

Toward automatically differentiable particle filters for phylodynamic inference

Workshop on Mathematical Ecology
Queen's University, Kingston, ON

Edward Ionides

pypomp developers: Kevin Tan, Jesse Wheeler, Jun Chen,
Bo Yang, Aaron Abkemeier, Aaron King, Kunyang He

University of Michigan

July 25, 2025

Outline

- ▶ What is automatic differentiation (AD)?
- ▶ Why is AD important for statistical inference? Existing evidence.
- ▶ What has delayed the impact of AD for inference on nonlinear stochastic dynamic systems?
- ▶ AD for simulation-based inference on mechanistic models: AD for particle filters using GPUs
- ▶ Toward AD for phylodynamic inference

Introduction to automatic differentiation (AD)

- ▶ AD is numerical differentiation where compiled code automatically computes a massive chain rule, taking advantage of exact expressions for the numerical derivative of basic operations ($+$, $-$, \times , \log , \exp , \sin , \cos , etc).
- ▶ AD does not involve symbolic algebra.
- ▶ Modern compilers make AD available for arbitrary programs.
- ▶ JAX, a Python extension, does AD with automatic parallelization for CPU or GPU.

Modern statistical methods powered by AD

- ▶ Deep learning (also uses GPU)
- ▶ Stan: Hamiltonian Markov chain Monte Carlo
- ▶ Prophet: automatic forecasting
- ▶ AD model builder (TMB): random effect models in ecology and fisheries

Conclusion: AD has allowed growth in model complexity and data size.

Mechanistic POMP models and the particle filter (PF, a.k.a. sequential Monte Carlo, SMC)

- ▶ Partially observed Markov process (POMP) models provide a framework for mechanistic modeling of dynamic systems with noisy measurements.
- ▶ PF uses a POMP model to generate an ensemble of forecasts for each successive observation.
- ▶ When the observation is recorded, it is “assimilated” into the forecast by resampling the ensemble members with weight proportional to the probability of the observation given the forecast for that ensemble member.
- ▶ Remarkably, these iterated forecast and correction steps provide an unbiased estimate of the likelihood for the model.
- ▶ PF just needs a simulator from the dynamic model, not transition probabilities. It is plug-and-play.

Why is AD not standard for the particle filter (PF) likelihood?

- ▶ Likelihood evaluation and maximization permits parameter estimation, confidence intervals, diagnostic testing, model selection.
- ▶ Derivatives help with maximization. They also help with Bayesian inference, via Hamiltonian Monte Carlo (HMC).
- ▶ So, why doesn't everyone use AD for PF?
- ▶ Technical difficulties
 1. Resampling is discontinuous: discrete and so piecewise constant.
 2. The derivative estimate can be numerically unstable even when the log-likelihood estimate is stable.

Hypothesis: PF + AD + GPU is a scientific opportunity

- ▶ Particle filter methods have proved effective where they are applicable because
 - ▶ They approximate the actual likelihood: statistical efficiency.
 - ▶ Simulation-based “plug-and-play” structure is convenient for scientific exploration.
- ▶ Deep learning methods use AD + GPU.
 - ▶ Artificial neurons with arbitrary weight propagation rules should be less pertinent than particles that directly simulate the dynamic process and are properly weighted to estimate the likelihood.

Previous attempts at AD for PF

- ▶ Approximating the particle filter using deep learning.
- ▶ Ignoring the resampling when calculating the derivative.
- ▶ Obtaining a smooth particle filter at the cost of inferior scaling and loss of the plug-and-play property.
- ▶ AD for smooth functions estimated by Monte Carlo expectation of non-smooth functions has been widely studied. Applied to PF, previous methods have struggled with stability for a long time series.

We looked for a plug-and-play particle filter equipped with automatic differentiation that maintains the guarantee of unbiased likelihood evaluation and provides numerically stable derivative estimates.

Differentiated Measurement Off Parameter (DMOP) filter

- ▶ Numerically stable automatic derivatives, with controlled bias and variance, from a plug-and-play particle filter with unbiased likelihood evaluation Tan et al. (2024).

Measurement Off-Parameter (MOP- α)

- ▶ A MOP particle filter at parameter θ is a smooth continuation of the filter at ϕ , fixing resampling decisions for ϕ and accounting for this via “off-parameter” measurement weights.
- ▶ *Off-parameter* resampling is analogous to *off-policy* reinforcement learning.
- ▶ Log-measurement weights are discounted by a factor $\alpha = 1 - \epsilon$ to add stability. At $\theta = \phi$, log-measurement weights are all 0, so discounting has no effect.

DMOP- α is differentiated MOP- α

- ▶ The derivative of MOP- α at $\theta = \phi$ can be computed using AD applied to an algorithm similar to a basic particle filter.

Measurement off-parameter filter, MOP- α

First pass: Set $\theta = \phi$, yielding $X_{n,j}^{P,\phi}$, $X_{n,j}^{F,\phi}$, $g_{n,j}^\phi$.

Second pass: Filter at $\theta \neq \phi$ using the same RNG seed

For $n = 1, \dots, N$:

1. $w_{n,j}^{P,\theta} = (w_{n-1,j}^{F,\theta})^\alpha$
2. $X_{n,j}^{P,\theta} \sim \text{process}_n(\cdot | X_{n-1,j}^{F,\theta}; \theta)$.
3. $g_{n,j}^\theta = f_{Y_n|X_n}(y_n^* | X_{n,j}^{P,\theta}; \theta)$.
4. **Likelihood:** $L_n^{\theta,\alpha} = \sum_{j=1}^J g_{n,j}^\theta w_{n,j}^{P,\theta} / \sum_{j=1}^J w_{n,j}^{P,\theta}$.
5. Draw $k_{1:J}$ with $P(k_j = m) \propto g_{n,m}^\theta$.
6. $X_{n,j}^{F,\theta} = X_{n,k_j}^{P,\theta}$.
7. $w_{n,j}^{F,\theta} = w_{n,k_j}^{P,\theta} g_{n,k_j}^\theta / g_{n,k_j}^\phi$.

Differentiated filter, DMOP- α

First pass: Evaluate and build computation graph

Second pass: Differentiate via the chain rule

For $n = 1, \dots, N$:

1. $w_{n,j}^{P,\theta} = (w_{n-1,j}^{F,\theta})^\alpha$
2. $X_{n,j}^{P,\theta} \sim \text{process}_n(\cdot | X_{n-1,j}^{F,\theta}; \theta)$.
3. $g_{n,j}^\theta = f_{Y_n|X_n}(y_n^* | X_{n,j}^{P,\theta}; \theta)$.
4. **Likelihood:** $L_n^{\theta,\alpha} = \sum_{j=1}^J g_{n,j}^\theta w_{n,j}^{P,\theta} / \sum_{j=1}^J w_{n,j}^{P,\theta}$.
5. Draw $k_{1:J}$ with $P(k_j = m) \propto g_{n,m}^\theta$.
6. $X_{n,j}^{F,\theta} = X_{n,k_j}^{P,\theta}$.
7. $w_{n,j}^{F,\theta} = w_{n,k_j}^{P,\theta} g_{n,k_j}^\theta / \text{stop_gradient}(g_{n,k_j}^\theta)$.

Comparing MOP- α and DMOP- α

- ▶ The two algorithms are similar
- ▶ DMOP- α only has θ , not ϕ .
- ▶ The role of ϕ in MOP- α is taken by the `stop_gradient` operator in DMOP- α
 - ▶ `stop_gradient(x)` evaluates to x on the forward pass but is not differentiated on the backward pass.

Software for DMOP: pypomp

First, an introduction to the R-pomp family of packages

- ▶ The pomp R package has provided a solid, principled, extendable framework for combining time series data with nonlinear stochastic partially observed mechanistic dynamic models.
- ▶ Extensions: panel time series data (panelPomp), spatiotemporal data (spatPomp), phylodynamic data (phyloPomp).
- ▶ Limitations: R is not convenient for AD or GPU computing. To incorporate these, it may be better to start afresh.

pypomp is a Python port of pomp enabling AD and GPU.

- ▶ Current status: pypomp is in alpha development at <https://github.com/pypomp>. The brave are welcome to join.

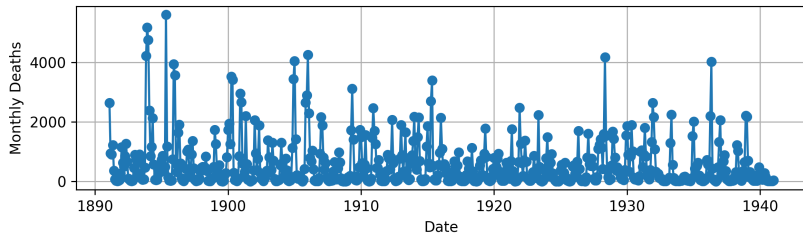
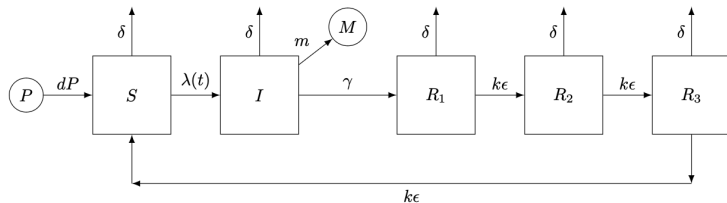
Why JAX?

- ▶ Actively developed by Google: “Google researchers have built and trained models like Gemini and Gemma on JAX, and it’s also used by researchers for a wide range of advanced applications.”
(<https://io.google/2025/explore/technical-session-1>)
- ▶ JAX code is Python, replacing numpy and scipy with `jax.numpy` and `jax.scipy`
- ▶ JAX composable transformations:
 - ▶ `vmap`: vectorization
 - ▶ `grad`: autodiff
 - ▶ `jit`: just-in-time compilation
- ▶ JAX uses XLA: accelerated linear algebra
 - ▶ distributes work across CPU/GPU/TPU cores

A previously difficult POMP likelihood maximization is easy with AD and quick/cheap with GPU

- ▶ We benchmark on an SIR³S model fitted to 50yr of monthly cholera mortality in historical Dacca, Bangladesh.
- ▶ Flexible modeling of seasonality and long-term trends lead to a highly parameterized model.
- ▶ Tractable, but difficult, using an early iterated filtering (IF1) algorithm in King et al. (2008). Possible with a large computing cluster and considerable dedication.
- ▶ Considerably easier, yet still hard, using the IF2 algorithm of Ionides et al. (2015).
- ▶ Here, we initialize with IF2 and then do stochastic gradient ascent with DMOP- α .

Model diagram (top) and data (bottom)



Results 1. Speed

- ▶ PF on 1 cpu core with 10^4 particles using R-pomp with the model compiled into C
 - ▶ 9.75 sec
- ▶ PF on a 10^4 core GPU with 10^4 particles using pypomp
 - ▶ 0.17 sec
- ▶ Here, a $10^4\%$ core GPU is worth 57 CPU cores.
- ▶ This is deliberately indirect: we are also comparing compiled C with jit-compiled Python.
- ▶ Additional vectorization helps for replicated GPU evaluations
 - ▶ For 360 replications of iterated filtering, a GPU is worth 450 CPU cores.

Results 2. The value of AD

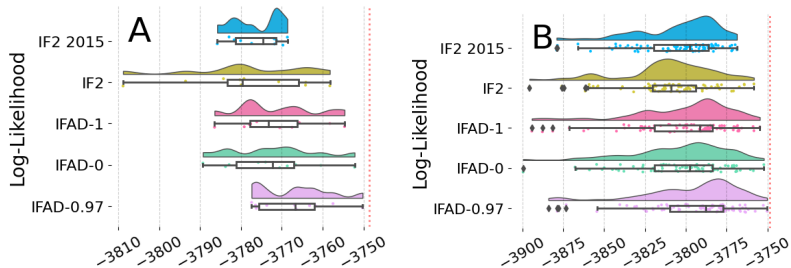


Figure 1: Performance of IFAD and IF2. **A:** the best out of every 10 runs. **B:** All searches. The dotted red line shows the true maximized log-likelihood.

pypomp.panel, pypomp.spat, pypomp.phylo

- ▶ There is no fundamental obstacle to including the R-pomp extensions in pypomp.
- ▶ DMOP in pypomp can improve maximization, and facilitate parallelization, but it does not deal with the “curse of dimensionality.”
- ▶ The basic particle filter can scale poorly with the dimension of the latent process. Intuition: importance sampling in a high-dimensional space is hard.
- ▶ In the basic phyloPomp particle filter, scaling can be poor with the number of demes: It is hard for the particles to maintain representative diversity for deme membership in a large tree.
- ▶ This is similar to the spatPomp problem of maintaining representative diversity at many spatial locations.

SpatPOMPs and the curse of dimensionality for SMC

Two approaches to scalability with (almost) full information.

1. **Guided particles.** Knowledge of the future can be built into proposal distributions to ensure that some particles are consistent with the data.
2. **Block particles.** Approximating conditional weak dependence by conditional independence.
 - ▶ Algorithmically, forecasts are made from the coupled spatiotemporal model, but data assimilation is carried out separately for each block.
 - ▶ Conditional on the coupled dynamics and the measurements, latent states for separate spatial blocks are approximated as conditionally independent.
 - ▶ This still permits dependence via the coupled dynamics.

Non-likelihood approaches could be implemented in phyloPomp: variational Bayes, ensemble Kalman filter, neural posterior estimation,

Summary

- ▶ Remarkably, likelihood-based inference is possible from noisy measurements of complex, high-dimensional, nonlinear dynamic systems.
- ▶ Examples follow.
- ▶ AD and massively parallel computing will push capabilities to a new level.

Thank you!

References I

- 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.
- King, A. A., Ionides, E. L., Pascual, M., and Bouma, M. J. (2008). Inapparent infections and cholera dynamics. *Nature*, 454:877–880.
- Tan, K., Hooker, G., and Ionides, E. L. (2024). Accelerated inference for partially observed Markov processes using automatic differentiation. *arXiv:2407.03085*.