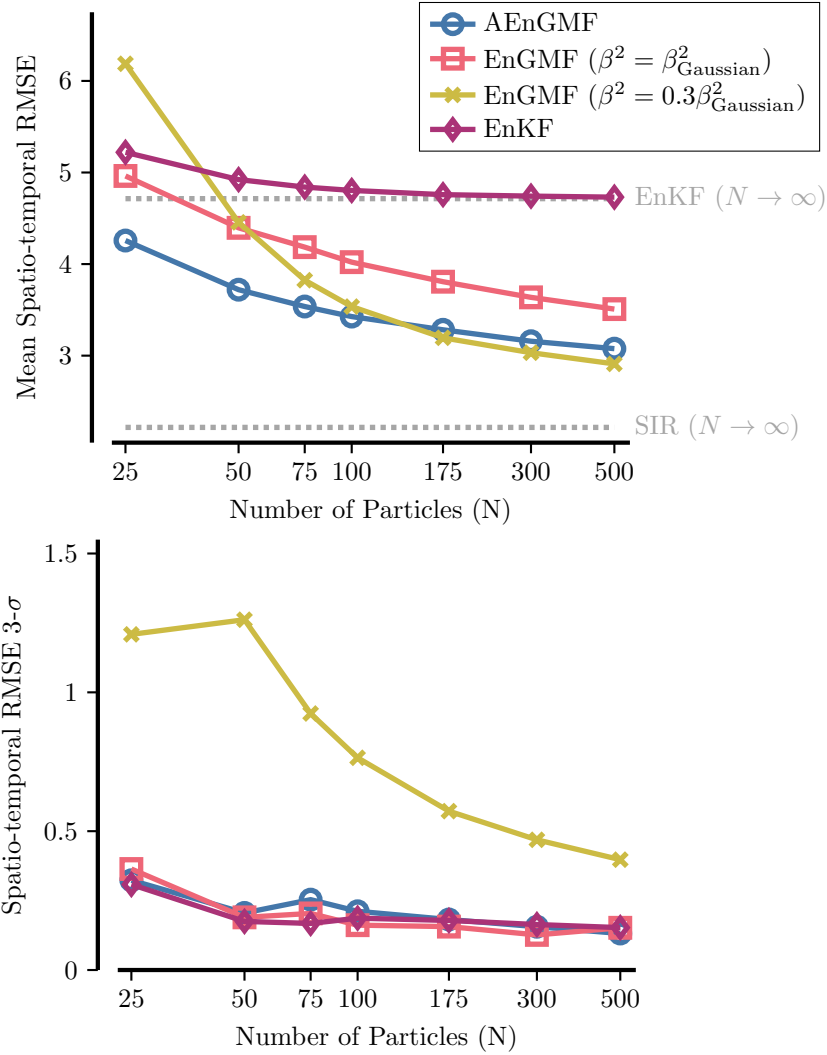


FIG. 2. RMSE simulation results for the Lorenz '63 equations for four different data assimilation algorithms. The top plot represents the mean RMSE value across all runs, while the bottom plot represents three standard deviations of error around the mean. The blue line with circular marks represents the AEnGMF, the red line with square marks represents the canonical EnGMF with Silverman's rule of thumb, the yellow line with x marks represents the EnGMF with a scaled Silverman's rule of thumb, and the Raspberry line with diamond marks represents the EnKF. Two baseline lines, running the EnKF and a particle filter (SIR) for a high particle number are also provided to provide theoretical bounds.

584 and in the center of each of the butterfly wings. The first non-zero critical point,

$$\begin{aligned}
 x_c &= \sqrt{\beta(\rho - 1)}, \\
 y_c &= \sqrt{\beta(\rho - 1)}, \\
 z_c &= \rho - 1,
 \end{aligned}
 \tag{4.2}$$

586 defines the center of one of the wings of the butterfly, with the other center being
 587 $(-x_c, -y_c, z_c)$ and the origin being $(0, 0, 0)$. For the non-linear observation operator

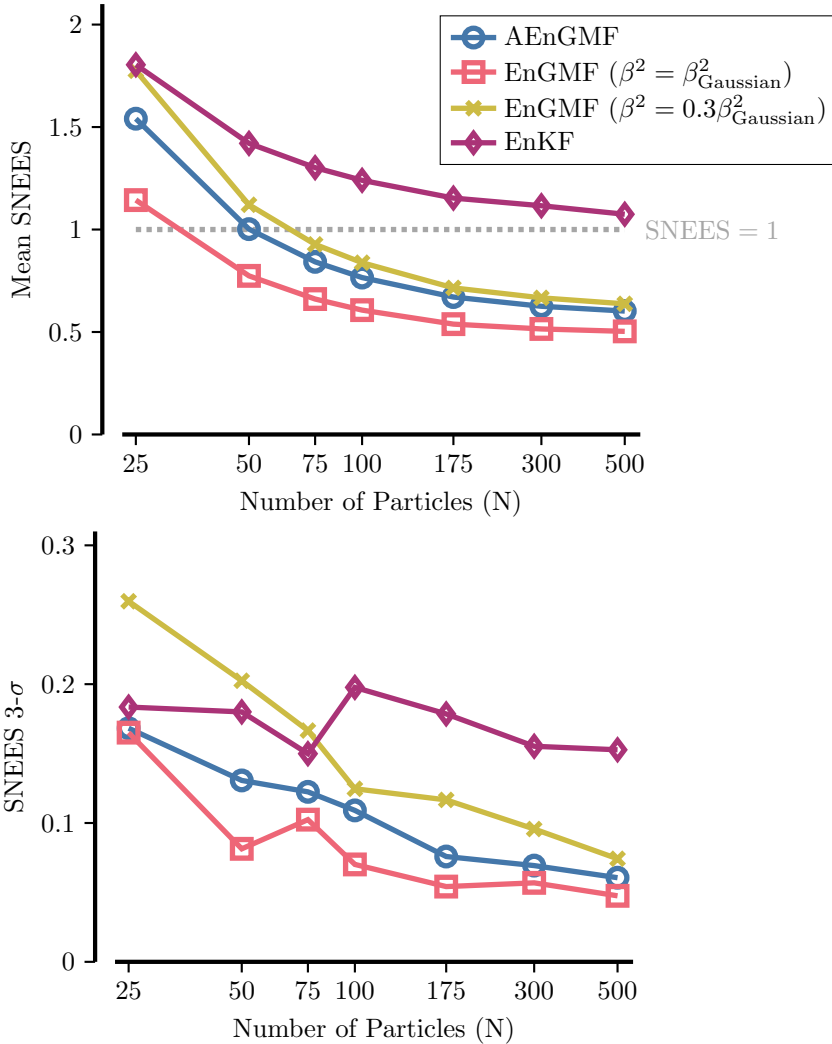


FIG. 3. SNEES simulation results for the Lorenz '63 equations for four different data assimilation algorithms. The top plot represents the mean SNEES value across all runs, while the bottom plot represents three standard deviations of error around the mean. The lines representing the algorithms are identical to those in Figure 2. The dashed lines around the main lines represent three standard deviations over the samples. The ideal SNEES of one is represented by the constant gray dashed line.

588 we measure the distance from the critical point (4.2) to the point being measured,

589 (4.3)
$$\mathcal{H}(x, y, z) = \sqrt{(x - x_c)^2 + (y - y_c)^2 + (z - z_c)^2},$$

590 as a scalar observation, with Gaussian error with an error variance of $\mathbf{R} = 1$.

591 The goal of this experiment is to show that the various variants of the EnGMF
 592 are superior to the ensemble Kalman filter (EnKF) and converge to exact Bayesian
 593 inference in the limit of particle number ($N \rightarrow \infty$). We therefore calculate two
 594 reference boundaries for this problem, one using the EnKF for a large ensemble size
 595 ($N = 1000$) and for the sequential importance resampling (SIR) particle filter with a

596 large number of particles ($N = 1000$), specifically the variant found in [37].

597 It is known in the literature [40, 21] that Silverman’s rule-of-thumb is usually
 598 an over-estimate of the optimal bandwidth term. We thus attempt to find a scaling
 599 factor $0 < s < 1$ such that the bandwidth parameter defined by the product,

$$600 \quad (4.4) \quad \beta_{i,N}^2 = s\beta_{\text{Gaussian}}^2,$$

601 would produce the minimal error for our choice of number of particles. From a quick
 602 parameter sweep, it was determined that $s = 0.3$ provides a good factor, that is
 603 optimal for a high number of particles.

604 Thus, we run four different algorithms, the EnKF, the AEnGMF, the EnGMF
 605 with Silverman’s rule of thumb, and the EnGMF with Silverman’s rule of thumb
 606 scaled by $s = 0.3$, for this model setup for various numbers of particles ranging from
 607 $N = 25$ to $N = 500$.

608 All experiments were run for 48 independent initial ensembles, with the same truth
 609 but different observations, for 5500 assimilation steps with the first 500 discarded for
 610 spinup, meaning that the first 500 steps do not count into the error calculations to
 611 let the filter reach a steady state. The mean of the spatio-temporal RMSE,

$$612 \quad (4.5) \quad \text{RMSE}(\bar{x}, x^t) = \sqrt{\frac{1}{nT} \sum_{i=1}^T \sum_{l=1}^n (\bar{x}_{i,l} - x_{i,l}^t)^2},$$

613 is calculated between the statistical mean \bar{x} and the truth x^t , over the 12 runs and
 614 is the metric by which the efficacy of the algorithms is determined. In order to check
 615 the consistency, the mean of the scaled normalized estimated error squared (SNEES)
 616 metric [44] is utilized,

$$617 \quad (4.6) \quad \text{SNEES}(\bar{x}, x^t) = \frac{1}{nT} \sum_{i=1}^T (\bar{x}_i - x_i^t)^T \mathbf{P}_i^{a,-1} (\bar{x}_i - x_i^t),$$

618 where \mathbf{P}_i^a is the estimate of the posterior covariance at time index i . A SNEES of
 619 one is considered to be ideal, as that means that the error predicted by the filter is
 620 in line with the actual error of the filter. Additionally, if the SNEES is not one, it
 621 is better for the filter to be more conservative, meaning the SNEES is less than one,
 622 that overconfident, meaning a SNEES greater than one.

623 For the choices of parameters in Algorithm 3.1, we choose $M = 5$ loops of the
 624 expectation maximization algorithm, $P = 1$ loops of sub-sampled Newton’s method,
 625 sampling $S = N$ exactly as many samples as there are particles, and a high learning
 626 rate of $\alpha = 1$. The previous parameter choices were hand-tuned to balance error and
 627 time to solution. The Rayleigh distribution with mean of Silverman’s rule-of-thumb
 628 is chosen for the bandwidth parameter just like in subsection 3.2.1.

629 The results of the RMSE experiments are visually demonstrated in Figure 2.
 630 As can be seen, the AEnGMF is consistently lower in error than the EnGMF with
 631 bandwidth equivalent to Silverman’s rule-of-thumb, and provides lower error in the
 632 particle number range of $N = 25$ to $N = 100$. The EnGMF with scaling factor $s = 0.3$
 633 is the only algorithm to perform worse than the EnKF for $N = 25$, but also produces
 634 the lowest error between $N = 300$ and $N = 500$. Crucially, almost all algorithms have
 635 the same error bounds, as shown by the $3\text{-}\sigma$ plot of the RMSE, except for the The
 636 EnGMF with scaling factor $s = 0.3$, which has a significantly higher error standard
 637 deviation. This means that the AEnGMF lowers the error without sacrificing stability.

638 For the SNEES of the experiments, demonstrated in [Figure 3](#), the AEnGMF lies
 639 between the EnGMF and the EnGMF with scaling facotr $s = 0.3$, both in terms
 640 of raw SNEES value and in terms of standard deviation. As both algorithms are
 641 conservative for higher particle numbers, this is not a particularly surprising result.

642 It can be seen that the EnGMF family of methods converges to the SIR limit
 643 very slowly, driven by the sublinear rate of convergence. The practical use of the
 644 Rayleigh distribution [subsection 3.2.1](#) for achieving this effect can be questioned as
 645 the convergence is clearly sub-optimal. A better choice of the distribution of the
 646 bandwidth parameter is required.

647 Finally, we report the practical computational time increase of the AEnGMF over
 648 the EnGMF. For the full forecast-analysis loop a 1.2 to 1.7 times increase in compu-
 649 tational cost was observed, meaning that when the computational time is dominated
 650 by the forward model runs, the cost of the AEnGMF is not an overly significant
 651 computational burden.

652 **4.2. Lorenz '96.** For the next set of experiments, we look at the case of an un-
 653 dersampled ($N \ll n$) estimate of the prior distribution. We look at two different types
 654 of covariance matrix parameterizations: one based on covariance shrinkage [\(3.22\)](#)
 655 with the target matrix being the diagonal of the statistical covariance [\(3.25\)](#), and a
 656 localization-based covariance matrix estimate [\(3.28\)](#).

657 For the model of interest we take the 40-variable Lorenz '96 equations [\[27\]](#),

658 (4.7)
$$x'_k = -x_{k-1}(x_{k-2} - x_{k+1}) - x_k + F \dots, \quad k = 1, \dots, 40,$$

659 with cyclic boundary conditions, and the forcing factor $F = 8$. For the time between
 660 assimilations we take one day of model time which is equivalent to a $\Delta t = 0.2$, leading
 661 to a high level of non-linearity in the system. The non-linear dynamics are propagated
 662 through time with an adaptive Runge-Kutta method [\[16\]](#) with absolute and relative
 663 tolerances of 10^{-6} .

664 We want to compare the AEnGMF approach of adaptively choosing the param-
 665 eters of the Kernel covariance with that of the more classic EnGMF approach where
 666 the parameters are determined by a known good heuristic. We also want to compare
 667 with a base-line state-of-the-art algorithm, the localized ensemble Kalman filter. To
 668 that end, we perform experiments on the following set of filters:

- 669 1. the shrinkage-based AEnGMF (Shr-AEnGMF), with parameters of β_i^2 for
 670 the bandwidth and $\zeta_i = \tanh^{-1} \gamma_i$, for an unbounded transformation of the
 671 shrinkage parameter $0 < \gamma_i < 1$,
- 672 2. the shrinkage-based EnGMF (Shr-EnGMF) with bandwidth defined by Silver-
 673 man's rule-of-thumb [\(3.12\)](#), $\beta_i^2 = \beta_{i,\text{Gaussian}}^2$ and the RBLW [\(3.24\)](#) shrinkage
 674 parameter $\gamma_i = \gamma_{i,\text{RBLW}}$,
- 675 3. the localized AEnGMF (LAEnGMF) with parameters of β_i^2 for the bandwidth
 676 and $\zeta_i = \sqrt{r_i}$ for an unbounded transformation of the localization radius
 677 $0 < r_i$,
- 678 4. the localized EnGMF (LEnGMF) with bandwidth defined by Silverman's
 679 rule-of-thumb [\(3.12\)](#), $\beta_i^2 = \beta_{i,\text{Gaussian}}^2$ and a fixed radius of $r_i = 4$,
- 680 5. and the localized EnKF (LEnKF) with fixed radius $r_i = 4$ for a useful com-
 681 parison with a state-of-the-art filter.

682 For the non-linear observation operator, we take the point-wise non-linear operator,

683 (4.8)
$$\mathcal{H}(x_i) = \frac{x_i}{2} \left[1 + \left(\frac{|x_i|}{10} \right)^{\omega-1} \right],$$

20

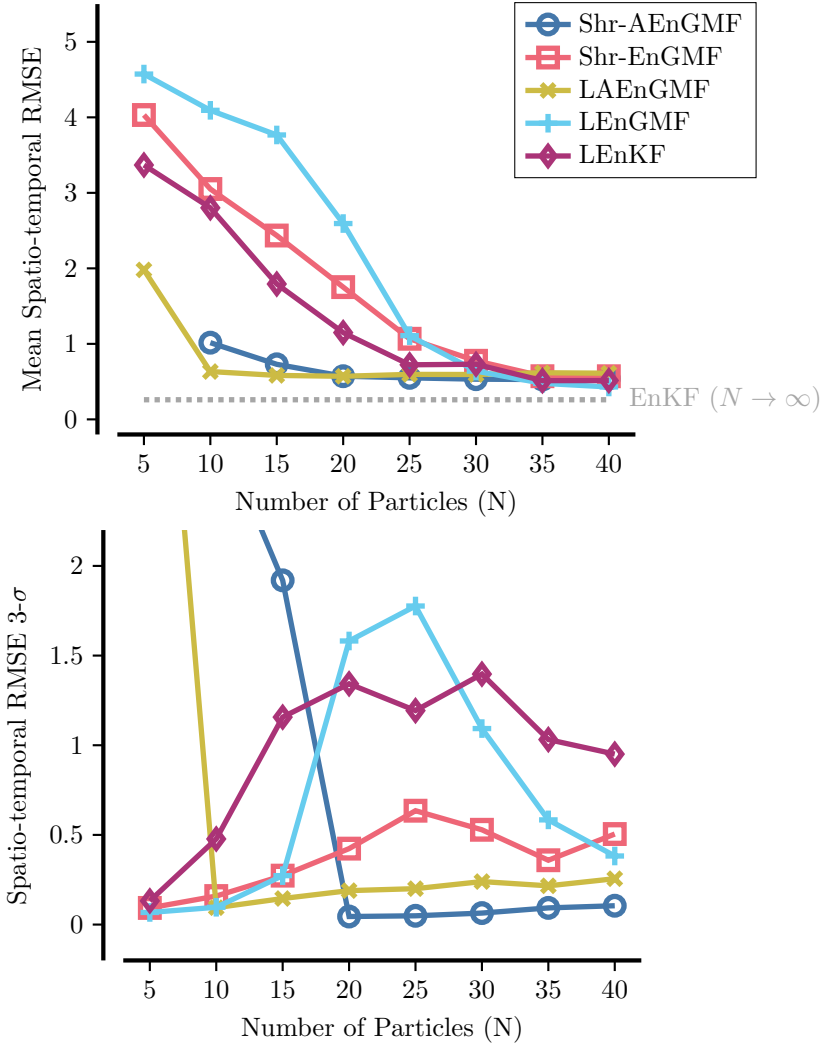


FIG. 4. Simulation results for the Lorenz '96 equations for five different data assimilation algorithms. The dark blue line with circle marks represents the AEnGMF with a shrinkage-based estimate to the covariance, with the light red line with square represents the standard EnGMF with a shrinkage-based estimate to the covariance, the yellow line with x marks represents a localized AEnGMF, the light-blue line with plus marks represents a localized EnGMF, and the raspberry line with diamond marks represents the localized EnKF.

684 as found in [5], with $\omega = 5$ for a medium level of non-linearity, with the observation
 685 covariance matrix being set to $\mathbf{R} = \frac{1}{4}\mathbf{I}_{40}$. The number of particles is taken to range
 686 from as little as $N = 5$ to as high as $N = 40$.

687 All experiments were run for 48 independent initial ensembles, with the same truth
 688 but different observations, for 5500 assimilation steps with the first 500 discarded for
 689 spinup, meaning that the first 500 steps do not count into the error calculations to let
 690 the filter reach a steady state. For our error metric we again take the spatio-temporal
 691 RMSE (4.5).

692 For the choices of parameters in Algorithm 3.1, we choose $M = 1$ loops of the

693 expectation maximization algorithm, $P = 1$ loops of sub-sampled Newton’s method,
694 sampling $S = 100$ to sample in the excess, and a low learning rate of $\alpha = 1e -$
695 2. The Rayleigh distribution with mean of Silverman’s rule-of-thumb is chosen for
696 the bandwidth parameter just like in [subsection 3.2.1](#), with both the radius r and
697 shrinkage parameter γ having distributions proportional to one along all of their
698 support as they are not required for convergence.

699 The results for this round of experiments can be seen in figure [Figure 4](#). At around
700 $N = 30$ particles, all algorithms perform roughly the same, thus the interesting be-
701 havior occurs when there are fewer particles. Both versions of the EnGMF without
702 adaptive covariance estimates (Shr-EnGMF and LEnGMF) perform worse than the
703 localized EnKF. The adaptive versions of the same algorithms (Shr-AEnGMF and
704 LAEnGMF) perform significantly better than all other tested algorithms. The the
705 LAEnGMF practically converges for $N = 10$ particles, and the Shr-AEnGMF practi-
706 cally converges for $N = 20$ particles. Additionally the Shr-AEnGMF and LAEnGMF
707 have tighter error bounds than all the other tested algorithms, potentially signifying
708 that the adaptive nature of the algorithm is better able to handle outlier scenar-
709 ios. These results highlight the need and utility of the adaptive covariance estimate
710 approach in the EnGMF presented in this paper.

711 **5. Conclusions.** By leveraging parameterized sample covariance estimates and
712 the expectation maximization algorithm, this work introduced the adaptive ensem-
713 ble Gaussian mixture filter (AEnGMF) as an extension of the ensemble Gaussian
714 mixture filter (EnGMF). Theoretical results about the convergence properties of this
715 filter were derived by making assumptions about the distribution of the kernel band-
716 width. Numerical results have verified the theoretical convergence properties of the
717 AEnGMF, and have shown that for a certain set of parameters the AEnGMF has
718 superior convergence to that of the EnGMF.

719 Future work could extend the AEnGMF to a smoothing [\[37\]](#) framework, a hybrid
720 filtering [\[18\]](#) framework, and to a multifidelity filtering [\[32\]](#) framework. An active
721 research direction is in applying the AEnGMF to a real-world orbit tracking prob-
722 lem [\[44\]](#). Work exploring practical consideration on choosing the parameters discussed
723 in [Remark 3.13](#) is also of independent interest.

724 **Acknowledgments.** The authors would like to thank the two anonymous re-
725 viewers for their very detailed and thorough comments that have helped significantly
726 increase the quality of this work.

727

REFERENCES

- 728 [1] C. C. AGGARWAL ET AL., *Neural networks and deep learning*, vol. 10, Springer, 2018.
729 [2] E. ALLEN, J. BAGLAMA, AND S. BOYD, *Numerical approximation of the product of the square*
730 *root of a matrix with a vector*, Linear Algebra and its Applications, 310 (2000), pp. 167–181.
731 [3] B. D. ANDERSON AND J. B. MOORE, *Optimal filtering*, Courier Corporation, 2012.
732 [4] J. L. ANDERSON AND S. L. ANDERSON, *A monte carlo implementation of the nonlinear filtering*
733 *problem to produce ensemble assimilations and forecasts*, Monthly weather review, 127
734 (1999), pp. 2741–2758.
735 [5] M. ASCH, M. BOCQUET, AND M. NODET, *Data assimilation: methods, algorithms, and appli-*
736 *cations*, SIAM, 2016.
737 [6] C. M. BISHOP AND N. M. NASRABADI, *Pattern recognition and machine learning*, vol. 4,
738 Springer, 2006.
739 [7] P. BLANCHARD, D. J. HIGHAM, AND N. J. HIGHAM, *Accurately computing the log-sum-exp and*
740 *softmax functions*, IMA Journal of Numerical Analysis, 41 (2021), pp. 2311–2330.
741 [8] M. BOCQUET, J. BRAJARD, A. CARRASSI, AND L. BERTINO, *Bayesian inference of chaotic*
742 *dynamics by merging data assimilation, machine learning and expectation-maximization*,

- 743 arXiv preprint arXiv:2001.06270, (2020).
- 744 [9] L. BREIMAN, W. MEISEL, AND E. PURCELL, *Variable kernel estimates of multivariate densities*,
745 *Technometrics*, 19 (1977), pp. 135–144.
- 746 [10] G. BURGERS, P. J. VAN LEEUWEN, AND G. EVENSEN, *Analysis scheme in the ensemble Kalman*
747 *filter*, *Monthly weather review*, 126 (1998), pp. 1719–1724.
- 748 [11] X. CHEN, Z. J. WANG, AND M. J. MCKEOWN, *Shrinkage-to-tapering estimation of large co-*
749 *variance matrices*, *IEEE Transactions on Signal Processing*, 60 (2012), pp. 5640–5656.
- 750 [12] Y. CHEN, A. WIESEL, Y. C. ELДАР, AND A. O. HERO, *Shrinkage algorithms for mmse covari-*
751 *ance estimation*, *IEEE Transactions on Signal Processing*, 58 (2010), pp. 5016–5029.
- 752 [13] Y. CHEN, A. WIESEL, AND A. O. HERO, *Shrinkage estimation of high dimensional covariance*
753 *matrices*, in *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*,
754 *IEEE*, 2009, pp. 2937–2940.
- 755 [14] Y. CHEN, A. WIESEL, AND A. O. HERO, *Robust shrinkage estimation of high-dimensional*
756 *covariance matrices*, *IEEE Transactions on Signal Processing*, 59 (2011), pp. 4097–4107.
- 757 [15] E. CHOW AND Y. SAAD, *Preconditioned Krylov subspace methods for sampling multivariate*
758 *gaussian distributions*, *SIAM Journal on Scientific Computing*, 36 (2014), pp. A588–A608.
- 759 [16] J. R. DORMAND AND P. J. PRINCE, *A family of embedded runge-kutta formulae*, *Journal of*
760 *computational and applied mathematics*, 6 (1980), pp. 19–26.
- 761 [17] G. EVENSEN, *Sequential data assimilation with a nonlinear quasi-geostrophic model using*
762 *monte carlo methods to forecast error statistics*, *Journal of Geophysical Research: Oceans*,
763 99 (1994), pp. 10143–10162.
- 764 [18] M. FREI AND H. R. KÜNSCH, *Bridging the ensemble Kalman and particle filters*, *Biometrika*,
765 100 (2013), pp. 781–800.
- 766 [19] S. GÜTTEL, *Rational Krylov approximation of matrix functions: Numerical methods and opti-*
767 *mal pole selection*, *GAMM-Mitteilungen*, 36 (2013), pp. 8–31.
- 768 [20] M. W. HIRSCH, S. SMALE, AND R. L. DEVANEY, *Differential equations, dynamical systems, and*
769 *an introduction to chaos*, Academic press, 2012.
- 770 [21] P. JANSSEN, J. S. MARRON, N. VERAVERBEKE, AND W. SARLE, *Scale measures for band-*
771 *width selection*, *Journal of Nonparametric Statistics*, 5 (1995), pp. 359–380, <https://doi.org/10.1080/10485259508832654>, <https://doi.org/10.1080/10485259508832654>, <https://arxiv.org/abs/https://doi.org/10.1080/10485259508832654>.
- 772 [22] E. T. JAYNES, *Probability theory: The logic of science*, Cambridge university press, 2003.
- 773 [23] D. P. KINGMA AND J. BA, *Adam: A method for stochastic optimization*, arXiv preprint
774 arXiv:1412.6980, (2014).
- 775 [24] O. LEDOIT AND M. WOLF, *A well conditioned estimator for large dimensional covariance ma-*
776 *trices*, *Journal of Multivariate Analysis*, 88 (2004), pp. 365–411, [10.1016/S0047-259X\(03\)](https://doi.org/10.1016/S0047-259X(03)00096-4)
777 [00096-4](https://doi.org/10.1016/S0047-259X(03)00096-4).
- 778 [25] B. LIU, B. AIT-EL-FQUIH, AND I. HOTEIT, *Efficient kernel-based ensemble gaussian mixture*
779 *filtering*, *Monthly Weather Review*, 144 (2016), pp. 781–800.
- 780 [26] E. N. LORENZ, *Deterministic nonperiodic flow*, *Journal of atmospheric sciences*, 20 (1963),
781 pp. 130–141.
- 782 [27] E. N. LORENZ, *Predictability: A problem partly solved*, in *Proc. Seminar on predictability*,
783 vol. 1, 1996.
- 784 [28] R. M. NEAL AND G. E. HINTON, *A view of the em algorithm that justifies incremental, sparse,*
785 *and other variants*, in *Learning in graphical models*, Springer, 1998, pp. 355–368.
- 786 [29] E. D. NINO-RUIZ AND A. SANDU, *Ensemble Kalman filter implementations based on shrinkage*
787 *covariance matrix estimation*, *Ocean Dynamics*, 65 (2015), pp. 1423–1439.
- 788 [30] J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer, 1999.
- 789 [31] A. PAPOULIS AND S. PILLA, *Probability, random variables, and stochastic processes*, 583 pp,
790 1965.
- 791 [32] A. A. POPOV, C. MOU, A. SANDU, AND T. ILIESCU, *A multifidelity ensemble Kalman filter*
792 *with reduced order control variates*, *SIAM Journal on Scientific Computing*, 43 (2021),
793 pp. A1134–A1162.
- 794 [33] A. A. POPOV AND A. SANDU, *A bayesian approach to multivariate adaptive localization in*
795 *ensemble-based data assimilation with time-dependent extensions*, *Nonlinear Processes in*
796 *Geophysics*, 26 (2019), pp. 109–122.
- 797 [34] A. A. POPOV, A. SANDU, E. D. NINO-RUIZ, AND G. EVENSEN, *A stochastic covariance shrinkage*
798 *approach in ensemble transform Kalman filtering*, arXiv preprint arXiv:2003.00354, (2020).
- 799 [35] A. A. POPOV, A. N. SUBRAHMANYA, AND A. SANDU, *A stochastic covariance shrinkage approach*
800 *to particle rejuvenation in the ensemble transform particle filter*, *Nonlinear Processes in*
801 *Geophysics*, 29 (2022), pp. 241–253.
- 802 [36] M. L. PSIAKI, *Gaussian mixture nonlinear filtering with resampling for mixand narrowing*,
803
804

- 805 IEEE Transactions on Signal Processing, 64 (2016), pp. 5499–5512.
- 806 [37] S. REICH AND C. COTTER, *Probabilistic forecasting and Bayesian data assimilation*, Cambridge
807 University Press, 2015.
- 808 [38] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, The annals of mathematical
809 statistics, (1951), pp. 400–407.
- 810 [39] F. ROOSTA-KHORASANI AND M. W. MAHONEY, *Sub-sampled newton methods*, Mathematical
811 Programming, 174 (2019), pp. 293–326.
- 812 [40] B. W. SILVERMAN, *Density estimation for statistics and data analysis*, Routledge, 2018.
- 813 [41] H. W. SORENSON AND D. L. ALSPACH, *Recursive bayesian estimation using Gaussian sums*,
814 Automatica, 7 (1971), pp. 465–479.
- 815 [42] A. S. STORDAL, H. A. KARLSEN, G. NÆVDAL, H. J. SKAUG, AND B. VALLÈS, *Bridging the*
816 *ensemble Kalman filter and particle filters: the adaptive gaussian mixture filter*, Compu-
817 tational Geosciences, 15 (2011), pp. 293–305.
- 818 [43] P. J. VAN LEEUWEN, H. R. KÜNSCH, L. NERGER, R. POTTHAST, AND S. REICH, *Particle filters*
819 *for high-dimensional geoscience applications: A review*, Quarterly Journal of the Royal
820 Meteorological Society, 145 (2019), pp. 2335–2365.
- 821 [44] S. YUN, R. ZANETTI, AND B. A. JONES, *Kernel-based ensemble gaussian mixture filtering for*
822 *orbit determination with sparse data*, Advances in Space Research, 69 (2022), pp. 4179–
823 4197.