

# Iterating marginalized Bayes maps for likelihood maximization with application to nonlinear panel models

Jesse Wheeler\*

Department of Mathematics and Statistics, Idaho State University

Aaron J. Abkemeier

Department of Statistics, University of Michigan

Edward L. Ionides

Department of Statistics, University of Michigan

November 24, 2025

## Abstract

Complex dynamic systems can be investigated by fitting mechanistic stochastic dynamic models to time series data. In this context, commonly used Monte Carlo inference procedures for model selection and parameter estimation quickly become computationally unfeasible as the system dimension grows. The increasing prevalence of panel data, characterized by multiple related time series, therefore necessitates the development of inference algorithms that are effective for this class of high-dimensional mechanistic models. Nonlinear, non-Gaussian mechanistic models are routinely fitted to time series data but seldom to panel data, despite its widespread availability, suggesting that the practical difficulties for existing procedures are prohibitive. We investigate the use of iterated filtering algorithms for this purpose. We introduce a novel algorithm that contains a marginalization step that mitigates issues arising from particle filtering in high dimensions. Our approach enables likelihood-based inference for models that were previously considered intractable, thus broadening the scope of dynamic models available for panel data analysis.

*Keywords:* Mechanistic Models, Nonlinear Dynamics, Particle Filters, High Dimensional Inference, Longitudinal data

---

\*Email: jessewheeler@isu.edu, ORCID: 0000-0003-3941-3884

# 1 Introduction

Panel data, otherwise known as longitudinal data, are a collection of related time series. Each time series is a measurement on a person, animal or object known as the *unit*. A unit may have its treatment assigned via a randomized experiment, or it may be measured in an observational study. The measurement on each unit at each observation time may be vector-valued. While each time series within a panel dataset could be analyzed individually, there is an advantage to analyzing the entire collection of time series simultaneously. For example, data from multiple measurement units during an infectious disease outbreak may reveal transmission dynamics not evident from individual units (Wiens et al., 2014; Ranjeva et al., 2019). Other examples include ecological experiments and observational studies in which data are collected over time across multiple measurement sites that may experience variation in covariates of interest (Searle et al., 2016; Hewitt et al., 2024). Modeling the data collectively allows for researchers to learn about processes that are shared across units, as well as identifying traits that are unique to each unit.

A common approach to modeling nonlinear dynamic systems is through the use of mechanistic models (Auger-Méthé et al., 2021). These models involve the proposal of a system of equations that describe how unobserved dynamic states evolve over time. When combined with a model relating the latent states to observable quantities, we obtain a partially observed Markov process (POMP) model, also known as a state space model (SSM) or a hidden Markov model. Once calibrated to data, these models provide a quantitative description of the observations while simultaneously adhering to a given scientific hypothesis about how the data are generated. Despite the growing ability of deep learning and other advanced machine learning techniques to extract useful information from complex data sets, mechanistic models continue to be an important tool for modern science (Baker et al.,

2018; Hogg and Villar, 2024). In addition to improved interpretability, a key advantage of these models is that they enable the estimation of counter-factual scenarios, such as impact of interventions on a dynamic system (Wheeler et al., 2024).

Various algorithms, both frequentist and Bayesian, are capable in principle of fitting mechanistic models to panel data. Some algorithms scale well as the model’s dimension grows, but they rely on unrealistic approximations (Evensen, 2009; Wigren and Lindsten, 2022), or they avoid the difficulties related to evaluating the model’s likelihood function by optimizing alternative measures of goodness-of-fit or approximations to the likelihood (Toni et al., 2009; Wood, 2010; Whitehouse et al., 2023; Häggström et al., 2025). Other algorithms can be applied in less restrictive cases, but they scale poorly as the model’s dimension increases (Andrieu et al., 2010; Ionides et al., 2015).

Contemporary applications of nonlinear mechanistic models for low-dimensional systems are abundant (e.g., Kramer et al., 2024; He et al., 2024; Newman et al., 2023), yet examples of higher-dimensional panel equivalents remain sparse. The lack of examples indicates that scientists have struggled, largely unsuccessfully, to use existing statistical methodology for nonlinear panel models. Our work addresses this possibility by proposing a new algorithm that is effective on benchmark problems, theoretically supported, and easy to implement using existing software such as the R package `panelPomp` (Bretó et al., 2025), or a newly developed python package called `pypomp`, available on GitHub and PyPI (Abkemeier et al., 2025).

An extensively used approach in ecology, known as *data cloning* (Lele et al., 2007), involves iteratively using Bayes’ rule, by recursively mapping a prior distribution to a posterior distribution using the same likelihood function, until the iterated posterior distribution converges to a point mass at the maximum likelihood estimate (MLE). Ionides et al.

(2015) developed an extension of data cloning for SSMS by treating the model parameters as latent states, and performing a random walk for these parameters at each observation time. By iteratively applying a particle filter (Arulampalam et al., 2002) to this extended model and decreasing the random walk standard deviations over time, it can be shown that the parameters will converge to the MLE after a sufficiently large number of iterations (Ionides et al., 2015; Chen et al., 2025). This algorithm, known as IF2, has seen extensive use for inference on low dimensional dynamic models, especially in epidemiological contexts (e.g., Pons-Salort and Grassly, 2018; Stocks et al., 2018; Subramanian et al., 2021; Fox et al., 2022).

In panel models, the iterated filtering algorithm must be extended to handle the Monte Carlo error that grows exponentially with the number of units. To address this, Bretó et al. (2020) proposed an algorithm called the panel iterated filter (PIF). As discussed in Appendix B, this algorithm is a special case of iterated filtering where the model structure and random walk sequence have been modified to reduce the loss of information that is described by Liu and West (2001). The PIF algorithm has been effective in obtaining maximum likelihood estimates for highly nonlinear, non-Gaussian panel models in previous works (Ranjeva et al., 2017, 2019; Wale et al., 2019; Lee et al., 2020; Domeyer et al., 2022). However, as we demonstrate here, the PIF algorithm can be inefficient when the number of units is large, due to the resampling of all parameters at each step of the iterated particle filters. In order to improve the computational inefficiencies of the PIF algorithm for high-dimensional panel models, we present a novel inference technique which we call the marginalized panel iterated filter (MPIF). The algorithm reduces the number of times particles are resampled, which results in increased diversity of the particles that represent the parameter distribution at each iteration and time step.

The article proceeds by introducing PanelPOMP models, the mathematical framework we use to discuss mechanistic models for panel data; elsewhere, these models have also been called multi-SSMs (Wigren and Lindsten, 2022). We then establish theoretical guarantees for the convergence of iterated filtering algorithms for these models (Theorem 1). Our theoretical framework extends the analysis by Chen et al. (2025) to panel data while requiring weaker conditions than those by Bretó et al. (2020). We then introduce the marginalization step of MPIF that facilitates high dimensional inference by describing a concept we call marginalized Bayes maps. The marginalization step introduces a nonlinearity in the Bayes map that existing theoretical approaches cannot immediately address. We provide theoretical guarantees of the MPIF algorithm for some special cases, and demonstrate that the algorithm remains effective in a more general setting via a simulation study. We then conduct a data analysis of a high-dimensional dataset of pre-vaccination measles case reports from 20 towns in the United Kingdom (UK).

## 2 Iterated filtering for panel models

A POMP model comprises an unobservable Markov process  $\{X(t), t \in \mathcal{T} \subset \mathbb{R}\}$  and an observable sequence  $Y_1, \dots, Y_N$ . We suppose that  $Y_n$  is a measurement of  $X(t_n)$ , with  $t_1 < \dots < t_N$ , formalized by a requirement that  $Y_n$  is conditionally independent of  $\{X(s), Y_k : k \neq n, s \neq t_n\}$  given  $X(t_n)$ . We assume that  $X(t)$  takes values in  $\mathcal{X} \subset \mathbb{R}^{d_x}$ , and  $Y_n$  takes values in  $\mathcal{Y} \subset \mathbb{R}^{d_y}$ . Data collected at time  $t_n$ , denoted by  $y_n^*$ , are modeled as a realization of  $Y_n$ .

While  $X(t)$  may be a continuous or discrete time process, the value of  $X(t)$  at observation times is of particular interest, and so we write  $X_n = X(t_n)$ . Initial values may be specified at a time  $t_0 < t_1$ , and we set  $X_0 = X(t_0)$ . We adopt the notation that for any

integers  $a$  and  $b$ ,  $a : b$  is the vector  $(a, a + 1, \dots, b - 1, b)$ , and use the convention that  $a : b = \emptyset$  if  $b < a$ . Similarly, we use the notation in subscripts to denote collections of random variables, such as  $X_{0:N} = (X_0, \dots, X_N)$ , and use the same basic notation for  $Y_{1:N}$ , and  $t_{0:N}$ . We assume that the joint probability density  $f_{X_{0:N}, Y_{1:N}}(x_{0:N}, y_{1:N}; \theta)$  exists, with respect to Lebesgue measure, or a counting measure if the set of possible values is discrete.

In a panel data analysis scenario, data are collected for each unit  $u \in 1 : U$ , where  $U$  is the number of units in the panel. Because the data generating process of each unit is assumed to be dynamically independent, we may model the panel data as a collection of related POMP models, which we call a PanelPOMP. To distinguish between each unit POMP model, we denote the measurement and latent process for each unit with subscript  $u \in 1 : U$ . Specifically, we write  $X_{u,n}, Y_{u,n}$  to denote the values of the latent process  $\{X_u(t), t \in \mathcal{T}_u \subset \mathbb{R}\}$  and measurement sequence for unit  $u$  at the measurement times  $t_{u,1} < t_{u,2} < \dots < t_{u,N_u}$ . The number of observations  $N_u$ , observations times  $t_{u,n}$ , and time domain  $\mathcal{T}_u$  do not need to match across units, though this is often the case. We use the shorthand  $\mathbf{X}, \mathbf{Y}$  to refer to the collection of all latent and observable process at observation times for all units in the panel. Due to the independence of units, the joint density of the PanelPOMP model can be written as

$$f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \theta) = \prod_{u=1}^U f_{X_{u,0:N_u}, Y_{u,1:N_u}}(x_{u,0:N_u}, y_{u,1:N_u}; \theta) \quad (1)$$

$$= \prod_{u=1}^U f_{X_{u,0}}(x_{u,0}; \theta) \prod_{n=1}^{N_u} f_{Y_{u,n}|X_{u,n}}(y_{u,n}|x_{u,n}; \theta) f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|x_{u,n-1}; \theta), \quad (2)$$

where Eq. 1 arises from the assumption of dynamically independent units, and Eq. 2 is a result of the Markov property of  $X_u(t)$  and the conditional independence of  $Y_{u,n}$ .

Our goal is to make inferences about the parameter  $\theta$  by maximizing the likelihood

function

$$L(\theta; \mathbf{y}^*) = \int f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}^*; \theta) d\mathbf{x}, \quad (3)$$

where  $\mathbf{y}^*$  denotes the collection observations from the entire panel. Although the unit densities in Eq. 1 can be factored as a result of the dynamic independence between measurement units, a defining feature of a panel model is that the parameter vector  $\theta$  remains relevant across all the dynamic systems. Therefore, inference methodology used to obtain estimates of  $\theta$  should use data from each unit in the PanelPOMP. A special case of particular interest arises when  $\theta = (\phi, \psi_{1:U})$ , where  $\phi \in \Theta_\phi \subset \mathbb{R}^{d_\phi}$  is a vector of parameters that are shared across each unit, and  $\psi_u \in \Theta_\psi \subset \mathbb{R}^{d_\psi}$  are parameters that are specific to unit  $u$ , formally,  $f_{X_{u,0:N_u}, Y_{u,1:N_u}}(x_{u,0:N_u}, y_{u,1:N}; \theta) = f_{X_{u,0:N_u}, Y_{u,1:N_u}}(x_{u,0:N_u}, y_{u,1:N}; \phi, \psi_u)$ .

Algorithm 1 describes an iterated filtering algorithm for PanelPOMP models that obtains the MLE for both shared ( $\phi$ ) and unit specific ( $\psi_u$ ) parameters. With `MARGINALIZE = FALSE`, this is the PIF algorithm of Bretó et al. (2020). Our innovation occurs when `MARGINALIZE = TRUE`, and we call this the marginalized panel iterated filter (MPIF). This small difference requires new approaches to theoretical analysis, and it is not immediately clear if and when this sequence of approximations may converge to the exact MLE. However, we demonstrate that this modification has dramatic consequences for scalability.

Mathematically, MPIF adds an additional step to PIF by marginalizing out the unit-specific parameters that are irrelevant for the unit currently being filtered. Because the parameter distributions are represented via Monte Carlo samples, marginalization is carried out by decoupling elements of the particles vectors representing parameters  $\psi_{-u} = \{\psi_k\}_{k \neq u}$  from the elements of the vector representing  $(\phi, \psi_u)$ . This decoupling can occur by not updating the particles representing  $\psi_{-u}$  when using weights obtained from data in unit  $u$ .

In other words, instead of resampling all parameter particles in line 13 of Algorithm 1, we only resample particles for the shared parameters  $\phi$  and the unit-specific parameters  $\psi_u$  related to the data of the specific unit under consideration.

---

**Algorithm 1: Iterated Filtering for PanelPOMP models**

---

**Inputs:**

Simulator of initial density,  $f_{X_{u,0}}(x_{u,0}; \theta)$  for  $u$  in  $1:U$ .  
 Simulator of transition density,  $f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|x_{u,n-1}; \theta)$  for  $u$  in  $1:U$ ,  $n$  in  $1:N_u$ .  
 Evaluator of measurement density,  $f_{Y_{u,n}|X_{u,n}}(y_{u,n}|x_{u,n}; \theta)$  for  $u$  in  $1:U$ ,  $n$  in  $1:N_u$ .  
 Data  $y_{u,n}^*$ , for  $u$  in  $1:U$ ,  $n$  in  $1:N_u$ .  
 Number of iterations,  $M$ .  
 Number of particles,  $J$ .  
 Starting parameter swarm,  $\Theta_j^0 = (\Phi_j^0, \Psi_{1:U,j}^0)$  for  $j \in 1:J$ ,  $u \in 1:U$ .  
 Simulator of perturbation densities,  $h_{u,n}(\cdot|\varphi; \sigma)$  for  $m \in 1:M$ ,  $u \in 1:U$ ,  $n \in 0:N_u$ .  
 Perturbation Sequence  $\sigma_{1:U,1:M}$ .  
 Logical variable determining marginalization, MARGINALIZE.

**Output:**

Final parameter swarm,  $\Theta_j^m = (\Phi_j^m, \Psi_{1:U,j}^m)$  for  $j \in 1:J$ ,  $u \in 1:U$ .

---

```

1 for  $m \in 1:M$  do
2   Set  $\Theta_{0,j}^{F,m} = \Theta_j^{m-1} = (\Phi_j^{m-1}, \Psi_{1:U,j}^{m-1})$  for  $j \in 1:J$ ;
3   for  $u \in 1:U$  do
4     Set  $\Theta_{u,0,j}^{F,m} = (\Phi_{u,0,j}^{F,m}, \Psi_{1:U,0,j}^{F,m}) \sim h_{u,0}(\cdot | \Theta_{u-1,j}^{F,m}; \sigma_{u,m})$  ;
5     Initialize  $X_{u,0,j}^{F,m} \sim f_{X_{u,0}}(x_{u,0}; \Phi_{u,0,j}^{F,m}, \Psi_{u,0,j}^{F,m})$  for  $j \in 1:J$ ;
6     for  $n \in 1:N_u$  do
7       Set  $\Theta_{u,n,j}^{P,m} = (\Phi_{u,n,j}^{P,m}, \Psi_{1:U,n,j}^{P,m}) \sim h_{u,n}(\cdot | \Theta_{u,n-1,j}^{F,m}; \sigma_{u,m})$  for  $j \in 1:J$  ;
8        $X_{u,n,j}^{P,m} \sim f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|X_{u,n-1,j}^{F,m}; \Phi_{u,n,j}^{P,m}, \Psi_{u,n,j}^{P,m})$  for  $j \in 1:J$  ;
9        $w_{u,n,j}^m = f_{Y_{u,n}|X_{u,n}}(y_{u,n}^* | X_{u,n,j}^{P,m}; \Phi_{u,n,j}^{P,m}, \Psi_{u,n,j}^{P,m})$  for  $j \in 1:J$  ;
10      Draw  $k_{1j}$  with  $P(k_j = i) = w_{u,n,i}^m / \sum_{v=1}^J w_{u,n,v}^m$  for  $i, j \in 1:J$ ;
11      Set  $X_{u,n,j}^{F,m} = X_{u,n,k_j}^{P,m}$ , and  $(\Phi_{u,n,j}^{F,m}, \Psi_{u,n,j}^{F,m}) = (\Phi_{u,n,k_j}^{P,m}, \Psi_{u,n,k_j}^{P,m})$  for  $j \in 1:J$ ;
12      if MARGINALIZE then
13         $\Psi_{\tilde{u},n,j}^{F,m} = \Psi_{\tilde{u},n,j}^{P,m}$  for all  $\tilde{u} \neq u$ ,  $j = 1:J$ 
14      else
15         $\Psi_{\tilde{u},n,j}^{F,m} = \Psi_{\tilde{u},n,k_j}^{P,m}$  for all  $\tilde{u} \neq u$ ,  $j = 1:J$ 
16      end
17    end
18    Set  $\Theta_{u,j}^{F,m} = (\Phi_{u,N_u,j}^{F,m}, \Psi_{u,N_u,j}^{F,m})$  for  $j \in 1:J$  ;
19  end
20  Set  $\Theta_j^{(m)} = \Theta_{U,j}^{F,m}$  for  $j \in 1:J$ ;
21 end

```

---

The algorithmic complexity of both PIF and MPIF is  $O(JMNU)$ , where  $M$  represents



the number of iterations,  $J$  is the number of particles,  $U$  the number of units in the panel, and  $N$  is the mean of  $\{N_1, \dots, N_U\}$ . Because MPIF does not require tracking the history of each particle, we achieve a minor reduction in computational overhead associated with the PIF algorithm. However, as discussed in Section 4.1, the primary advantage of MPIF is that it typically exhibits lower Monte Carlo uncertainty (requiring smaller  $J$ ) and generally converges in fewer iterations (requiring smaller  $M$ ).

Central to the success of iterated filtering algorithms is the fact that the repeated use of the posterior distribution from one Bayesian update as the prior distribution for the next iteration ultimately leads to a degenerate distribution centered at the MLE (see Section 3). This idea is key to deriving the proof of Theorem 1, which extends existing theory (Chen et al., 2025) for the convergence of the PIF algorithm.

**Theorem 1.** *Consider a PanelPOMP model defined by Eq. 1, and let  $\Theta \subset \mathbb{R}^{d_\phi + U d_\psi}$  be a compact set that satisfies condition (A1) in Appendix A, and assume there exists a  $\delta > 0$  such that  $\{\theta \in \Theta : |\theta - \hat{\theta}|_2 < \delta\} \subseteq \Theta$ , where  $\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta; \mathbf{y}^*)$ . Denote the output of Algorithm 1 without marginalization as  $\Theta_{1:J}^{(M)}$ , and assume that the model satisfies conditions (B1)–(B3), and the sequence of perturbation defined by  $h_{u,n}(\cdot | \varphi, \sigma_{1:U, 1:M})$  satisfy conditions (C1)–(C4). Then there exists some positive sequences  $\{C_M\}_{M \geq 1}$  and  $\{\epsilon_M\}_{M \geq 1}$  where  $\lim_{M \rightarrow \infty} \epsilon_M = 0$  such that for all  $(J, M) \in \mathbb{N}^2$ ,*

$$E \left[ \left\| \frac{1}{J} \sum_{i=1}^J \Theta_i^{(M)} - \hat{\theta} \right\|_2 \right] \leq \frac{C_M}{\sqrt{J}} + \epsilon_M$$

*Proof outline.* This theorem is an application of Theorem 4 of Chen et al. (2025) to PanelPOMP models, following the approach of Bretó et al. (2020), who showed that PanelPOMP models can be expressed as a lower-dimensional POMP model. A full proof is provided in Appendix B. □

Theorem 1 formally provides guarantees for a variety of iterated filtering algorithms for PanelPOMP models, but as we demonstrate in Section 4, the applicability of these algorithms often have scalability issues. The marginalization step in Algorithm 1 mitigates these issues but introduces a nonlinearity that invalidates the data cloning principle that is key to previous theoretical work on iterated filtering algorithms. In the following section, we briefly discuss data cloning and its relationship to the MPIF algorithm.

### 3 Iterating marginalized Bayes maps

If we denote  $\pi_i(\theta)$  as the posterior distribution of the parameter vector  $\theta$  after the  $i$ th Bayesian update, and  $\mathbf{y}^*$  as the observed data, we can express an iterated Bayesian update as the following:

$$\begin{aligned}\pi_1(\theta) &\propto f(\mathbf{y}^*; \theta) \pi_0(\theta), \\ \pi_2(\theta) &\propto f(\mathbf{y}^*; \theta) \pi_1(\theta) \propto f^2(\mathbf{y}^*; \theta) \pi_0(\theta), \\ &\vdots \\ \pi_m(\theta) &\propto f^m(\mathbf{y}^*; \theta) \pi_0(\theta).\end{aligned}$$

In this representation,  $f(\mathbf{y}^*; \theta)$  is the likelihood function (Eq. 3), and  $\pi_0(\theta)$  is the original prior distribution for  $\theta$ . If we let  $m \rightarrow \infty$ , the effect of the initial prior distribution diminishes, and the  $m$ th posterior has all of its mass centered at the MLE. This can be shown by taking the limit of  $\pi_m(\theta)/\pi_m(\hat{\theta})$  as  $m$  goes to infinity: if  $\theta = \hat{\theta}$ , then the limit is one, and zero otherwise (Lele et al., 2007).

The data cloning algorithm is useful for estimating the MLE in situations where the likelihood function is known or readily evaluated up to a constant of proportionality. In this scenario, practitioners can leverage existing Bayesian software in order to obtain a

maximum likelihood estimate (MLE) (Auger-Méthé et al., 2021; Ponciano et al., 2009). However, for nonlinear non-Gaussian state-space models, the likelihood function is generally inaccessible (Häggström et al., 2025). In these cases, simulation-based inference techniques such as the particle filter (Arulampalam et al., 2002) can reliably approximate the likelihood function. An iterated Bayes procedure could in theory be used in conjunction with a particle filter to estimate the MLE. However, a well-known issue with particle filters is the difficulty in accurately sampling the posterior distribution of fixed model parameters due to particle depletion (Liu and West, 2001). Consequently, a direct application of the iterated Bayes approach using particle filters for MLE estimation is impractical.

Iterated filtering algorithms overcome the issues associated with particle depletion by introducing a random walk for model parameters, thereby rescuing the degenerate particle representation of model parameters. Early analysis of this procedure showed that if the random walk standard deviations are small, then this modification still leads to an approximation of the iterated Bayes algorithm outlined above, where the final particle mass is still centered at the MLE (Ionides et al., 2015). The introduced perturbations of parameter values are necessary for inference, but they also introduce a loss of information (Liu and West, 2001). Therefore, in practice, the random perturbations are reduced as the number of iterations increases. Recent theoretical analysis demonstrates that the algorithm still concentrates on the MLE when perturbations are reduced over time (Chen et al., 2025).

Because PanelPOMP models are a special case of POMP models, this same approach can theoretically be used to estimate the MLE. However, the particle filter famously suffers from the curse-of-dimensionality. That is, the approximation error of the particle filter grows exponentially with the number of units in a panel model (Bengtsson et al., 2008; Snyder et al., 2008). The panel iterated filter (PIF) of Bretó et al. (2020) partially addresses

this issue by stacking the individual time series into a single long time series, and modifying the perturbation kernel for parameter particles to mitigate loss of information, as discussed in Section 6. Importantly, the PIF algorithm still approximates the iterated Bayes map when the random walk standard deviations are small (Theorem 1 of Bretó et al. (2020)), and can be shown to converge when the perturbations shrink over time (our Theorem 1).

The additional marginalization step in the MPIF algorithm introduces a nonlinear transformation at each iteration, changing the distribution that is approximated by the algorithm. As such, existing theoretic approaches for iterated filtering algorithms (Chen et al., 2025) for this class of models are insufficient to demonstrate convergence of the algorithm. We explore why this is the case in the context of *marginalized data cloning*. Let  $\theta = (\phi, \psi_{1:U})$  denote the parameter vector for the panel model, where  $\phi$  denotes the parameters that are shared by all  $U$  units, and  $\psi_u$  are the parameters that are only relevant to unit  $u$ . Similarly, we write  $y_u^* = y_{u,1:N_u}^*$  to denote the time series data for unit  $u$ . By stacking the times series into a single time series and iteratively filtering one at a time and ignoring parameter perturbations, the Bayes map that is approximated by a single sub-iteration of PIF can be written sequentially as:

$$\pi_{m,u}(\theta) \propto f_u(y_u^*; \theta) \pi_{m,u-1}(\theta) = f_u(y_u^*; \phi, \psi_u) \pi_{m,u-1}(\theta), \quad (4)$$

where  $\pi_{m,u}(\theta)$  is the parameter distribution at step  $(m, u)$ , and we adopt the convention that  $\pi_{0,0}(\theta) = \pi_0(\theta)$  is the initial prior density and  $\pi_{m,U} = \pi_{m+1,0} = \pi_{m+1}$ . This update is completed for each unit  $u \in 1:U$  which we call *unit* iterations, then for iterated for  $m \in 1:M$ , which we call *complete* or *full* iterations.

An important observation regarding the representation in Eq. 4 is that each unit-specific likelihood function  $f_u(y_u^*; \phi, \psi_u)$  contributes directly to the information about the shared

parameter vector  $\phi$  and its respective unit-specific parameter vector  $\psi_u$ , but not that of the unit-specific parameters  $\psi_{-u}$ . Despite this, the Bayes update in the PIF algorithm necessitates updating the posterior distribution of all parameters at each unit-iteration. Because the parameter distributions are represented via Monte Carlo samples (called particles), this implies that the PIF algorithm re-weights particles representing  $\psi_{-u}$  based on a likelihood that does not contain direct information about the parameters. For example, within a single  $m$ -iteration, the particles representing the distribution of  $\psi_U$  have been resampled  $\sum_{u=1}^{U-1} N_u$  times before encountering the data  $y_U^*$ , the only subset of data containing direct information about  $\psi_U$ . This process can lead to significant particle depletion (see Figure 2), particularly if  $U$  or  $N_u$  are large.

If the prior distribution used in Eq. 4 is independent across parameters, then the posterior distribution of the sub-vector  $\psi_{-u}$  will be unchanged. In this case, the particles representing the density for these parameters do not need to be resampled, which would avoid the issue of particle depletion. The use of independent priors is common practice in Bayesian statistics, but each complete iteration of Eq. 4 introduces parameter dependence via the likelihood function. This observation leads to the proposal of the marginalized PIF algorithm (MPIF), where the intermediate posterior distributions are made independent by marginalization before use as a prior distribution in the subsequent unit-iteration. A representation of a unit-iteration following this approach is given in Eqs. 5 and 6.

$$\tilde{\pi}_{m,u}(\theta) \propto f_u(y_u^*; \phi, \psi_u) \pi_{m,u-1}(\theta) \quad (5)$$

$$\pi_{m,u}(\theta) \propto \int \tilde{\pi}_{m,u}(\theta) d\phi d\psi_u \times \int \tilde{\pi}_{m,u}(\theta) d\psi_{-u}. \quad (6)$$

Here we have described iterated filtering algorithms while ignoring parameter perturbations and Monte Carlo evaluations of the likelihood function. Each of these components

play an important role in the practicality of iterated filtering, but existing theoretical justifications rely on the convergence of data cloning to the MLE, and show that the convergence still holds in spite of the additional complexities. A natural question is whether iterating Eqs. 5–6 results in a probability distribution with all mass centered at the MLE(s), similar to the case without marginalization. The nonlinearization introduced by the marginalization, however, adds difficulty to the task of calculating, or bounding, the density. In particular, previous approaches that rely on the linearization of unnormalized Bayes updates (e.g., Ionides et al., 2015) are no longer applicable. In the following subsection, we show that iterating Eqs. 5–6 does converge to the MLE when the likelihood is Gaussian.

### 3.1 Consistency for Gaussian models

For Gaussian models, conditioning and marginalization can be carried out exactly. The properties of this analytically tractable special case is relevant to the broader class of models that is well approximated by Gaussian models, for example, models satisfying the widely studied property of local asymptotic normality (LAN) (Le Cam and Yang, 2000). We show in Theorem 2 that MPIF for a Gaussian model converges to the exact MLE as long as unit-specific parameters are not highly informative about shared parameters.

As before, we assume there are  $U \geq 1$  units, and the likelihood of each unit is defined by  $L_u(\theta; y_u^*) = f_u(y_u^*; \phi, \psi_u)$ , and the likelihood of the entire model is  $L(\theta; \mathbf{y}^*) = \prod_{u=1}^U L_u(\theta; y_u^*)$ . In what follows, we assume  $\phi \in \mathbb{R}$  and  $\psi_u \in \mathbb{R}$  for all  $u \in 1:U$  in order to ease the notation and analysis.

**Theorem 2.** *Let  $f_u(y_u^*; \phi, \psi_u)$  be the density that corresponds to a Gaussian distribution*

with mean  $(\phi^*, \psi_u^*)$  and precision  $\Lambda_u^* \in \mathbb{R}^{2 \times 2}$ . Assume that  $\Lambda_u^*$  satisfies

$$[\Lambda_u^*]_{1,2}^2 < \frac{4[\Lambda_u^*]_{1,1} [\Lambda_u^*]_{2,2} \sum_{k=1}^U [\Lambda_k^*]_{1,1}}{\left([\Lambda_u^*]_{2,2} + \sum_{k=1}^U [\Lambda_k^*]_{1,1}\right)^2}. \quad (7)$$

If the initial prior density  $\pi_0(\theta) = \pi_{0,0}(\theta)$  corresponds to a Gaussian distribution with mean  $\mu_0$  and covariance  $\Sigma_0$ , then the density of the  $m$ th iteration of Eq. 6 corresponds to a Gaussian distribution with mean  $\mu_m \in \mathbb{R}^{U+1}$  and covariance  $\Sigma_m \in \mathbb{R}^{(U+1) \times (U+1)}$  such that  $\mu_m \rightarrow (\phi^*, \psi_1^*, \dots, \psi_U^*)$  and  $\|\Sigma_m\|_2 \rightarrow 0$ . That is, the algorithm converges in probability to the MLE.

The proof of Theorem 2 is included in Appendix C. At each iteration of the algorithm, the marginalization step results in a loss of information about the likelihood surface. In the Gaussian setting, this equates to setting the covariance term between the shared and unit-specific parameters to be zero before performing a Bayes update. The assumption in Eq. (7) therefore helps mitigate this loss of information by controlling the size of the covariance in the likelihood surface. If the data are transformed to ensure that the likelihood covariance matrix has ones on the diagonal, that is,  $\Sigma_u^* = (\Lambda_u^*)^{(-1)} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ , then  $[\Lambda_u^*]_{1,1} = [\Lambda_u^*]_{2,2} = 1/(1 - \rho^2)$  and  $[\Lambda_u^*]_{1,2} = [\Lambda_u^*]_{2,1} = -\rho/(1 - \rho^2)$ . In this case, the convergence condition in Eq. (7) becomes  $\rho < 2/(\sqrt{U}(1 + 1/U))$ .

The proof of Theorem 2 shows that the condition in Eq. (7) is sufficient for a convergence guarantee, but it may not be necessary. Furthermore, even some asymptotic bias may be tolerable compared to alternative algorithms that fail to scale. As demonstrated in Section 4.1, the particle depletion suffered by PIF can result in MPIF obtaining a better approximation of the unmarginalized map than the PIF algorithm. In this case, MPIF is preferable to PIF even if the marginalized map results in a small bias.

Theorem 2 provides convergence results for the algorithm in the absence of parameter perturbations. Using a similar setup, we can now consider the behavior of the algorithm with perturbations added to model parameters at each step. Let  $f * g$  denote the convolution of probability densities  $f$  and  $g$ . We assume that  $h_{u,m}(\theta)$  is some perturbation density, and we modify the marginalized Bayes maps by adding this random noise at each unit-iteration.

$$\tilde{\pi}'_{m,u}(\theta) \propto f_u(y_u^*; \phi, \psi_u) (\pi'_{m,u-1} * h_{u,m})(\theta) \quad (8)$$

$$\pi'_{m,u}(\theta) \propto \int \tilde{\pi}'_{m,u}(\theta) d\phi d\psi_u \times \int \tilde{\pi}'_{m,u}(\theta) d\psi_{-u}, \quad (9)$$

Corollary 1 shows that, under similar conditions as Theorem 2, marginalized data cloning with perturbations also converges to a point mass at the MLE if the likelihood is Gaussian.

**Corollary 1.** *Consider the setup of Theorem 2. If the parameter perturbations are Gaussian with covariance matrix  $\sigma_m^2 \Sigma_0$  for some initial covariance matrix  $\Sigma_0 \in \mathbb{R}^{(U+1) \times (U+1)}$  and sequence  $\sigma_m^2 = o(1/m)$ , then the  $m$ th iteration of Eqs. 8 and 9 corresponds to a Gaussian distribution with mean  $\mu'_m \in \mathbb{R}^{U+1}$  and covariance  $\Sigma'_m \in \mathbb{R}^{U+1 \times U+1}$ . If  $\hat{\theta} = (\phi^*, \psi_1^*, \dots, \psi_U^*)$  denotes the MLE, then  $|\mu'_m - \hat{\theta}|_2 \rightarrow 0$  and  $\|\Sigma'_m\|_2 \rightarrow 0$ .*

The convergence of the Eqs. 8–9 can be partially explained using a common heuristic in Bayesian analysis: a more dispersed prior typically results in a posterior distribution that more closely resembles the likelihood function. In an iterated Gaussian setting, adding noise at each step results in intermediate prior distributions that have the same mean, but larger variance. Therefore each iteration of Eq. 8 is expected to result in a mean closer to the MLE than the case without perturbations (Eq. 5). Following this logic, if the perturbations are chosen to ensure that they eventually approach zero, then the convergence of the unperturbed marginalized data cloning algorithm heuristically implies



the convergence of the perturbed version of the algorithm, as the perturbations to the prior densities at each step result in larger movements of the posterior density toward the MLE. The proof of Corollary 1 in Appendix C.1 demonstrates that this is true for the Gaussian model and the chosen perturbation schedule.

In principle, the marginalization procedure can be applied at various steps in the data cloning process and one can obtain similar convergence results. In Figure 1, we demonstrate the difference between data cloning and marginalized data cloning for a two parameter model with Gaussian likelihoods and priors with only a single unit but applying the marginalization for all parameters. This useful visualization demonstrates how, even when the likelihoods can be computed exactly, the marginalization only makes small modifications to the intermediate distributions.

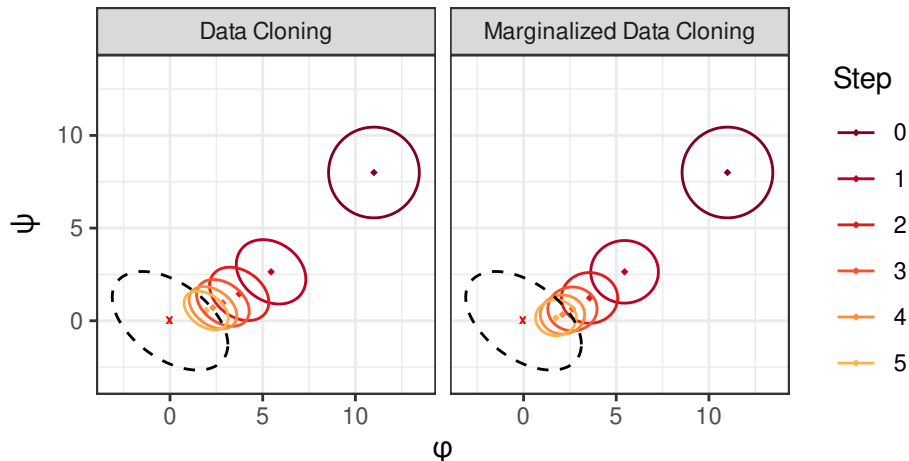


Figure 1: Data cloning and marginalized data cloning for two parameter model with Gaussian likelihoods and priors. The ellipses show the region of the parameter distribution that contains 95% of the probability mass of the distribution. The black dashed line shows this region for the likelihood surface, and the red “x” marks the MLE. Theorem 2 implies that the intermediate posterior densities will converge to a point mass at the MLE.

## 4 Simulation Studies

### 4.1 Marginalization to reduce particle depletion

The primary benefit and motivation of the marginalization step is improving the particle representations of the intermediate parameter distributions. In this sense, the marginalization step can be viewed as an attempt to take advantage of a bias-variance tradeoff. The marginalization procedure introduces a small amount of bias in the Bayesian posterior at each step in order to greatly reduce the variance of the particle representations of the distribution.

We demonstrate this idea via a simple simulation study that explores the particle representation of parameter distributions with and without marginalization for only a single unit-iteration within Algorithm 1. For our model, we suppose  $Y_{u,n}$  are independent and identically distributed (iid) from a normal  $\mathcal{N}[\psi_u, 1]$  distribution, and do not specify a latent process model as it is irrelevant for this model. We consider only  $U = 2$  units, and  $N_u = N = 100$  for all  $u$ . For our prior distribution, we let  $\Theta_{1:J}^{(0)} \stackrel{\text{iid}}{\sim} \mathcal{N}[\mu_0, \Sigma_0]$ , and use  $J = 1000$  particles to represent the joint parameter density. This simple model and setup is selected so that the priors and likelihoods can be exactly calculated; we can compare this to their particle representations using both versions of a single  $u = 1$  iteration of Algorithm 1 (Lines 4–18).

When iterating through unit  $u = 1$ , the Bayes map that is approximated by the unmarginalized filter requires an update to the particles that correspond to all model parameters for each time step  $n \in 1:N_1$ . This reduces the number of unique particles that represent the intermediate posterior distributions for parameter  $\Psi_2$  (Figure 2A). The number of unique particles representing  $\Psi_1$  remains high as a result of the added parameter

perturbations (line 7). On the other hand, the MPIF algorithm does not require resampling the  $\Psi_2$  particles during the  $u = 1$  iteration, and thus maintains the same number of unique particles during this update (dashed horizontal line in Figure 2A).

Figure 2B shows the filtered parameter particle swarm  $\Theta_{1,1:J}^{F,1}$  after the single unit update under both versions of the algorithm compared to the Bayes posterior distribution that PIF approximates. Although the marginalized version of the algorithm does not directly approximate this distribution, the particle swarm suffers less from particle depletion. This results in a better approximation of the intermediate posterior, despite introducing a small amount of bias. Theorem 2 provides sufficient conditions where the added bias at each step is negligible enough for the algorithm to converge to the MLE.

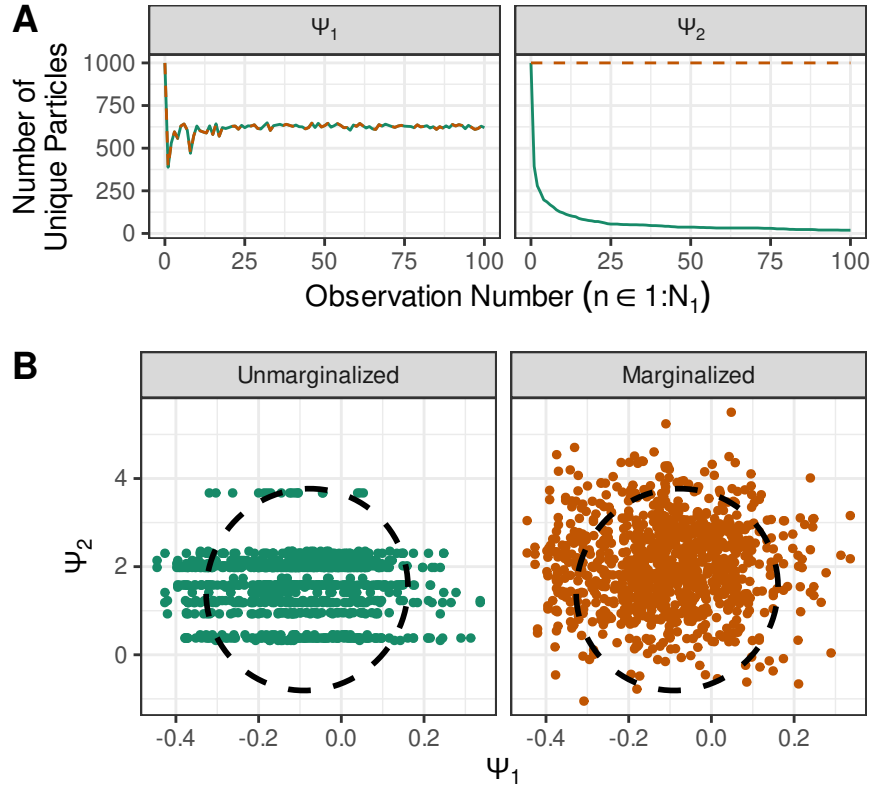


Figure 2: Updating parameter distributions with a single  $u = 1$  iteration of both versions of Algorithm 1. (A) The total number of unique particles representing each parameter. The dashed horizontal line shows that MPIF maintains the number of particles for  $\Psi_2$  over time. (B) Parameter particle swarm of a single update with and without marginalization compared to the true posterior distribution.

## 4.2 Stochastic Gompertz Population Model

We demonstrate the efficacy of the MPIF algorithm by estimating the model parameters of a high-dimensional, nonlinear PanelPOMP model by fitting data simulated from a collection of stochastic Gompertz population models. This model is commonly used to describe population dynamics in Ecology and has been used as a benchmark for comparison between various algorithms in previous studies (Bretó et al., 2020). The model assumes a latent state vector  $X_{u,n}$  for each unit  $u \in 1 : U$  and  $n \in 0 : N$ . For each unit  $u$ , the latent state has a one-step transition density that satisfies  $X_{u,n+1} = K_u^{1-e^{r_u}} X_{u,n}^{e^{-r_u}} \epsilon_{u,n}$ , where  $K_u$  is the carrying capacity for the population in unit  $u$ ,  $r_u$  is a positive parameter that corresponds to the growth rate, and  $\epsilon_{u,n}$  are iid log-normal random variables such that  $\log \epsilon_{u,n} \stackrel{\text{iid}}{\sim} \mathcal{N}[0, \sigma_u^2]$ .

Measurements of the population are obtained via the density  $\log Y_{u,n} \stackrel{\text{iid}}{\sim} \mathcal{N}[\log X_{u,n}, \tau_u^2]$  where  $\tau_u$  is a positive variance parameter. This model is a convenient nonlinear non-Gaussian PanelPOMP model because a logarithmic transformation makes the model a linear Gaussian process, and as such the exact likelihood of the model can be calculated by the Kalman filter (Kalman, 1960).

For this simulation study, we generate data from several Gompertz population models, with equal number of observations  $N$  in unit  $u$ , with values of  $N \in \{20, 50, 100\}$  and values of  $U$  ranging from  $U = 5$  to  $U = 2500$ . To generate data from this model, we fix  $K_u = 1$  and  $X_{u,0} = 1$  for all  $u$  and treat these parameters as known constants. We then set  $\sigma_u^2 = 0.01$  and  $r_u = 0.1$  for all values of  $u$  to generate data, but treat these parameters as unknown shared parameters that need to be estimated from the data. Finally, we set  $\tau_u^2 = 0.01$  for all  $u$  and treat these parameters as unknown unit-specific parameters. These values were chosen to obtain comparable simulations and results as Bretó et al. (2020).

The models were fit using the MPIF algorithm with the number of iterations  $M = 50$ ,

and the number of particles  $J = 1000$ . For this analysis, intermediate parameter estimates are obtained every five iterations of the MPIF algorithm, and the likelihood of the intermediate parameter values are obtained. The goal of calculating the intermediate likelihood values is to demonstrate how many iterations are needed to obtain model convergence, and to compare the algorithms performance against that of the PIF algorithm for each time step. Because the maximization procedure is inherently stochastic, it is recommended to try multiple starting parameter values. Therefore for each model, 50 unique starting points are used; these starting points are randomly sampled from the hypercube with lower-bounds for each parameter determined by the generating parameter value divided by two, and the upper bound for the parameter defined as the generating parameter value multiplied by two.

The results of this simulation study with  $N = 50$  are shown in Figure 3. For every combination of  $\{U, N\}$  considered, and for all numbers of MIF iterations, the maximum likelihood obtained using MPIF was higher than the maximum obtained using PIF. For large  $U$ , we also found that the worst performing Monte Carlo replicate of the MPIF algorithm often did better than the best performing replication of the PIF algorithm. In the next section, we find similar results using a more complicated model to describe data that have previously been used as a means of comparing inference procedures.

## 5 Measles in the United Kingdom

We show how MPIF and PIF compare in an epidemiological model for weekly reported measles cases for 20 different UK cities from 1950 to 1964 (Korevaar et al., 2020). Pre-vaccination UK measles data has been used extensively to motivate innovative methods for inference on stochastic processes since Bartlett (1960), yet, fitting nonlinear Markov

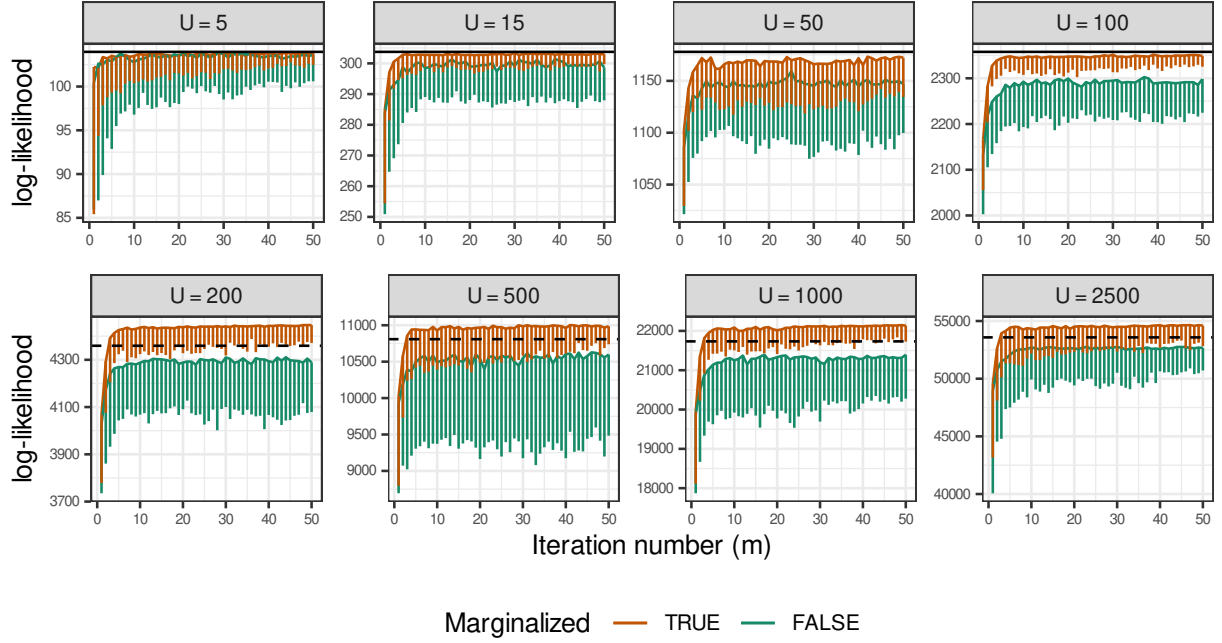


Figure 3: Comparison of the MPIF and PIF algorithms for fitting the stochastic Gompertz population model. The solid horizontal line shows the true maximum likelihood, determined via the Kalman filter and a numeric optimizer, an intractable approach for high-dimensional parameter spaces ( $U > 100$ ); in these cases, the dashed line indicates the likelihood at the data-generating parameters. Each algorithm used 50 unique starting points. Vertical bars span the tenth percentile to the maximum likelihood values

process models simultaneously to multiple cities with shared and unit-specific parameters remains a challenge, leading practitioners to consider linearizations that have uncertain consequences (Korevaar et al., 2020). We fit three different models to the data, all based on the susceptible-exposed-infectious-recovered (SEIR) model of He et al. (2010). The state process,  $X_u(t) = (S_t^{(u)}, E_t^{(u)}, I_t^{(u)}, R_t^{(u)})$ , tracks the number of susceptible, exposed, infected and recovered individuals in each unit  $u$ . The total population size,  $\text{pop}^{(u)}(t)$ , is treated as known, interpolated from census data using cubic splines, and we use this constraint to avoid explicit specification of  $R_t^{(u)}$ . State transitions are generated using an Euler approximation to a continuous-time Markov chain, with time step  $\Delta = 1$  day, as

follows:

$$\begin{aligned}
A_{BS,t}^{(u)} &\sim \text{Pois}(\mu_{BS}^{(u)}(t) \Delta) \\
(A_{SE,t}^{(u)}, A_{SD,t}^{(u)}) &\sim \text{Eulermultinom}\left(S_t^{(u)}, \bar{\mu}_{SE}^{(u)}(t)(\Gamma^{(u)}(t + \Delta) - \Gamma(t))/\Delta, \mu_{SD}^{(u)}, \Delta\right) \\
(A_{EI,t}^{(u)}, A_{ED,t}^{(u)}) &\sim \text{Eulermultinom}(E_t^{(u)}, \mu_{EI}^{(u)}(t), \mu_{ED}^{(u)}, \Delta) \\
(A_{IR,t}^{(u)}, A_{ID,t}^{(u)}) &\sim \text{Eulermultinom}(I_t^{(u)}, \mu_{IR}^{(u)}(t), \mu_{ID}^{(u)}, \Delta), \\
S_{t+\Delta}^{(u)} &= S_t^{(u)} + A_{BS,t}^{(u)} - A_{SE,t}^{(u)} - A_{SD,t}^{(u)}, \\
E_{t+\Delta}^{(u)} &= E_t^{(u)} + A_{SE,t}^{(u)} - A_{EI,t}^{(u)} - A_{ED,t}^{(u)}, \\
I_{t+\Delta}^{(u)} &= I_t^{(u)} + A_{EI,t}^{(u)} - A_{IR,t}^{(u)} - A_{ID,t}^{(u)}.
\end{aligned}$$

Here,  $A_{BC,t}^{(u)}$  counts transitions from compartment  $B$  into  $C$  for unit  $u$  between time  $t$  and  $t + \Delta$ ,  $\text{Pois}(\lambda)$  is a Poisson distribution with mean  $\lambda$ ,  $\Gamma^{(u)}(t)$  is a gamma process with mean  $t$  and intensity  $\sigma_{SE}^{(u)}$  (Bretó et al., 2009), and  $\text{Eulermultinom}(n, \mu_1, \mu_2, \Delta)$  is a multinomial distribution with  $n$  independent trials and event probabilities given by  $p_0 = \exp\{- (\mu_1 + \mu_2)\Delta\}$  and  $p_i = \frac{\mu_i}{\mu_1 + \mu_2}(1 - p_0)$ . The Eulermultinom outcome is the number of events of type  $p_1$  and  $p_2$ , which correspond to transitions out of the source compartment. The remaining events, of type  $p_0$ , correspond to individuals remaining in the source compartment. The rate of arrivals into the susceptible class,  $\mu_{BS}^{(u)}(t)$ , is given by

$$\mu_{BS}^{(u)}(t) = (1 - c^{(u)}) b^{(u)}(t - \tau_d) + c^{(u)} \delta((t - \tau_c) \bmod \tau_y) \int_{t-\tau_y}^t b^{(u)}(t - \tau_d - s) ds,$$

where  $\delta$  is the Dirac delta function,  $b^{(u)}(t)$  is the births per year at time  $t$  interpolated from annual data using cubic splines,  $\tau_d$  is the delay between when the births take place and when they actually contribute to the transition rate,  $\tau_y = 1$  year, and  $\tau_c$  is the school

admission day, i.e., the 251st day of the year. We set  $\tau_d = 4$  years to describe the age at which individuals typically enter a high-transmission environment.

The rate at which individuals enter state  $E$  is  $\mu_{SE}^{(u)}(t) = \bar{\mu}_{SE}^{(u)}(t) \frac{d\Gamma^{(u)}(t)}{dt}$ , where  $\bar{\mu}_{SE}^{(u)}(t)$  is the mean rate given by  $\bar{\mu}_{SE}^{(u)}(t) = \frac{\beta^{(u)}(t) (I_t^{(u)} + \iota^{(u)})}{\text{pop}^{(u)}(t)}$ , where  $\text{pop}^{(u)}(t)$  is the city population at time  $t$  interpolated from annual data, and  $\iota^{(u)}$  is the average number of infected individuals visiting the city at any time. The force of infection,  $\beta^{(u)}(t)$ , is given by

$$\beta^{(u)}(t) = \begin{cases} \beta_0^{(u)} (1 + a^{(u)}(1 - p)/p) & \text{during school term} \\ \beta_0^{(u)} (1 - a^{(u)}) & \text{during vacation} \end{cases}$$

$$\beta_0^{(u)} = \mathcal{R}_0^{(u)} (1 - \exp \{ - (\mu_{IR}^{(u)}(t) + \mu_{ID}^{(u)}) \Delta \}) / \Delta.$$

where  $p = 0.7589$  is the proportion of the year occupied by the school term. The remaining transition rates are assumed to be constant:  $\mu_{EI}^{(u)}(t) = \sigma^{(u)}$ ,  $\mu_{IR}^{(u)}(t) = \gamma^{(u)}$ , and  $\mu_{SD}^{(u)} = \mu_{ED}^{(u)} = \mu_{ID}^{(u)} = 0.02 \text{ yr}^{-1}$ . We follow He et al. (2010) by using a discretized normal distribution for the number of cases reported:

$$P(Y_n^{(u)} = y^{(u)} | Z_n^{(u)} = z^{(u)}) = \Phi(y^{(u)} + 0.5, \rho^{(u)} z^{(u)}, \rho^{(u)}(1 - \rho^{(u)}) z^{(u)} + [\psi^{(u)} \rho^{(u)} z^{(u)}]^2) \\ - \Phi(y^{(u)} - 0.5, \rho^{(u)} z^{(u)}, \rho^{(u)}(1 - \rho^{(u)}) z^{(u)} + [\psi^{(u)} \rho^{(u)} z^{(u)}]^2)$$

where  $\Phi(\cdot; \mu, \sigma^2)$  is the CDF for a  $\mathcal{N}[\mu, \sigma^2]$  random variable and  $Z_n$  is the number of people transitioning from compartment  $I$  to  $R$  between the  $(n - 1)th$  and  $nth$  observation times. Lastly, we estimate the proportion of individuals in the first three states at time  $t_0$ ,  $S_0^{(u)}$ ,  $E_0^{(u)}$ , and  $I_0^{(u)}$ . Given that these proportions along with  $R_0^{(u)}$  must add up to 1, there is no need to actually estimate  $R_0^{(u)}$ . Consequently, the total number of parameters per unit is 12.



We use three different models based on the He et al. (2010) model, with the key difference of investigating subsets of the parameters that might be well modeled as shared.

1. The “ $c$ -shared” model uses a shared parameter for  $c$ .
2. The “ $\iota$ -shared” model uses a log-log linear model between  $\iota$  and the city population for year 1950, specifically  $\log(\iota^{(u)}) = \iota_1 + \iota_2 \cdot \log(\text{pop}^{(u)}(1950))$ .
3. The “7-shared” model uses the log-log linear model for  $\iota$  and shared parameters for  $c$ ,  $\mathcal{R}_0$ ,  $\gamma$ ,  $\sigma$ ,  $\sigma_{SE}$ , and  $a$ .

For each model, we run MPIF and PIF for 200 iterations, each search starting from one of 36 different parameter vectors randomly sampled from a hypercube where each dimension is slightly larger than the range spanned by the corresponding unit-specific MLE’s in He et al. (2010). We perform these searches using 500, 5000, and 10000 particles to discern how the fits yielded by MPIF and PIF differ based on the selected particle count. For the present data set, 500 particles is too low for proper optimization, 5000 is adequate, and 10000 enables a thorough parameter search. The log-likelihood for each fit is evaluated using the average of replicated particle filter evaluations with 10000 particles at evenly-spaced iterations: 20, 56, 92, 128, 164, and 200. Figure 4 summarizes the output of this Monte Carlo optimization search. Similar to the Gompertz population model, the MPIF algorithm consistently yields parameter estimates corresponding to higher likelihoods than the PIF algorithm.

In practice, because iterated filtering algorithms are stochastic optimization algorithms, many Monte Carlo replicates are conducted from various initialization points, and final parameter estimates correspond to the search with the highest likelihood. Because of this, we are primarily interested in how in the maximum estimate from each algorithm compares.

In Figure 4, we note that MPIF consistently yields larger sample maximums of the log-likelihood evaluations irrespective of the model or particle count we use. By 200 iterations, we see that maximum log-likelihoods obtained using MPIF are 70 to 105 units higher than PIF for the  $\iota$ -shared and  $c$ -shared models, or about 4 to 5 units per city. MPIF has about half the advantage for the 7-shared model by 200 iterations. Close inspection of the results show that for nearly all combinations of number of particles, iterations, and model specification, MPIF outperforms the PIF algorithm. In cases where this is not true, the distribution for PIF-generated log-likelihoods have especially long right-tails, suggesting that the observed advantage for PIF in these scenarios is a result of large variance working in its favor. Increasing the number of particles and iterations reduces the variance among Monte Carlo replications, and when the number of particles ( $J$ ) is largest, the advantage of MPIF over PIF becomes more evident for all models variations.

The boxplots in Figure 4 demonstrate that the standard deviation of the log-likelihood across Monte Carlo replicates is consistently lower for MPIF. In addition to suggesting that even poor performing Monte Carlo replicates of the MPIF algorithm have better outcomes than the best replicates for the PIF algorithm, low standard deviation across initializations is useful in practice as the convergence of iterated filtering algorithms is often judged by whether or not distinct Monte Carlo replicates finish at the same maximum, within suitable Monte Carlo error. Reducing the variance between Monte Carlo replicates also results in tighter confidence intervals when computing Monte Carlo adjusted profile confidence intervals (Ionides et al., 2017).

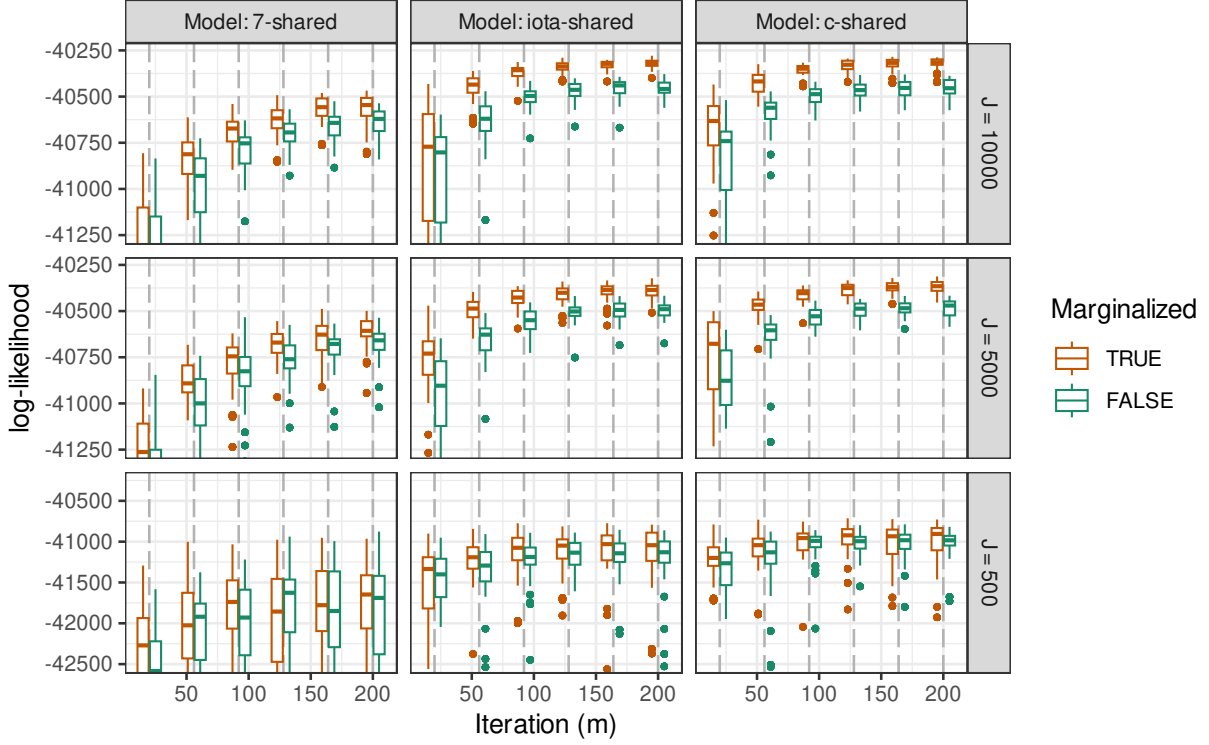


Figure 4: Log-likelihoods yielded by PIF and MPIF for the mechanistic measles model to the UK data. Rows correspond to a different number of particles  $J$  used in Algorithm 1. The log-likelihood is evaluated at iterations 20, 56, 92, 128, 164, and 200.

## 6 Discussion

The issue of particle depletion when estimating the posterior distribution of parameter values using a particle filter has previously been noted (e.g., Chen et al., 2025, Section 1.2). A key innovation of iterated filtering algorithms is that adding parameter perturbations helps revive the particle representations from a depleted state by adding artificial noise to the parameter particles. These perturbations are precisely why the number of unique particles representing the distribution of  $\Psi_1$  in Panel A of Figure 2 does not degenerate to zero over time. Given this observation, an alternative approach to solving the particle depletion issue that arises in higher dimensions would be to perturb all model parameters at each time-step. This approach is supported by Theorem 1, and avoids the analytic challenges associated with adding the marginalization step. However, perturbing all model

parameters at each step is equivalent to applying a vanilla iterated filtering algorithm to a PanelPOMP model, which generally performs worse on panel models than the PIF variation (Bretó et al., 2020).

Liu and West (2001) note that the artificial noise introduced from parameter perturbations results in loss of information. This loss of information motivates the common practice in applications of iterated filtering of only perturbing model parameters when the data used to calculate the particle weights has direct relation with the parameters that are being perturbed. For instance, it is common practice to only perturb parameters that are unique to the initialization density  $f_{X_{u,0}}(x_{u,0}; \theta)$  at the first available observations, as these observations have the strongest signal for the initialization parameters. If initial parameters are perturbed at all observation times, then the signal from these initial observations gets lost over time. This same principle serves as a motivator for the PIF algorithm and describes why PIF is more successful than vanilla iterated filtering algorithms applied to PanelPOMPs: perturbing unit-specific parameters while considering data from other units results in a significant loss of information.

The challenge of iterating filtering for PanelPOMP models can be described as a tradeoff between particle depletion and a loss of signal due to parameter perturbations. The MPIF algorithm introduced here is designed to address both of these challenges simultaneously by avoiding perturbations when the signal is weak, but not resampling unit-specific parameters using weights calculated with data from other units. The cost of this modification is a small amount of bias for the particle representation of the posterior distribution at each step. Theorem 2 formally demonstrates that the affect of the bias can be completely negated if the log-likelihood is quadratic, which is the limiting behavior of all likelihood surfaces.

Finally, Theorem 1 does provide some stronger guarantees for MPIF under additional

constraints on the PanelPOMP model that we have not yet mentioned. In a model with only unit-specific parameters ( $\Phi = \emptyset$ ), for instance, MPIF is equivalent to conducting IF2 independently on each unit, and therefore the convergence results of Theorem 1 apply. Similarly, if the model only contains shared parameters ( $\Psi_{1:U} = \emptyset$ ), then MPIF is equivalent to PIF, and once again stronger theoretical guarantees are available. This equivalency also provides some intuition as to when MPIF will outperform its alternatives. In a model consisting primarily of shared parameters, MPIF behaves very similarly to PIF and adds a smaller advantage relative to the case when there are more unit-specific parameters. We have seen this pattern in our results, as the gain in log-likelihood obtained via MPIF was smallest for the measles model with only one unit specific parameter.

## References

- Abkemeier, A., Chen, J., Ionides, E., Wheeler, J., and Tan, K. (2025), “pypomp,” URL <https://github.com/pypomp/pypomp>.
- Andrieu, C., Doucet, A., and Holenstein, R. (2010), “Particle Markov chain Monte Carlo methods,” *Journal of the Royal Statistical Society, Series B*, 72, 269–342.
- Arulampalam, M. S., Maskell, S., Gordon, N., and Clapp, T. (2002), “A tutorial on particle filters for online nonlinear, non-Gaussian Bayesian tracking,” *IEEE Transactions on Signal Processing*, 50, 174 – 188.
- Auger-Méthé, M., Newman, K., Cole, D., Empacher, F., Gryba, R., King, A. A., Leos-Barajas, V., Mills Flemming, J., Nielsen, A., Petris, G., and Thomas, L. (2021), “A guide to state-space modeling of ecological time series,” *Ecological Monographs*, 91, e01470.
- Baker, R. E., Pena, J.-M., Jayamohan, J., and Jérusalem, A. (2018), “Mechanistic models

- versus machine learning, a fight worth fighting for the biological community?” *Biology letters*, 14, 20170660.
- Bartlett, M. S. (1960), *Stochastic population models in ecology and epidemiology*, volume 4 of *Methuen’s Monographs on Applied Probability and Statistics*, London: Spottiswoode, Ballantyne & Co. Ltd.
- Bengtsson, T., Bickel, P., and Li, B. (2008), “Curse-of-dimensionality revisited: Collapse of the particle filter in very large scale systems,” in *Probability and statistics: Essays in honor of David A. Freedman*, volume 2, Institute of Mathematical Statistics, 316–335.
- Bretó, C., He, D., Ionides, E. L., and King, A. A. (2009), “Time series analysis via mechanistic models,” *The Annals of Applied Statistics*, 319–348.
- Bretó, C., Ionides, E. L., and King, A. A. (2020), “Panel data analysis via mechanistic models,” *Journal of the American Statistical Association*, 115, 1178–1188.
- Bretó, C., Wheeler, J., King, A. A., and Ionides, E. L. (2025), “panelPomp: Analysis of panel data via partially observed Markov processes in R,” *The R Journal*, 17, 180–199.
- Chen, Y., Gerber, M., Andrieu, C., and Douc, R. (2025), “Self-organizing state-space models with artificial dynamics,” *arXiv:2409.08928*.
- Crauel, H. (2002), *Random probability measures on Polish spaces*, volume 11, CRC press.
- Domeyer, J. E., Lee, J. D., Toyoda, H., Mehler, B., and Reimer, B. (2022), “Driver-pedestrian perceptual models demonstrate coupling: Implications for vehicle automation,” *IEEE Transactions on Human-Machine Systems*, 52, 557–566.
- Evensen, G. (2009), “The ensemble Kalman filter for combined state and parameter estimation,” *IEEE Transactions on Control Systems*, 29, 83–104.

- Fox, S. J., Lachmann, M., Tec, M., Pasco, R., Woody, S., Du, Z., Wang, X., Ingle, T. A., Javan, E., Dahan, M., Gaither, K., Escott, M. E., Adler, S. I., Johnston, S. C., Scott, J. G., and Meyers, L. A. (2022), “Real-time pandemic surveillance using hospital admissions and mobility data,” *Proceedings of the National Academy of Sciences of the USA*, 119, e2111870119.
- Häggström, H., Persson, S., Cvijovic, M., and Picchini, U. (2025), “Simulation-based inference for stochastic nonlinear mixed-effects models with applications in systems biology,” *arXiv:2504.11279*.
- He, D., Artzy-Randrup, Y., Musa, S. S., Gräf, T., Naveca, F., and Stone, L. (2024), “Modelling the unexpected dynamics of COVID-19 in Manaus, Brazil,” *Infectious Disease Modelling*, 9, 557–568.
- He, D., Ionides, E. L., and King, A. A. (2010), “Plug-and-play inference for disease dynamics: Measles in large and small populations as a case study,” *Journal of The Royal Society Interface*, 7, 271–283.
- Hewitt, J., Wilson-Henjum, G., Collins, D. T., Linder, T. J., Lenocho, J. B., Heale, J. D., Quintanal, C. A., Pleszewski, R., McBride, D. S., Bowman, A. S., Chandler, J. C., Shriner, S. A., Bevins, S. N., Kohler, D. J., Chipman, R. B., Gosser, A. L., Bergman, D. L., DeLiberto, T. J., and Pepin, K. M. (2024), “Landscape-scale epidemiological dynamics of SARS-CoV-2 in white-tailed deer,” *Transboundary and Emerging Diseases*, 2024, 7589509.
- Hogg, D. W. and Villar, S. (2024), “Position: Is machine learning good or bad for the natural sciences?” in Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (editors), *Proceedings of the 41st International Confer-*

- ence on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, PMLR.
- Ionides, E. L., Breto, C., Park, J., Smith, R. A., and King, A. A. (2017), “Monte Carlo profile confidence intervals for dynamic systems,” *Journal of the Royal Society Interface*, 14, 1–10.
- Ionides, E. L., Nguyen, D., Atchadé, Y., Stoev, S., and King, A. A. (2015), “Inference for dynamic and latent variable models via iterated, perturbed Bayes maps,” *Proceedings of the National Academy of Sciences of the USA*, 112, 719–724.
- Ionides, E. L., Ning, N., and Wheeler, J. (2024), “An iterated block particle filter for inference on coupled dynamic systems with shared and unit-specific parameters,” *Statistica Sinica*, 34, 1241–1262.
- Kalman, R. E. (1960), “A new approach to linear filtering and prediction problems,” *Journal of Basic Engineering*, 82, 35–45.
- Korevaar, H., Metcalf, C. J., and Grenfell, B. T. (2020), “Structure, space and size: Competing drivers of variation in urban and rural measles transmission,” *Journal of the Royal Society Interface*, 17, 20200010.
- Kramer, S. C., Pirikahu, S., Casalegno, J.-S., and Domenech de Cellès, M. (2024), “Characterizing the interactions between influenza and respiratory syncytial viruses and their implications for epidemic control,” *Nature Communications*, 15, 10066.
- Le Cam, L. and Yang, G. L. (2000), *Asymptotics in statistics*, New York: Springer, 2nd edition.



- Lee, E. C., Chao, D. L., Lemaitre, J. C., Matrajt, L., Pasetto, D., Perez-Saez, J., Finger, F., Rinaldo, A., Sugimoto, J. D., Halloran, M. E., Longini, I. M., Ternier, R., Vissieres, K., Azman, A. S., Lessler, J., and Ivers, L. C. (2020), “Achieving coordinated national immunity and cholera elimination in Haiti through vaccination: A modelling study,” *The Lancet Global Health*, 8, e1081–e1089.
- Lele, S. R., Dennis, B., and Lutscher, F. (2007), “Data cloning: Easy maximum likelihood estimation for complex ecological models using Bayesian Markov chain Monte Carlo methods,” *Ecology Letters*, 10, 551–563.
- Liu, J. and West, M. (2001), “Combined parameter and state estimation in simulation-based filtering,” in *Sequential Monte Carlo Methods in Practice*, New York, NY: Springer New York, 197–223.
- Newman, K., King, R., Elvira, V., de Valpine, P., McCrea, R. S., and Morgan, B. J. (2023), “State-space models for ecological time-series data: Practical model-fitting,” *Methods in Ecology and Evolution*, 14, 26–42.
- Ning, N. and Ionides, E. L. (2023), “Iterated block particle filter for high-dimensional parameter learning: Beating the curse of dimensionality,” *Journal of Machine Learning Research*, 24, 1–76.
- Ponciano, J. M., Taper, M. L., Dennis, B., and Lele, S. R. (2009), “Hierarchical models in ecology: Confidence intervals, hypothesis testing, and model selection using data cloning,” *Ecology*, 90, 356–362.
- Pons-Salort, M. and Grassly, N. C. (2018), “Serotype-specific immunity explains the incidence of diseases caused by human enteroviruses,” *Science*, 361, 800–803.

- Ranjeva, S., Subramanian, R., Fang, V. J., Leung, G. M., Ip, D. K., Perera, R. A., Peiris, J. M., Cowling, B. J., and Cobey, S. (2019), “Age-specific differences in the dynamics of protective immunity to influenza,” *Nature Communications*, 10, 1660.
- Ranjeva, S. L., Baskerville, E. B., Dukic, V., Villa, L. L., Lazcano-Ponce, E., Giuliano, A. R., Dwyer, G., and Cobey, S. (2017), “Recurring infection with ecologically distinct HPV types can explain high prevalence and diversity,” *Proceedings of the National Academy of Sciences of the USA*, 114, 13573–13578.
- Rebeschini, P. and Van Handel, R. (2015), “Can local particle filters beat the curse of dimensionality?” *Annals of Applied Probability*, 25, 2809–2866.
- Searle, C. L., Cortez, M. H., Hunsberger, K. K., Grippi, D. C., Oleksy, I. A., Shaw, C. L., de la Serna, S. B., Lash, C. L., Dhir, K. L., and Duffy, M. A. (2016), “Population density, not host competence, drives patterns of disease in an invaded community,” *The American Naturalist*, 188, 554–566.
- Snyder, C., Bengtsson, T., Bickel, P., and Anderson, J. (2008), “Obstacles to high-dimensional particle filtering,” *Monthly Weather Review*, 136, 4629–4640.
- Stocks, T., Britton, T., and Höhle, M. (2018), “Model selection and parameter estimation for dynamic epidemic models via iterated filtering: Application to rotavirus in Germany,” *Biostatistics*, 21, 400–416.
- Subramanian, R., He, Q., and Pascual, M. (2021), “Quantifying asymptomatic infection and transmission of COVID-19 in New York City using observed cases, serology, and testing capacity,” *Proceedings of the National Academy of Sciences of the USA*, 118.
- Toni, T., Welch, D., Strelkowa, N., Ipsen, A., and Stumpf, M. P. (2009), “Approximate

- Bayesian computation scheme for parameter inference and model selection in dynamical systems,” *Journal of the Royal Society Interface*, 6, 187–202.
- Wale, N., Jones, M. J., Sim, D. G., Read, A. F., and King, A. A. (2019), “The contribution of host cell-directed vs. parasite-directed immunity to the disease and dynamics of malaria infections,” *Proceedings of the National Academy of Sciences*, 116, 22386–22392.
- Wheeler, J., Rosengart, A., Jiang, Z., Tan, K., Treutle, N., and Ionides, E. L. (2024), “Informing policy via dynamic models: Cholera in Haiti,” *PLOS Computational Biology*, 20, e1012032.
- Whitehouse, M., Whiteley, N., and Rimella, L. (2023), “Consistent and fast inference in compartmental models of epidemics using Poisson approximate likelihoods,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 85, 1173–1203.
- Wiens, J., Gutttag, J., and Horvitz, E. (2014), “A study in transfer learning: Leveraging data from multiple hospitals to enhance hospital-specific predictions,” *Journal of the American Medical Informatics Association*, 21, 699–706.
- Wigren, A. and Lindsten, F. (2022), “Marginalized particle Gibbs for multiple state-space models coupled through shared parameters,” *arXiv:2210.07379*.
- Wood, S. N. (2010), “Statistical inference for noisy nonlinear ecological dynamic systems,” *Nature*, 466, 1102–1104.

## A Assumptions for Theorem 1

For all constants  $(\epsilon, r) \in (0, \infty) \times \mathbb{R}^d$ , we define  $B_\epsilon^d(r) = \{r' \in \mathbb{R}^d : |r - r'|_2 < \epsilon\}$ . Let  $\mathcal{B}^{d_\phi + Ud_\psi}$  be the Borel  $\sigma$ -algebra on the set of real numbers  $\mathbb{R}^{d_\phi + Ud_\psi}$ . We assume that  $\Theta \in \mathcal{B}^{d_\phi + Ud_\psi}$  is a compact set that satisfies (A1). Informally, this ensures that the corners of the set are not too sharp, and directly follows definition of a regular compact set from Chen et al. (2025).

(A1) There exists a continuous function  $\gamma: [0, \infty) \rightarrow [0, \infty)$  such that  $\lim_{x \downarrow 0} \gamma(x) = 0$ , and for all  $(\epsilon, x) \in (0, \infty) \times \Theta$ , there exists  $x' \in \Theta$  such that

$$B_{\gamma(\epsilon)}^{d_\phi + Ud_\psi}(x') \subseteq B_\epsilon^{d_\phi + Ud_\psi}(x) \cap \Theta.$$

We make the following assumptions on the probability densities that are used to define a PanelPOMP model described in Section 2.

(B1) There exists a compact set  $E \subset X$  such that

$$\inf_{(\theta, x) \in \Theta \times X} \int_E f_{X_n|X_{n-1}}(x_n|x_{n-1}; \theta) dx_n > 0,$$

for all  $n \in \{1, 2, \dots, N\}$ .

(B2)  $L(\theta; \mathbf{y}^*) > 0$  for all  $\theta \in \Theta$  and  $\sup_{(\theta, x_{u,n}) \in (\Theta, \mathcal{X})} f_{Y_{u,n}|X_{u,n}}(y_{u,n}^*|x_{u,n}; \theta) < \infty$  for all  $u \in 1:U$  and  $n \in 1:N_u$ .

(B3) The transition and measurement densities are sufficiently smooth functions of  $\theta$ , in the sense that for any  $\theta, \theta' \in \Theta$ , there exists a continuous and strictly increasing function  $g: [0, \infty) \rightarrow [0, \infty)$  and sequence of measurable functions  $\varphi_{u,0:N_u}: \mathcal{X}^2 \rightarrow \mathbb{R}$ ,

that satisfy:

$$|\log (f_{Y_{u,n}|X_{u,n}}(y_{u,n}^*|x_{u,n}; \theta)f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|x_{u,n-1}; \theta)) \\ - \log (f_{Y_{u,n}|X_{u,n}}(y_{u,n}^*|x_{u,n}; \theta')f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|x_{u,n-1}; \theta'))| \leq g(\|\theta - \theta'\|)\varphi_{u,n}(x_{u,n-1}, x_{u,n}),$$

With the constraint that for all  $u \in 1:U$ , there exists a  $\delta_u \in (0, \infty)$  such that

$$\int \exp \left\{ \delta_u \sum_{n=0}^{N_u} \varphi_{u,n}(x_{u,n-1}, x_{u,n}) \right\} f_{X_{u,0}}(x_{u,0}; \theta) \times \\ \prod_{n=1}^{N_u} f_{Y_{u,n}|X_{u,n}}(y_{u,n}^*|x_{u,n}; \theta) f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|x_{u,n-1}; \theta) dx_{u,0:N_u} < \infty,$$

using the convention that if  $n = 0$ , then

$$\begin{aligned} \varphi_{u,n}(x_{u,n-1}, x_{u,n}) &= \varphi_{u,0}(x_{u,0}, x_{u,0}) \\ f_{Y_{u,n}|X_{u,n}}(\cdot|x_{u,n}; \theta) &= 1 \\ f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|x_{u,n-1}; \theta) &= f_{X_{u,0}}(x_{u,0}; \theta) \end{aligned}$$

Finally, the following assumptions are made about the random perturbations in lines 4 and 7 of Algorithm 1. Let  $\mu_0$  denote the probability measure on  $(\mathbb{R}^{d_\phi+Ud_\psi}, \mathcal{B}^{d_\phi+Ud_\psi})$  that defines the distribution of the initial particle swarm, i.e.,  $\Theta_{1:J}^0 \stackrel{\text{iid}}{\sim} \mu_0$ . If  $\{\mu_n\}_{n \geq 1}$  is a sequence of random probability measures (Crauel, 2002) on  $(\mathbb{R}^{d_\phi+Ud_\psi}, \mathcal{B}^{d_\phi+Ud_\psi})$  with  $\mathcal{F}_n$  denoting the corresponding filtration, then we denote  $K_{\mu_n}$  to be the Markov kernel such that  $\theta \sim K_{\mu_n}(\theta', d\theta) \iff \theta \stackrel{\text{dist}}{=} \theta' + Z$ , where  $Z|\mathcal{F}_n \sim \mu_n$ .

Let  $S_{\tilde{u}} = \sum_{k=1}^{\tilde{u}} (N_k + 1)$ . We define the Markov kernel as a sequence in  $\tilde{n} \in \mathbb{N}$ , where  $\tilde{n}$  defines the values  $(m, u, n)$  via the equation  $\tilde{n} = (m - 1)S_U + S_{u-1} + n + 1$ , such that  $K_{\tilde{n}}(\theta_{\tilde{n}-1}, d\theta_{\tilde{n}}) = h_{u,n}(\theta_{\tilde{n}}|\theta_{\tilde{n}-1}; \sigma_{u,m})d\theta_{\tilde{n}}$ . Let  $\{U_{\tilde{n}}\}_{\tilde{n} \geq 1}$  be a sequence of  $\Theta$ -valued random

variables such that for all  $\tilde{n} \geq 1$  and sets  $A_1, \dots, A_{\tilde{n}} \in \mathcal{B}(\Theta)$ , where  $\mathcal{B}(\Theta)$  is the Borel  $\sigma$ -algebra on  $\Theta$ .

$$\mathbb{P}(Z_k \in A_k, k \in \{1, \dots, \tilde{n}\} | \mathcal{F}_{\tilde{n}}) = \prod_{k=1}^{\tilde{n}} \mu_k(A_k).$$

Then we assume the sequence of probability measures  $\{\mu_{\tilde{n}}\}_{\tilde{n} \geq 0}$  satisfy:

(C1)  $\mu_0(B_\epsilon(\theta)) > 0$  for all  $\theta \in \Theta$  and  $\epsilon \in (0, \infty)$ .

(C2) There exists a family of  $(0, 1]$ -valued random variables  $(\Gamma_\delta^\mu)_{\delta \in (0, \infty)}$  such that, for all

$\delta \in (0, \infty)$ , we have  $\mathbb{P}(\inf_{\tilde{n} \geq 1} \mu_{\tilde{n}}(B_\delta(0)) \geq \Gamma_\delta^\mu) = 1$ .

(C3)  $\mathbb{P}(\inf_{\tilde{n} \geq 1} \inf_{\theta' \in \Theta} \int_{\Theta} K_{\mu_{\tilde{n}}}(\theta', d\theta) \geq \Gamma^\mu) = 1$  for some  $(0, 1]$ -valued random variable  $\Gamma^\mu$ .

(C4) There exists a sequence of natural numbers  $\{k_{\tilde{n}}\}_{\tilde{n} \geq 1}$  and, for all  $l \in \mathbb{N}_0$ , a sequence of  $(0, \infty]$  valued functions  $\{g_{\tilde{n}, l}\}_{\tilde{n} \geq 1}$  defined on  $(0, \infty)$  such that

(a)  $\lim_{\tilde{n} \rightarrow \infty} k_{\tilde{n}}/\tilde{n} = \lim_{\tilde{n} \rightarrow \infty} 1/k_{\tilde{n}} = 0$ , and for all  $\epsilon \in (0, \infty)$ ,  $\lim_{\tilde{n} \rightarrow \infty} g_{\tilde{n}, l}(\epsilon) = 0$ .

(b) For all  $\tilde{n} \geq 1$  such that  $\tilde{n} > 2k_{\tilde{n}}$ ,  $k_{\tilde{n}}^* \in \{k_{\tilde{n}}, \dots, 2k_{\tilde{n}}\}$ , and for all  $\epsilon \in (0, \infty)$ ,

$$\frac{1}{\tilde{n} - k_{\tilde{n}}^*} \log \mathbb{P}\left(\exists s \in \{k_{\tilde{n}}^* + 1, \dots, \tilde{n}\} : \sum_{i=k_{\tilde{n}}^*+1}^s Z_i \notin B_\epsilon(0) | \mathcal{F}_{\tilde{n}}\right) \leq -\frac{1}{g_{\tilde{n}, l}(\epsilon)}.$$

(c) For all  $\tilde{n} > l$  and  $\epsilon \in (0, \infty)$ , we have

$$\frac{1}{\tilde{n} - l} \log \mathbb{P}\left(\sum_{i=l+1}^s Z_i \in B_\epsilon(0), \forall s \in \{l+1, \dots, \tilde{n}\} | \mathcal{F}_{\tilde{n}}\right) \geq -g_{\tilde{n}, l}(\epsilon).$$

## B Proof and discussion of Theorem 1

Theorem 1 is a straightforward extension of Theorem 4 of Chen et al. (2025) to PanelPOMP models. Our approach is to express an arbitrary PanelPOMP model as a POMP model

using a long format, where the latent and observed processes are stacked one unit after another to describe a single long time series that arises by stacking unit time series data one after another. Then, Algorithm 1 is equivalent to applying the theory developed in the appendix of Chen et al. (2025) to this long POMP model. For instance, if the perturbation schedule in Algorithm 1 is chosen to follow the dynamic approach introduced by the authors, then Theorem 1 is just an application of Theorem 4 of Chen et al. (2025) to a panel version of the model.

The stacking argument to prove Theorem 1 follows the approach of Bretó et al. (2020), who previously introduced the PIF algorithm and extended an existing theory for low-dimensional POMP models (Ionides et al., 2015) to derive theoretical properties of PIF. However, the recent work of Chen et al. (2025) provides convergence results for iterated filtering algorithms under weaker assumptions than Ionides et al. (2015). Most notably, Chen et al. (2025) prove convergence of iterated filtering algorithms as the random walk standard deviation for parameter perturbations decreases over time rather than being fixed at a small constant. Thus, this approach allows us to provide stronger theoretical results for panel iterated filtering than obtained by Bretó et al. (2020).

*Proof.* Recall the basic definition of the joint density of a PanelPOMP model:

$$f_{X_{0:N}, Y_{1:N}}(x_{0:N}, y_{1:N}; \theta) = \prod_{u=1}^U f_{X_{u,0}}(x_{u,0}; \theta) \prod_{n=1}^{N_u} f_{Y_{u,n}|X_{u,n}}(y_{u,n}|x_{u,n}; \theta) f_{X_{u,n}|X_{u,n-1}}(x_{u,n}|x_{u,n-1}; \theta). \quad (10)$$

We would like to pivot the unit specific processes into a single long process, avoiding the need to loop over the unit  $u$ . We define  $S_{\tilde{u}} = \sum_{k=1}^{\tilde{u}} (N_k + 1)$  to be the sum total number of time-steps (observations + initialization) for units 1 up to unit  $\tilde{u}$ , defining  $S_0 = 0$ . We use

the sequence  $n' \in \mathbb{N}$  to map time points and states indexed with  $(u, n) \in 1:U \times 0:N_u$  to a new model with only a single index  $n' \in 1:S_U$ , defined by the equation  $n' = S_{u-1} + n + 1$

Now let  $T_{\tilde{u}} = \sum_{k=1}^{\tilde{u}} (t_{k,N_k} - t_{k,0})$  for  $\tilde{u} \in 1:U$ , with  $T_0 = 0$ . We now define new latent processes  $\tilde{X}$  and  $\tilde{Y}$ . As the original model allows for continuous time latent process, we write First,

$$\tilde{X}(t) = X_u(t_{u,0} + (t - T_u)), \text{ for } T_{u-1} \leq t \leq T_u. \quad (11)$$

Letting  $\tau_{n'} = t_{u,n}$ , the latent and observable processes are equivalent to

$$\tilde{X}_{n'} = \tilde{X}(\tau_{n'}) = X_{u,n}, \quad \text{and} \quad \tilde{Y}_{n'} = Y_{u,n}.$$

With this new definition of  $\tilde{X}$  and  $\tilde{Y}$ , we have a new POMP model which has joint density

$$\begin{aligned} f_{\tilde{X}_{1:S_U}, \tilde{Y}_{1:S_U}}(x_{1:S_U}, y_{1:S_U}; \theta) = \\ f_{\tilde{X}_0}(x_0; \theta) \prod_{n'=1}^{S_U} f_{\tilde{Y}_{n'}|\tilde{X}_{n'}}(y_{n'}|x_{n'}; \theta) f_{\tilde{X}_{n'}|\tilde{X}_{n'-1}}(x_{n'}|x_{n'-1}; \theta), \end{aligned} \quad (12)$$

Using the convention that  $n' \mapsto (u, 0)$  for any  $u$ , then  $f_{\tilde{Y}_{n'}|\tilde{X}_{n'}}(y_{n'}|x_{n'}; \theta) = 1$  and  $f_{\tilde{X}_{n'}|\tilde{X}_{n'-1}}(x_{n'}|x_{n'-1}; \theta) = f_{X_{u,0}}(x_{u,0}; \theta)$ .

As defined, Eqs. 12 and 10 describe the same model, but 12 has been written to match the SSM of Chen et al. (2025), after adjusting for the choice of initializing at  $\tilde{X}_1$  rather than  $\tilde{X}_0$ . We refer to Eq 12 as the *long* format, which is indicative that the model has been expressed as a low-dimensional POMP model with observation times ranging from 1 to  $S_U$ , rather than a collection (or product) of POMP models, each with observation times  $N_u + 1$  for  $u \in 1:U$ .



This representation allows us to naturally extend the theoretical framework of Chen et al. (2025) to PanelPOMP models. Specifically, Assumption (A1) imply that the set  $\Theta$  is a regular compact set (See Definition 1 of Chen et al., 2025); Assumptions (B2)–(B3) ensure that the density in Eq. 12 corresponds to a state-space model that satisfies the MLE conditions of Chen et al. (2025), and Assumption (B1) is used as a uniformity condition across particle representations. Finally, Assumptions (C1)–(C4) are applied to the cloned version of Model 12. The process is cloned such that we have a new SSM  $\tilde{\tilde{Y}}_{\tilde{n}} = Y_{u,n}$ ,  $\tilde{\tilde{X}}_{\tilde{n}} = X_{u,n}$ , using the mapping  $\tilde{n} = (m - 1)S_U + S_{u-1} + n + 1$ . Assumptions (C1)–(C4) therefore imply assumptions C1-C3 and C4' of Chen et al. (2025) for the perturbation kernels of the self organized SSM defined via  $(\tilde{\tilde{Y}}_{\tilde{n}}, \tilde{\tilde{X}}_{\tilde{n}})$ . An alternative version of Assumption (C4) is also stated in Appendix A.5 of Chen et al. (2025), but here we only present one version for brevity. Together, the conditions stated in Appendix A allow for a direct application of of the theory developed by Chen et al. (2025, Appendix A), with precise details of applying this theorem to the cloned model given in Chen et al. (Supplement S7, 2025).  $\square$

Formally, Theorem 1 provides guarantees for several variants of iterated filtering algorithms applied to panel models, where each variant is a change to the perturbation kernel that satisfy the conditions (C1)–(C4). For example, Theorem 4 of Chen et al. (2025) is stated for specific perturbation kernels, such as those based on a normal distribution—a choice that has been used by default in most iterated filtering applications (e.g., Ionides et al., 2015; Fox et al., 2022; Subramanian et al., 2021; Wheeler et al., 2024). This theory also applies to the variant of iterated filtering proposed by Chen et al. (2025) where perturbations are only applied when needed in a dynamic fashion.

Conditions (C1)–(C4) are difficult to verify in practice, and not all variants of the algorithm that satisfy this condition are useful for panel models. For instance, multivariate

normal perturbations applied at each time point is equivalent to applying IF2 to panel models; this approach generally leads to worse results than the PIF algorithm applied to the same model (Bretó et al., 2020). As pointed out in Section 6, the effect of the perturbation kernel is twofold. First, it helps revive particle representations of the intermediate parameter distributions by adding random noise. Second, the random noise results in a loss of information by masking the signal from the observed data. These competing interests are not addressed in theorems involving iterated filtering algorithms. In this case, heuristics are useful for determining suitable perturbation densities.

The modification of Chen et al. (2025) that applies perturbations only when needed is an effective approach to address the tradeoff between these competing interests. For high-dimensional panel models, however, this modification still leads to a large number of updates of the particle representation of the parameter vector  $\psi_{-u}$  using data from unit  $u$ , which is the primary reason that algorithms like PIF struggle in higher dimensional settings. This observation leads to the proposal of the MPIF algorithm, which avoids resampling parameters if little information via the likelihood function. Combining the dynamic perturbations strategies of Chen et al. (2025) with a marginalization step may also result in an improved algorithm in some cases. However, we find that the default multivariate normal perturbations with singular covariance matrix to avoid perturbing all parameters at each time step to be sufficient for parameter estimation in practice.

One reason that Theorem 1 cannot be used directly to infer the practicality of iterated filtering algorithms for panel data is because the behavior of  $C_M$  as  $M \rightarrow \infty$  is unknown. Previous works on high-dimensional particle filtering—without adding perturbations or performing data cloning—suggest that the sequence  $C_M$  scales exponentially with the number of units (Snyder et al., 2008; Bengtsson et al., 2008). While this problem has partially been

avoided by writing the PanelPOMP in a long format, which reduces the size of both the latent and observed spaces at each time point, the parameter space for  $\Theta$  remains large. Specifically, the total dimension of  $\Theta$  is  $d_\phi + Ud_\psi$ , where  $d_\phi$  and  $d_\psi$  are the number of shared and unit-specific parameters, respectively.

Because particle filters are known to perform poorly in high-dimensions, alternative filtering algorithms should be used. One such example is called the block particle filter (Rebeschini and Van Handel, 2015); this algorithm breaks the state-space into separate units called *blocks*, and performs the update step in the filtering equation independently on each block. Block particle filters have been found to be effective in high-dimensional settings where the units can be treated as approximately independent. These results also provides an alternative motivation and justification for the MPIF algorithm, and a potential avenue for expanding Theorem 2 to more general state space models. In the panel model setting, the MPIF algorithm takes a similar approach to the block particle filter by updating the filtering distribution over the independent units, though the blocking occurs in the parameters space  $\Theta$  rather than the latent space  $\mathcal{X}$ . Thus, the MPIF algorithm has many similarities to other iterated block particle filtering algorithms that have been effective for moderately sized dynamic systems with spatial coupling (Ning and Ionides, 2023; Ionides et al., 2024).

## C Proof of Theorem 2

As a reminder, we assume that there are  $U \geq 2$  units, labeled  $1:U$ , with the data for unit  $u$  being denoted as  $\mathbf{y}_u^*$ . We write  $\theta = (\phi, \psi_1, \dots, \psi_U)$ , and recall that the unit likelihood  $L_u(\theta; \mathbf{y}_u^*)$  for unit  $u \in 1:U$  depends only on the shared parameter  $\phi$  and the unit-specific parameter  $\psi_u$ . The panel assumption implies that, conditioned on the parameter vector, the

units are dynamically independent. Therefore, we can decompose the likelihood function for the entire collection of data  $L(\theta; \mathbf{y}^*)$  as:

$$L(\theta; \mathbf{y}^*) = \prod_{u=1}^U L_u(\theta; \mathbf{y}_u^*) = \prod_{u=1}^U L_u(\phi, \psi_u; \mathbf{y}_u^*).$$

Each iteration of the Eqs. 5 and 6 corresponds to a Bayes update followed by a marginalization. Under the statement of Theorem 2, the prior and likelihood are assumed to be Gaussian. In this setting, it is well known that the resulting Bayes posterior also corresponds to a Gaussian distribution. Similarly, the marginalization of a multivariate Gaussian distribution also results in multivariate Gaussian distributions. Thus, each iteration of Eqs. 5 and 6 results in a density that corresponds to a Gaussian distribution.

We now show that the mean of the resulting Gaussian distribution converges to the MLE, while the covariance matrix converges to the zero matrix, resulting in a density with all mass centered at the MLE. To aid this calculation, we introduce the following lemma.

**Lemma 1.** *Let  $d$  be a positive integer, and let  $B_k \in \mathbb{R}^{2 \times 2}$  for  $k \in 1:d$  be a collection of real-valued matrices. We construct a sequence of matrices  $A_k \in \mathbb{R}^{d+1 \times d+1}$  such that:*

$$[A_k]_{i,j} = \begin{cases} [B_k]_{1,1}, & i = j = 1 \\ [B_k]_{1,2}, & i = 1, j = k + 1 \\ [B_k]_{2,1}, & i = k + 1, j = 1 \\ [B_k]_{2,2}, & i = j = k + 1 \\ 1, & i = j, i \notin \{1, k + 1\} \\ 0, & \text{otherwise} \end{cases}$$

*That is,  $A_k$  is a perturbation of the identity matrix, where the first and  $(k + 1)$ th row and column have been modified on the diagonal and on their off-diagonal intersection to match*

the matrix  $B_k$ . If for all  $k \in 1:d$ ,  $\|B_k\|_2 \leq c$  for some constant  $0 < c \leq 1$ , then

$$\left\| \prod_{k=1}^d A_k \right\|_2 \leq c. \quad (13)$$

*Proof of Lemma 1.* For  $i, j \in 1:(d+1)$ , denote  $\mu_{(i,j)} \in \mathbb{R}^2$  as the sub-vector of  $\mu \in \mathbb{R}^{d+1}$  that contains only the  $i$  and  $j$ th elements, and write  $\mu_{-(i,j)} \in \mathbb{R}^{d-1}$  to be the sub-vector of  $\mu$  after removing these elements. By design, the matrix  $A_k$  operates only on the sub-vector  $\mu_{(1,k+1)}$ . That is, if  $B_k \mu_{(1,k+1)} = (\tilde{\mu}_{(1)}, \tilde{\mu}_{(k+1)})^T$ , then  $A_k \mu = (\tilde{\mu}_{(1)}, \mu_{(2)}, \dots, \tilde{\mu}_{(k+1)}, \dots, \mu_{(d+1)})$ . From this, we see that for positive integer  $m \leq d$ , the first  $m$  factors of the product  $\prod_{k=1}^d A_k$  modify only the first  $m+1$  dimensions of a vector  $\mu \in \mathbb{R}^{d+1}$ .

We proceed by mathematical induction on the dimension size. Let  $\{A_k^{(d)}, k \in 1:d-1\}$  be a collection of matrices satisfying the condition of Lemma 1 for each  $d$ . Setting  $P_d = \left\| \prod_{k=1}^{d-1} A_k^{(d)} \right\|_2$ , we first observe that Eq. (13) holds for  $d=2$  as a direct consequence of the condition  $\|B_k\|_2 \leq c$ . Suppose inductively that the lemma holds for  $d$ , so that  $P_d \leq c$ . We wish to bound  $P_{d+1}$ . We can choose  $A_k^{(d)}$  to be the  $(d+1) \times (d+1)$  sub-matrix of  $A_k^{(d+1)}$  omitting row and column  $(d+2)$ . Thus, for  $k \in 1:d$ ,

$$A_k^{(d+1)} = \begin{pmatrix} A_k^{(d)} & 0 \\ 0 & 1 \end{pmatrix}.$$

Consider a vector  $\mu \in \mathbb{R}^{d+1}$ , such that  $\|\mu\|_2 \leq 1$ . Let  $\tilde{\mu} \in \mathbb{R}^d$  be defined by

$$\tilde{\mu} = \left[ \left( \prod_{k=1}^{d-1} A_k^{(d+1)} \right) \mu \right]_{(1:d)} \quad (14)$$

and notice that we have

$$\tilde{\mu} = \left( \prod_{k=1}^{d-1} A_k^{(d)} \right) \mu_{(1:d)}. \quad (15)$$

By construction, we have

$$\begin{aligned}
\left| \left( \prod_{k=1}^d A_k^{(d+1)} \right) \mu \right|_2^2 &= \left| \left( A_d^{(d+1)} \prod_{k=1}^{d-1} A_k^{(d+1)} \right) \mu \right|_2^2 \\
&= \left| A_d^{(d+1)} (\tilde{\mu}_{(1:d)}, \mu_{(d+1)})^T \right|_2^2 \\
&= \left| B_d^{(d+1)} (\tilde{\mu}_1, \mu_{(d+1)})^T \right|_2^2 + \left| \tilde{\mu}_{(1:d)} \right|_2^2 - \tilde{\mu}_1^2.
\end{aligned}$$

Because  $\|B_d^{(d+1)}\|_2 \leq c$ ,  $|B_d^{(d+1)}(\tilde{\mu}_1, \mu_{(d+1)})^T|_2^2 \leq c^2 |\tilde{\mu}_1, \mu_{(d+1)}|^2_2$ . Furthermore, our inductive hypothesis applied to Eq. (15) implies that  $|\tilde{\mu}_{(1:d)}|_2^2 \leq c^2 |\mu_{(1:d)}|_2^2 \leq c^2 |\mu|_2^2$ . Therefore,

$$\left| \left( \prod_{k=1}^d A_k \right) \mu \right|_2^2 \leq c^2 (\tilde{\mu}_1^2 + \mu_{(d+1)}^2) + c^2 (|\mu_{(1:d)}|_2^2) - \tilde{\mu}_1^2 \tag{16}$$

$$= (c^2 - 1) \tilde{\mu}_1^2 + c^2 |\mu|_2^2 \tag{17}$$

$$\leq c^2 |\mu|_2^2, \tag{18}$$

with Eq. (18) implied by  $c \leq 1$ , and hence  $(c^2 - 1) \leq 0$ . It follows immediately from Eq. (18) that  $P_{d+1} \leq c$ , completing the proof.  $\square$

We now return to the main argument.

*Proof of Theorem 2.* Using the transformation invariance of the MLE, we can suppose without loss of generality that the maximum of the marginal likelihood for unit  $u$  is at  $\phi, \psi_u = 0$ . To help ease notation, we write the covariance matrix as

$$\Lambda_u^* = \begin{pmatrix} \Lambda_\phi^{(u)} & \Lambda_{\phi,u} \\ \Lambda_{\phi,u} & \Lambda_u \end{pmatrix}.$$

Let the prior density  $\pi_0(\theta)$  correspond to a Gaussian distribution with mean  $\mu_0 \in R^{U+1}$

and precision  $\Gamma_0 = \Sigma_0^{-1} \in \mathbb{R}^{U+1 \times U+1}$ ,

$$\mu_0 = \begin{pmatrix} \mu_0^{(\phi)} \\ \mu_0^{(1)} \\ \vdots \\ \mu_0^{(U)} \end{pmatrix}, \quad \Gamma_0 = \begin{pmatrix} \tau_0^{(\phi)} & 0 & \dots & 0 \\ 0 & \tau_0^{(1)} & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \tau_0^{(U)} \end{pmatrix}.$$

Eqs. 5 and 6 contain two indices  $(m, u)$  for the intermediate density function  $\pi_{m,u}(\theta)$ . The first index ( $m \in \mathbb{N}$ ) counts the number of times data from all units has been used in the Bayes update (Eq. 5), and this corresponds to the number of complete iterations of the  $m$  loop in Algorithm 1. The second index ( $u \in 1:U$ ) denotes the data from which unit is currently being used to update the density, and we refer to this as a sub-iteration. In what follows, we assume that  $\pi_{m,u-1}(\theta) = \pi_{m-1,U}(\theta)$  if  $u = 1$ , and use a similar convention for the corresponding mean and covariance that correspond to this intermediate density.

The density after each sub-iteration of Eqs. 5 and 6 is Gaussian, and the marginalization step ensures that the precision matrix from previous sub-iterations is diagonal. Using  $\mu_{m,u-1}$  and  $\Gamma_{m,u-1}$  to denote the mean and precision after  $m$  complete iterations and the  $(u-1)$ th unit-iteration, we write:

$$\mu_{m,u-1} = \begin{pmatrix} \mu_{m,u-1}^{(\phi)} \\ \mu_{m,u-1}^{(1)} \\ \vdots \\ \mu_{m,u-1}^{(U)} \end{pmatrix}, \quad \Gamma_{m,u-1} = \begin{pmatrix} \tau_{m,u-1}^{(\phi)} & 0 & \dots & 0 \\ 0 & \tau_{m,u-1}^{(1)} & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \tau_{m,u-1}^{(U)} \end{pmatrix}.$$

By design, each sub-iteration  $u$  only modifies the elements of the mean and precision that correspond to the shared parameter  $\phi$  and the unit specific parameter  $\psi_u$ . Thus, we use

the superscript notation  $\mu_{m,u}^{(1,k)}$  and  $\Gamma_{m,u}^{(1,k)}$  to denote the components of the vector  $\mu_{m,u}$  corresponding to parameters  $\phi$  and  $\psi_k$ , and the  $2 \times 2$  submatrix of  $\Gamma_{m,u}$  with the elements corresponding to parameters  $\phi$  and  $\psi_k$ .

Performing the Bayes update to this distribution (Eq. 5) results in an unmarginalized precision matrix, which we denote as  $\tilde{\Gamma}_{m,u}$ , where the only modified components are

$$\tilde{\Gamma}_{m,u}^{(m,u)} = \Gamma_{m,u-1}^{(m,u)} + \Lambda_u^* = \begin{pmatrix} \tau_{m,u-1}^{(\phi)} + \Lambda_\phi^{(u)} & \Lambda_{\phi,u} \\ \Lambda_{\phi,u} & \tau_{m,u-1}^{(u)} + \Lambda_u \end{pmatrix}.$$

The corresponding mean  $\tilde{\mu}_{(m,u)} = \mu_{(m,u)}$ , noting that the marginalization procedure (Eq. 6) does not affect the mean, remains the same as the previous iteration except for the components

$$\begin{aligned} \tilde{\mu}_{m,u}^{(\phi,u)} &= (\tilde{\Gamma}_{m,u}^{(m,u)})^{-1} (\Gamma_{m,u-1}^{(m,u)} \mu_{m,u-1}^{(\phi,u)} + \Lambda_u^* (0, 0)^T) \\ \mu_{m,u}^{(\phi,u)} &= (\tilde{\Gamma}_{m,u}^{(\phi,u)})^{-1} \Gamma_{m,u-1}^{(\phi,u)} \mu_{m,u-1}^{(\phi,u)} \\ &= B_{m,u} \mu_{m,u-1}^{(\phi,u)}, \end{aligned} \tag{19}$$

where  $B_{m,u} = (\tilde{\Gamma}_{m,u}^{(\phi,1)})^{-1} \Gamma_{m,u-1}^{(\phi,u)}$  is a  $2 \times 2$  matrix that provides the update to the first and  $(u+1)$ th components of a vector  $\mu \in \mathbb{R}^{U+1}$ . Now by defining matrix  $A_{m,u}$  to be a perturbation of the  $U+1$  identity matrix, such that

$$[A_{m,u}]_{i,j} = \begin{cases} [B_{m,u}]_{1,1} & i = j = 1 \\ [B_{m,u}]_{2,2} & i = j = u+1 \\ [B_{m,u}]_{1,2} & i = 1, j = k+1 \\ [B_{m,u}]_{2,1} & i = k+1, j = 1 \\ 1 & i = j, i \notin \{1, k+1\} \\ 0 & \text{otherwise} \end{cases},$$



We see that for a given precision matrix  $\Gamma_{m,u-1}$ , the a single sub-iteration of Eqs. 5 and 6 corresponds to a linear update of the mean vector  $\mu_{m,u-1}$ , defined by

$$\mu_{m,u} = A_{m,u} \mu_{m,u-1}.$$

Thus, the mean after  $M$  complete iterations, denoted by  $\mu_M$  is calculated by the product of these linear transformations

$$\mu_M = \left( \prod_{m=1}^M \prod_{u=1}^U A_{m,u} \right) \mu_0.$$

Thus, the long-term value of  $\mu_M$  depends on the long-term behavior of the product of matrices  $A_{m,u}$ .

We now consider the limiting behavior of the precision matrices, which determines the behavior of the matrices  $A_{m,u}$ . First, we recall that the unmarginalized precision matrix is of the form:

$$\tilde{\Gamma}_{m,u}^{(m,u)} = \begin{pmatrix} \tau_{m,u-1}^{(\phi)} + \Lambda_{\phi}^{(u)} & \Lambda_{\phi,u} \\ \Lambda_{\phi,u} & \tau_{m,u-1}^{(u)} + \Lambda_u \end{pmatrix}.$$

The marginalization step (Eq. 6) for Gaussian densities is represented by setting the off-diagonal elements of the covariance matrix to be zero. Thus, if we generalize this step by writing  $\tilde{\Gamma}_{m,u}^{(m,u)} = \begin{pmatrix} a & b \\ b & d \end{pmatrix}$ , then  $(\tilde{\Gamma}_{m,u}^{(m,u)})^{-1} = \frac{1}{ad-b^2} \begin{pmatrix} d & -b \\ -b & a \end{pmatrix}$  and the marginalized version of the matrix is  $(\Gamma_{m,u}^{(m,u)})^{-1} = \begin{pmatrix} \frac{d}{ad-b^2} & 0 \\ 0 & \frac{a}{ad-b^2} \end{pmatrix}$ ; by taking the inverse of this matrix, we get the precision matrix after a unit-update to be  $\Gamma_{m,u}^{(m,u)} = \begin{pmatrix} a - \frac{b^2}{d} & 0 \\ 0 & d - \frac{b^2}{a} \end{pmatrix}$ .

Using this general calculation, the precision matrix after the marginalization step is

$$\begin{aligned}\Gamma_{m,u}^{(\phi,u)} &= \begin{pmatrix} \tau_{m,u-1}^{(\phi)} + \Lambda_{\phi}^{(u)} - \frac{\Lambda_{\phi,u}^2}{\tau_{m,u-1}^{(u)} + \Lambda_u} & 0 \\ 0 & \tau_{m,u-1}^{(u)} + \Lambda_u - \frac{\Lambda_{\phi,u}^2}{\tau_{m,u-1}^{(\phi)} + \Lambda_{\phi}^{(u)}} \end{pmatrix} \\ &= \begin{pmatrix} \tau_{m,u-1}^{(\phi)} + \Lambda_{\phi}^{(u)} - \alpha_{m,u} & 0 \\ 0 & \tau_{m,u-1}^{(u)} + \Lambda_u - \beta_{m,u} \end{pmatrix},\end{aligned}\tag{20}$$

with  $\alpha_{m,u} = \frac{\Lambda_{\phi,u}^2}{\tau_{m,u-1}^{(u)} + \Lambda_u}$  and  $\beta_{m,u} = \frac{\Lambda_{\phi,u}^2}{\tau_{m,u-1}^{(\phi)} + \Lambda_{\phi}^{(u)}}$ . Because  $\Lambda_u^*$  is a positive definite matrix, for all real numbers  $c > 0$ ,

$$\Lambda_{\phi}^{(u)} > \frac{\Lambda_{\phi,u}^2}{\Lambda_u} > \frac{\Lambda_{\phi,u}^2}{c + \Lambda_u}$$

And therefore by letting  $\alpha = \min_u (\Lambda_u - \frac{\Lambda_{\phi,u}^2}{\Lambda_{\phi}^{(u)}}) > 0$ , we have  $\tau_{m,u}^{(\phi)} > \tau_{m,u-1}^{(\phi)} + \alpha$  for all  $m \in \mathbb{N}$  and  $u \in 1:U$ . By iterating this inequality, we have  $\tau_{m,u}^{(\phi)} > \tau_0^{(\phi)} + m\alpha$ , and thus we see that  $\tau_{m,u}^{(\phi)} = O(m)$ . Therefore as  $m \rightarrow \infty$ ,  $\tau_{m,u}^{(\phi)} \rightarrow \infty$ . A similar calculation shows that  $\tau_{m,k}^{(u)} = O(m)$  for all  $k, u \in 1:U$ . Thus, the covariance matrix after  $m$  complete iterations has zeros on off-diagonal elements, and the diagonal elements are of order  $O(1/m)$ , proving the statement in Theorem 2 that  $\|\Sigma_m\|_2 \rightarrow 0$ . To finish the proof, we need to show that  $|\mu_m|_2 \rightarrow 0$ . To do this, we need more precise descriptions for the rates at which the precision grows.

Our previous calculations establish that both  $\tau_{m,k}^{(\phi)}$  and  $\tau_{m,k}^{(u)}$  are strictly increasing sequences in both  $m$  and  $k$ , and are unbounded. Thus, the correction terms  $\alpha_{m,u}$  and  $\beta_{m,u}$  converge to zero, as the only moving parts are the precision terms that are in the denomi-

nator of each of these sequences. Furthermore, the sequence  $\{\alpha_i\}$ , defined as

$$\alpha_i = \sum_{u=1}^U \alpha_{i,u}$$

satisfies  $\alpha_i \rightarrow 0$ . Thus, by iterating Eq. 20, we can express the precision matrix after  $m$  complete iterations as

$$\Gamma_m^{(\phi,u)} = \Gamma_{m-1,U}^{(\phi,u)} = \begin{pmatrix} \tau_0^{(\phi)} + m \sum_{k=1}^U \Lambda_\phi^{(k)} - \sum_{i=1}^m \alpha_i & 0 \\ 0 & \tau_0^{(u)} + m \Lambda_k - \sum_{i=1}^m \beta_{i,u} \end{pmatrix}.$$

Using the fact that the Cesàro mean of a convergent sequence converges to the same limit as the sequence, we have  $\frac{1}{m} \sum_{i=1}^m \alpha_i \rightarrow 0$ , and therefore

$$\Gamma_{m,u}^{(\phi,k)} = \begin{pmatrix} m \sum_{k=1}^U \Lambda_\phi^{(k)} + o(m) & 0 \\ 0 & m \Lambda_k + o(m) \end{pmatrix}. \quad (21)$$

Using the rates established in Eq. 21, we now investigate the long-term behavior of the mean vector  $\mu_m$  by considering the spectral norm of the matrices  $B_{m,u}$  that define the linear updates to the sub-vector  $\mu_m^{(\phi,u)}$ . Recall from Eq. 19 that  $B_{m,u} = (\Gamma_{m,u-1}^{(\phi,u)} + \Lambda_u^*)^{-1} \Gamma_{m,u-1}^{(\phi,u)}$ .

We can take advantage of the fact that  $\Gamma_{m,u-1}^{(\phi,u)}$  is a diagonal matrix to write

$$\begin{aligned}
B_{m,u} &= (\Gamma_{m,u-1}^{(\phi,u)} + \Lambda_u^*)^{-1} \Gamma_{m,u-1}^{(\phi,u)} \\
&= \left( I + (\Gamma_{m,u-1}^{(\phi,u)})^{-1} \Lambda_u^* \right)^{-1} \\
&= \begin{pmatrix} 1 + \frac{\Lambda_\phi^{(u)}}{m \sum_{k=1}^U \Lambda_\phi^{(k)} + o(m)} & \frac{\Lambda_{\phi,u}}{m \sum_{k=1}^U \Lambda_\phi^{(k)} + o(m)} \\ \frac{\Lambda_{\phi,u}}{m \Lambda_u + o(m)} & 1 + \frac{\Lambda_u}{m \Lambda_u + o(m)} \end{pmatrix}^{-1} \\
&= \begin{pmatrix} 1 + \frac{\Lambda_\phi^{(u)}}{m \sum_{k=1}^U \Lambda_\phi^{(k)}} + o(1/m) & \frac{\Lambda_{\phi,u}}{m \sum_{k=1}^U \Lambda_\phi^{(k)}} + o(1/m) \\ \frac{\Lambda_{\phi,u}}{m \Lambda_u} + o(1/m) & 1 + \frac{1}{m} + o(1/m) \end{pmatrix}^{-1}.
\end{aligned}$$

Next we would like to calculate the spectral norm of  $B_{m,u}$ . Recall that for an invertible matrix  $A$ ,  $\|A^{-1}\|_2 = \sigma_{\max}(A^{-1}) = 1/\sigma_{\min}(A)$ , where  $\sigma_{\max}$  and  $\sigma_{\min}$  correspond to the maximum and minimum singular values of  $A$ . Therefore in order to calculate  $\|B_{m,u}\|_2$ , we need to find the minimum singular value of  $B_{m,u}^{-1}$ , which is equal to the minimum eigenvalue of the matrix  $B_{m,u}^{-T} B_{m,u}^{-1}$ :

$$\begin{aligned}
B_{m,u}^{-T} B_{m,u}^{-1} &= \begin{pmatrix} 1 + \frac{2\Lambda_\phi^{(u)}}{m \sum_{k=1}^U \Lambda_\phi^{(k)}} + o(1/m) & \frac{\Lambda_{\phi,u}}{m} \left( \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right) + o(1/m) \\ \frac{\Lambda_{\phi,u}}{m} \left( \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right) + o(1/m) & 1 + \frac{2}{m} + o(1/m) \end{pmatrix} \\
&= I + \frac{1}{m} \begin{pmatrix} \frac{2\Lambda_\phi^{(u)}}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + o(1) & \Lambda_{\phi,u} \left( \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right) + o(1) \\ \Lambda_{\phi,u} \left( \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right) + o(1) & 2 + o(1) \end{pmatrix} \\
&= I + \frac{1}{m} C_{m,u}.
\end{aligned}$$

Using this expression, the eigenvalues of  $B_{m,u}^{-T} B_{m,u}^{-1}$  are equal to one plus the eigenvalues of  $\frac{1}{m} C_{m,u}$ . Thus, we consider the characteristic polynomial defined by  $\det(C_{m,u} - \lambda I)$  to

calculate the eigenvalues of  $C_{m,u}$ . The resulting polynomial is

$$\begin{aligned} f_{m,u}(\lambda) &= \left( \frac{2\Lambda_\phi^{(u)}}{\sum_{k=1}^U \Lambda_\phi^{(k)}} - \lambda \right) (2 - \lambda) - \Lambda_{\phi,u}^2 \left( \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right)^2 + o(1) \\ &= \lambda^2 - 2 \left( 1 + \frac{\Lambda_\phi^{(u)}}{\sum_{k=1}^U \Lambda_\phi^{(k)}} \right) \lambda + \left\{ \frac{4\Lambda_\phi^{(u)}}{\sum_{k=1}^U \Lambda_\phi^{(k)}} - \Lambda_{\phi,u}^2 \left( \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right)^2 \right\} + o(1). \end{aligned}$$

We now proceed by showing that, under the condition specified in Assumption 7, all roots of the polynomial  $f_{m,u}(\lambda)$  are strictly positive. That is, for some  $0 < \lambda_{m,u}^{(1)} < \lambda_{m,u}^{(2)}$ , we have the eigenvalues of  $C_{m,u}$  are  $\lambda_{m,u}^{(1)} + o(1)$  and  $\lambda_{m,u}^{(2)} + o(1)$ .

For this to hold, we need first that the discriminant of the resulting quadratic equation to be strictly positive. This condition holds without any assumptions, as we can write the condition as

$$\begin{aligned} \left( 2 \left[ 1 + \frac{\Lambda_\phi^{(u)}}{\sum_{k=1}^U \Lambda_\phi^{(k)}} \right] \right)^2 - 4 \left( \frac{4\Lambda_\phi^{(u)}}{\sum_{k=1}^U \Lambda_\phi^{(k)}} - \Lambda_{\phi,u}^2 \left[ \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right]^2 \right) &> 0 \iff \\ \Lambda_{\phi,u}^2 \left( \frac{1}{\sum_{k=1}^U \Lambda_\phi^{(k)}} + \frac{1}{\Lambda_u} \right)^2 &> - \frac{(\sum_{k=1}^U \Lambda_\phi^{(k)} - \Lambda_\phi^{(u)})^2}{(\sum_{k=1}^U \Lambda_\phi^{(k)})}. \end{aligned}$$

This condition always holds as the left-hand side of the inequality is strictly positive, and the right-hand side of the inequality is strictly negative. Thus, the polynomial  $f_{m,u}(\lambda)$  has two unique real roots.

Next, by taking the derivative and setting equal to zero, we note that  $\arg \min_\lambda f_{m,u}(\lambda) = \lambda^* > 0$ , as the coefficient on the linear term of the polynomial is always negative. This combined with the previous condition ensures that  $\lambda_{m,u}^{(2)} > 0$ .

Finally, if  $f_{m,u}(0) > 0$ , then the intermediate value theorem implies that because  $f_{m,u}(\lambda^*) < 0$ , then there exists some  $\lambda_{m,u}^{(1)} \in (0, \lambda^*)$  such that  $f_{m,u}(\lambda_{m,u}^{(1)}) = 0$ , implying

that  $\lambda_{m,u}^{(1)}$  is a positive root. Therefore we need

$$\frac{4\Lambda_{\phi}^{(u)}}{\sum_{k=1}^U \Lambda_{\phi}^{(k)}} - \Lambda_{\phi,u}^2 \left( \frac{1}{\sum_{k=1}^U \Lambda_{\phi}^{(k)}} + \frac{1}{\Lambda_u} \right)^2 > 0,$$

which is equivalent to the condition

$$\Lambda_{\phi,u}^2 < \frac{4\Lambda_{\phi}^{(u)} \Lambda_u^2 \sum_{k=1}^U \Lambda_{\phi}^{(k)}}{(\Lambda_u + \sum_{k=1}^U \Lambda_{\phi}^{(k)})^2}.$$

This condition is guaranteed by Assumption 7, ensuring that the eigenvalues of  $C_{m,u}$  are  $\lambda_{m,u}^{(1)} + o(1)$  and  $\lambda_{m,u}^{(1)} + o(1)$  for two positive numbers  $\lambda_{m,u}^{(1)}$  and  $\lambda_{m,u}^{(1)}$ . As a result of this computation, The minimum eigenvalue of  $B_{m,u}^{-T} B_{m,u}^{-1}$  is equal to  $1 + \frac{\lambda_{m,u}^{(1)}}{m} + o(1/m)$ , implying that

$$\begin{aligned} \|B_{m,u}\|_2 &= 1/\sigma_{\min}(B^{-1}) \\ &= \frac{1}{1 + 1 + \frac{\lambda_{m,u}^{(1)}}{m} + o(1/m)} \\ &= 1 - \frac{\lambda_{m,u}^{(1)}}{m} + o(1/m). \end{aligned}$$

By Lemma 1, we have for all  $m \in \mathbb{N}$ ,  $\|\prod_{u=1}^U A_{m,u}\|_2 \leq 1 - \frac{\max_u \lambda_{m,u}^{(1)}}{m} + o(1/m)$ , and therefore by the sub-multiplicative property of the spectral norm,

$$\begin{aligned} \left\| \prod_{i=1}^m \prod_{u=1}^U A_{i,u} \right\|_2 &\leq \prod_{i=1}^m \left\| \prod_{u=1}^U A_{i,u} \right\|_2 \\ &< \prod_{i=1}^m \left( 1 - \frac{\max_u \lambda_{i,u}^{(1)}}{i} + o(1/i) \right) \rightarrow 0. \end{aligned}$$

Therefore for any arbitrary initial conditions  $\mu_0$  and  $\Gamma_0$ , we have

$$|\mu_m|_2 = \left| \left( \prod_{i=1}^m \prod_{u=1}^U A_{i,u} \right) \mu_0 \right|_2 \rightarrow 0$$

as  $m \rightarrow \infty$ , completing the proof.  $\square$

## C.1 Proof of Corollary 1

For this corollary, we use the same setup as Appendix C. At each step, we add independent random perturbations to the current parameter distribution via a convolution operation before updating using Bayes rule. Because Gaussian density convolved with another Gaussian density is still Gaussian, we again only record the mean and precision corresponding to the resulting Gaussian distribution at each step, which we denote  $\mu'_m \in \mathbb{R}^{U+1}$  and  $\Gamma'_m \in \mathbb{R}^{(U+1) \times (U+1)}$ .

*Proof.* Let  $\mu_{m,u} \in \mathbb{R}^{U+1}$  and  $\Gamma_{m,u} = \text{diag}(\tau_{m,u}^{(\phi)}, \tau_{(m,u)}^{(1)}, \dots, \tau_{(m,u)}^{(U)}) \in \mathbb{R}^{(U+1) \times (U+1)}$  denote the mean and precision matrices of the multivariate Gaussian after the  $(m, u)$ th iteration. Consider the  $(m, u)$ th update of Eqs. 8–9. Using well-known results on the convolution of Gaussian densities and the conjugate prior identities, the unmarginalized precision is evaluated as

$$\tilde{\Gamma}'_{m,u+1} = \begin{pmatrix} \frac{\tau_{m,u}^{(\phi)} \tau_{m,u}^{(\sigma)}}{\tau_{m,u}^{(\phi)} + \tau_{m,u}^{(\sigma)}} + \Lambda_{\phi}^{(u)} & \Lambda_{\phi,u} \\ \Lambda_{\phi,u} & \frac{\tau_{m,u}^{(u)} \tau_{m,u}^{(\sigma)}}{\tau_{m,u}^{(u)} + \tau_{m,u}^{(\sigma)}} + \Lambda_u \end{pmatrix}.$$

Following the formula from Appendix C for finding the marginalized precision matrices,

the marginalized precision  $\Gamma'_{m,u}$  can be expressed as

$$\Gamma'_{m,u+1} = \begin{pmatrix} \frac{\tau_m^{(\phi)} \tau_m^{(\sigma)}}{\tau_m^{(\phi)} + \tau_m^{(\sigma)}} + \Lambda_\phi^{(u)} - \frac{\Lambda_{\phi,u}^2 (\tau_m^{(\sigma)} + \tau_m^{(u)})}{\tau_m^{(\sigma)} \tau_m^{(u)} + \Lambda_u (\tau_m^{(\sigma)} + \tau_m^{(u)})} & 0 \\ 0 & \frac{\tau_m^{(u)} \tau_m^{(\sigma)}}{\tau_m^{(u)} + \tau_m^{(\sigma)}} + \Lambda_u - \frac{\Lambda_{\phi,u}^2 (\tau_m^{(\sigma)} + \tau_m^{(\phi)})}{\tau_m^{(\sigma)} \tau_m^{(\phi)} + \Lambda_\phi^{(u)} (\tau_m^{(\sigma)} + \tau_m^{(\phi)})} \end{pmatrix}.$$

Because  $\Lambda_u^*$  is positive definite, the same argument in the unperturbed case (Appendix C)

gives the existence of positive constants  $\{\beta_u, \alpha_u\}_{u=1}^U$  such that

$$\begin{aligned} \tau_{m,u+1}^{(\phi)} &> \frac{\tau_m^{(\phi)} \tau_m^{(\sigma)}}{\tau_m^{(\phi)} + \tau_m^{(\sigma)}} + \alpha_u \\ \tau_{m,u+1}^{(u)} &> \frac{\tau_m^{(u)} \tau_m^{(\sigma)}}{\tau_m^{(u)} + \tau_m^{(\sigma)}} + \beta_u. \end{aligned}$$

Furthermore, the sequence defined by iterating the right hand side of these inequalities is unbounded, which is easily demonstrated by assuming it is bounded, and noting that this leads to a contradiction because  $1/\tau_m^{(\sigma)} = o(1/m)$ , and therefore each update grows by an amount arbitrarily close to the constants  $\alpha_u$  and  $\beta_u$ .

Immediately, this result implies that  $\tau_{m,u}^{(\phi)}, \tau_{m,u}^{(u)} \rightarrow \infty$  as  $m \rightarrow \infty$ . That is, the covariance matrix convergence to the zero matrix, even in the presence of parameter perturbations. In particular, we can write the correction term

$$\epsilon_{m,u} = \frac{\Lambda_{\phi,u}^2 (\tau_m^{(\sigma)} + \tau_m^{(u)})}{\tau_m^{(\sigma)} \tau_m^{(u)} + \Lambda_u (\tau_m^{(\sigma)} + \tau_m^{(u)})} = o(1),$$



and, using the fact that the sequence  $\tau_m^{(\sigma)}$  is defined independently of  $\tau_{m,u}^\phi$ , we can write

$$\begin{aligned}\frac{\tau_{m,u}^{(\phi)} \tau_m^{(\sigma)}}{\tau_{m,u}^{(\phi)} + \tau_m^{(\sigma)}} &= \frac{\tau_{m,u}^{(\phi)}}{\frac{1}{\tau_m^{(\sigma)}} (\tau_{m,u}^{(\phi)} + \tau_m^{(\sigma)})} \\ &= \tau_{m,u}^{(\phi)} \frac{1}{1 + o(1/m)} \\ &= \tau_{m,u}^{(\phi)} + o(1/m).\end{aligned}$$

Thus, the updated precision after a full iteration can be approximated as

$$\tau_{m+1}^{(\phi)} = \tau_m^\phi + m \sum_{u=1}^U \Lambda_\phi^{(u)} - \sum_{u=1}^U \epsilon_{m,u} + o(1/m),$$

and therefore

$$\tau_{m+1}^{(\phi)} = o(m) + m \sum_{u=1}^U \Lambda_\phi^{(u)},$$

and the same basic calculation shows that for all  $u$ ,

$$\tau_{m+1}^{(u)} = o(m) + m \Lambda_u.$$

Similar to the unperturbed case, the update to the mean can be expressed as  $\mu_{m+1} = (\prod_{u=1}^U A'_{m,u}) \mu_m$ , where the matrix  $A'_{m,u}$  is defined in the same way as  $A_{m,u}$  from Appendix C, after replacing the covariance matrices with the perturbed versions:

$$A'_{m,u} = (\Gamma'_{m,u-1} + \Lambda_u^*)^{-1} \Gamma'_{m,u-1}.$$

We have shown above that  $A'_{m,u}$  and  $A_{m,u}$  are asymptotically equivalent. Therefore the

same argument in the proof of Theorem 2 immediately implies

$$|\mu'_m|_2 = \left| \left( \prod_{n=1}^m \prod_{u=1}^U A'_{n,u} \right) \mu_0 \right|_2 \rightarrow 0.$$

□