

1 EnsembleKalmanProcesses.jl: Derivative-free 2 ensemble-based model calibration

3 **Oliver R. A. Dunbar** ^{1*}[¶], **Ignacio Lopez-Gomez** ^{1*}, **Alfredo**
4 **Garbuno-Iñigo** ², **Daniel Zhengyu Huang** ¹, **Eviatar Bach** ¹, and **Jin-long**
5 **Wu** ³

6 ¹ Division of Geological and Planetary Sciences, California Institute of Technology ² Department of
7 Statistics, Mexico Autonomous Institute of Technology ³ Department of Mechanical Engineering,
8 University of Wisconsin-Madison ¶ Corresponding author * These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](#))

9 Summary

10 EnsembleKalmanProcesses.jl is a Julia-based toolbox that can be used for a broad class of
11 black-box gradient-free optimization problems. Specifically, the tools enable the optimization,
12 or calibration, of parameters within a computer model in order to best match user-defined
13 outputs of the model with available observed data ([Kennedy & O'Hagan, 2001](#)). Some of the
14 tools can also approximately quantify parametric uncertainty ([Huang, Huang, et al., 2022](#)).
15 Though the package is written in Julia, a read-write TOML-file interface is provided so that
16 the tools can be applied to computer models implemented in any language. Furthermore, the
17 calibration tools are non-intrusive, relying only on the ability of users to compute an output of
18 their model given a parameter value.

19 As the package name suggests, the tools are inspired by the well-established class of ensemble
Kalman methods. Ensemble Kalman filters are currently one of the only practical ways to
20 assimilate large volumes of observational data into models for operational weather forecasting
21 ([Evensen, 1994](#); [Houtekamer & Mitchell, 1998, 2001](#)). In the data assimilation setting, a
22 computational weather model is integrated for a short time over a collection, or ensemble,
23 of initial conditions, and the ensemble is updated frequently by a variety of atmospheric
24 observations, allowing the forecasts to keep track of the real system.

25
26 The workflow is similar for ensemble Kalman processes. Here, a computer code is run (in
27 parallel) for an ensemble of different values of the parameters that require calibration, producing
28 an ensemble of outputs. This ensemble of outputs is then compared to observed data, and
29 the parameters are updated to a new set of values which reduce the output-data misfit. The
30 computer model is then evaluated for the new ensemble values and the outputs, under certain
31 conditions, are guaranteed to lie closer to the observed data. The process is iterated until a
32 user-defined criterion of convergence is met. Optimal values are selected from statistics of the
33 final ensemble.

34 Features

35 There are different ensemble Kalman algorithms in the literature, which differ in the way that
36 the ensemble update is performed. The following ensemble Kalman processes are implemented
37 tools in our package, and we provide published references for detailed descriptions and evidence
38 of their efficacy:

- Ensemble Kalman Inversion (EKI, [Iglesias et al. \(2013\)](#)),

- 40 ▪ Ensemble Kalman Sampler (EKS, Garbuno-Inigo, Hoffmann, et al. (2020); Garbuno-Inigo,
41 Nüsken, et al. (2020)),
- 42 ▪ Unscented Kalman Inversion (UKI, Huang, Schneider, et al. (2022)),
- 43 ▪ Sparse Ensemble Kalman Inversion (SEKI, Schneider, Stuart, et al. (2022)).

44 We also implement some features to improve robustness and flexibility of the ensemble
45 algorithms:

- 46 ▪ The `ParameterDistribution` structure allows us to perform calibrations for parameters
47 with known constraints. It does so by defining transformation maps under-the-hood from
48 the constrained space to an unconstrained space where the optimization problem can be
49 suitably defined. Constrained optimization using this framework has been successfully
50 demonstrated in a variety of settings (Dunbar et al., 2022; Lopez-Gomez et al., 2022;
51 Schneider, Dunbar, et al., 2022).
- 52 ▪ The `FailureHandler` structure allows calibrations to continue when several ensemble
53 members fail. Common reasons for failure could be, for instance, simulation blow-up for
54 certain parameter configurations, user termination of slow computations, data corruption,
55 or bad nodes in a high-performance computing facility. This methodology is demonstrated
56 in Lopez-Gomez et al. (2022).
- 57 ▪ The `Localizer` structure allows us to overcome the restriction of the solution of the
58 calibration to the linear span of the initial ensemble, and to reduce sampling errors due
59 to the finite size of the ensemble. Various such localization and sampling error correction
60 methods are implemented in `EnsembleKalmanProcesses.jl` (Lee, 2021; Tong & Morzfeld,
61 2022).
- 62 ▪ The TOML-file interface defined in the `UQParameters` module allows non-intrusive
63 use of `EnsembleKalmanProcesses.jl` through TOML files, which are widely used for
64 configuration files and easily read in any programming language. Given the computer
65 model to calibrate and prior distributions on the parameters, `EnsembleKalmanProcesses.jl`
66 reads these distributions from a file and, after an iteration of the ensemble Kalman
67 algorithm, writes each member of the updated ensemble to a parameter file. Each
68 of these parameter files can be then read individually to initiate the ensemble of the
69 computer model for the next iteration.

70 Statement of need

71 The task of estimating parameters of a computer model or simulator such that its outputs
72 fit with data is ubiquitous in science and engineering, coming under many names such as
73 calibration, inverse problems, and parameter estimation. In statistics and machine learning,
74 when closed-form estimators of parameters of a model are unavailable, similar approaches may
75 need to be employed to fit the model to data. There is a wide variety of algorithms to suit
76 these applications; however, there are many bottlenecks in the practical application of such
77 methods to computer codes:

- 78 ▪ Legacy codes: Often code is old, and written in different languages than the packages
79 implementing the calibration algorithms, requiring elaborate interfaces.
- 80 ▪ Complex codes: Often large complex codes are difficult to change, so application of
81 intrusive calibration tools to models can be challenging.
- 82 ▪ Derivatives: When derivatives of a model output can be taken with respect to parameters,
83 they can often improve the rate of convergence. But in many practical cases, these
84 parameter-to-output maps are not differentiable; they may be chaotic or stochastic. Here
85 one should not – or cannot – apply gradient-based methods.

- 86 ▪ Lack of parallelism: There is now widespread access to high-performance computing
87 clusters, cloud computing, and local multi-threading, and such facilities should be
88 exploited where possible.

89 EnsembleKalmanProcesses.jl aims to provide a flexible and comprehensive solution to address
90 these challenges:

- 91 1. It is embarrassingly parallel with respect to the ensemble; therefore, all computer model
92 evaluations within an ensemble can happen simultaneously within an iteration.
- 93 2. It is derivative-free, and so is appropriate for computer codes for which derivatives are
94 not available. The optimal updates are robust to noise.
- 95 3. It is non-intrusive and so can be applied to black-box computer codes written in any
96 language or style, or to computer models for which the source code is not available to
97 the user.
- 98 4. With scalability enhancements, such as the ones provided by the Localizer structure, it
99 can be applied to high-dimensional problems.

100 Research projects using the package

- 101 ▪ EnsembleKalmanProcesses.jl has been used to train physics-based and machine-learning
102 models of atmospheric turbulence and convection, implemented using Flux.jl and
103 TurbulenceConvection.jl (Lopez-Gomez et al., 2022). In this application, the available
104 model outputs are not differentiable with respect to the learnable parameters, so gradient-
105 based optimization was not an option. In addition, the unscented Kalman inversion
106 algorithm was used to approximately quantify parameter uncertainty.
- 107 ▪ EnsembleKalmanProcesses.jl features within Calibrate-Emulate-Sample (CES, Cleary et
108 al. (2021)), a pipeline used to accelerate parameter uncertainty quantification (by a
109 factor of 10^3 - 10^4 with respect to Monte Carlo methods) by using statistical emulators.
110 EnsembleKalmanProcesses.jl is used to choose training points for these emulators. The
111 training points are naturally concentrated by the ensemble Kalman processes into areas
112 of high posterior probability mass. Within CES, the trained emulators are used to sample
113 this probability distribution, and by design are most accurate where they need to be. CES
114 has been successfully used to quantify parameter uncertainty within the moist convection
115 scheme of a simplified climate model (Dunbar et al., 2021, 2022; Howland et al., 2022),
116 within a droplet collision-coalescence scheme for cloud microphysics (Bieli et al., 2022),
117 and within boundary layer turbulence schemes for ocean modeling (Hillier, 2022).
- 118 ▪ EnsembleKalmanProcesses.jl has been used to learn hyperparameters within a machine
119 learning tool known as Random Features within a julia package RandomFeatures.jl.
120 Here, the hyperparameters characterize an infinite family of functions, from which a
121 finite sample is drawn to use as a basis in regression problems. The objective for learning
122 the parameters is noisy and non-differentiable due to the random sampling, so ensemble
123 Kalman processes naturally perform well in this setting.

124 Acknowledgements

125 We acknowledge contributions from several others who played a role in the evolution of this
126 package. These include Jake Bolewski, Navid Constantinou, Gregory L. Wagner, Thomas
127 Jackson, Michael Howland, Melanie Bieli, and Adeline Hillier. The development of this package
128 was supported by the generosity of Eric and Wendy Schmidt by recommendation of the
129 Schmidt Futures program, and by the Defense Advanced Research Projects Agency (Agreement
130 No. HR00112290030).

131 **References**

- 132 Bieli, M., Dunbar, O. R. A., Jong, E. K. de, Jaruga, A., Schneider, T., & Bischoff, T. (2022).
133 An efficient Bayesian approach to learning droplet collision kernels: Proof of concept using
134 “Cloudy,” a new n-moment bulk microphysics scheme. *Journal of Advances in Modeling*
135 *Earth Systems*, 14(8), e2022MS002994. <https://doi.org/10.1029/2022MS002994>
- 136 Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2021). Calibrate,
137 emulate, sample. *Journal of Computational Physics*, 424, 109716. <https://doi.org/10.1016/j.jcp.2020.109716>
- 138
- 139 Dunbar, O. R. A., Garbuno-Inigo, A., Schneider, T., & Stuart, A. M. (2021). Calibration
140 and uncertainty quantification of convective parameters in an idealized GCM. *Journal of*
141 *Advances in Modeling Earth Systems*, 13(9), e2020MS002454. <https://doi.org/10.1029/2020MS002454>
- 142
- 143 Dunbar, O. R. A., Howland, M. F., Schneider, T., & Stuart, A. M. (2022). Ensemble-based
144 experimental design for targeting data acquisition to inform climate models. *Journal of*
145 *Advances in Modeling Earth Systems*, 14(9), e2022MS002997. <https://doi.org/10.1029/2022MS002997>
- 146
- 147 Evensen, G. (1994). Sequential data assimilation with a nonlinear quasi-geostrophic model
148 using Monte Carlo methods to forecast error statistics. *Journal of Geophysical Research: Oceans*, 99, 10143–10162. <https://doi.org/10.1029/94JC00572>
- 149
- 150 Garbuno-Inigo, A., Hoffmann, F., Li, W., & Stuart, A. M. (2020). Interacting Langevin
151 diffusions: Gradient structure and ensemble Kalman sampler. *SIAM Journal on Applied*
152 *Dynamical Systems*, 19(1), 412–441. <https://doi.org/10.1137/19M1251655>
- 153 Garbuno-Inigo, A., Nüsken, N., & Reich, S. (2020). Affine invariant interacting Langevin
154 dynamics for Bayesian inference. *SIAM Journal on Applied Dynamical Systems*, 19(3),
155 1633–1658. <https://doi.org/10.1137/19M1304891>
- 156 Hillier, A. (2022). *Supervised calibration and uncertainty quantification of subgrid closure*
157 *parameters using ensemble Kalman inversion* [Master’s thesis, Massachusetts Institute of
158 Technology. Department of Electrical Engineering; Computer Science]. <https://hdl.handle.net/1721.1/145140>
- 159
- 160 Houtekamer, P. L., & Mitchell, H. L. (1998). Data assimilation using an ensemble Kalman
161 filter technique. *Monthly Weather Review*, 126, 796–811. [https://doi.org/10.1175/1520-0493\(1998\)126%3C0796:DAJAEK%3E2.0.CO;2](https://doi.org/10.1175/1520-0493(1998)126%3C0796:DAJAEK%3E2.0.CO;2)
- 162
- 163 Houtekamer, P. L., & Mitchell, H. L. (2001). A sequential ensemble Kalman filter for
164 atmospheric data assimilation. *Monthly Weather Review*, 129, 123–137. [https://doi.org/10.1175/1520-0493\(2001\)129%3C0123:ASEKFF%3E2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129%3C0123:ASEKFF%3E2.0.CO;2)
- 165
- 166 Howland, M. F., Dunbar, O. R. A., & Schneider, T. (2022). Parameter uncertainty quantifica-
167 tion in an idealized GCM with a seasonal cycle. *Journal of Advances in Modeling Earth*
168 *Systems*, 14(3), e2021MS002735. <https://doi.org/10.1029/2021MS002735>
- 169 Huang, D. Z., Huang, J., Reich, S., & Stuart, A. M. (2022). Efficient derivative-free
170 Bayesian inference for large-scale inverse problems. *arXiv Preprint arXiv:2204.04386*.
171 <https://doi.org/10.48550/arXiv.2204.04386>
- 172 Huang, D. Z., Schneider, T., & Stuart, A. M. (2022). Iterated Kalman methodology for inverse
173 problems. *Journal of Computational Physics*, 463, 111262. <https://doi.org/10.1016/j.jcp.2022.111262>
- 174
- 175 Iglesias, M. A., Law, K. J., & Stuart, A. M. (2013). Ensemble Kalman methods for inverse
176 problems. *Inverse Problems*, 29(4), 045001. <https://doi.org/10.1088/0266-5611/29/4/045001>
- 177

- 178 Kennedy, M., & O'Hagan, A. (2001). Bayesian calibration of computer models. *Journal of the*
179 *Royal Statistical Society Series B*, 63, 425–464. <https://doi.org/10.1111/1467-9868.00294>
- 180 Lee, Y. (2021). *Sampling error correction in ensemble Kalman inversion*. [https://doi.org/10.](https://doi.org/10.48550/arxiv.2105.11341)
181 [48550/arxiv.2105.11341](https://doi.org/10.48550/arxiv.2105.11341)
- 182 Lopez-Gomez, I., Christopoulos, C., Langeland Ervik, H. L., Dunbar, O. R. A., Cohen, Y.,
183 & Schneider, T. (2022). Training physics-based machine-learning parameterizations with
184 gradient-free ensemble Kalman methods. *Journal of Advances in Modeling Earth Systems*,
185 14(8), e2022MS003105. <https://doi.org/10.1029/2022MS003105>
- 186 Schneider, T., Dunbar, O. R. A., Wu, J., Böttcher, L., Burov, D., Garbuno-Inigo, A., Wagner,
187 G. L., Pei, S., Daraio, C., Ferrari, R., & Shaman, J. (2022). Epidemic management and
188 control through risk-dependent individual contact interventions. *PLOS Computational*
189 *Biology*, 18(6), e1010171. <https://doi.org/10.1371/journal.pcbi.1010171>
- 190 Schneider, T., Stuart, A. M., & Wu, J.-L. (2022). Ensemble Kalman inversion for sparse
191 learning of dynamical systems from time-averaged data. *Journal of Computational Physics*,
192 111559. <https://doi.org/10.1016/j.jcp.2022.111559>
- 193 Tong, X. T., & Morzfeld, M. (2022). *Localization in ensemble Kalman inversion*. <https://doi.org/10.48550/arXiv.2201.10821>
194

DRAFT