

Anytime Parallel Tempering for Approximate Bayesian Computation

Alix Marie d’Avigneau · Sumeetpal Singh · Lawrence M. Murray

Received: date / Accepted: date

Abstract The development of scalable Markov chain Monte Carlo (MCMC) algorithms has recently gathered a lot of attention in the Bayesian inference literature. Indeed, they are the methods best equipped to deal with the ever increasing complexity of models. However, the cost of computing the likelihood of the data for these models tends to quickly become prohibitive, in particular when dealing with big data. A possible solution is the use of parallel tempering, which improves the performance of the algorithm by providing a more efficient exploration of the state space. For that, multiple Markov chain replicas are updated individually in parallel, often on multiple processors, and exchange moves are performed between the chains. Another problem is then encountered: all chains must be simultaneously ready for exchange moves to occur, and the real time taken for computations to complete may vary between chains. Also, factors such as competing jobs or unforeseen interruptions on even a single processor may delay the whole process. To solve this problem, an anytime Monte Carlo framework has been proposed, which imposes a real-time budget on the computations.

The first part of this paper presents and demonstrates the gain in efficiency provided by a new algorithm which combines parallel tempering with the anytime Monte Carlo framework, imposing a real-time deadline on the within-chain updates that occur in

parallel tempering to ensure all processors are simultaneously ready for the exchange steps. The second part adapts the new algorithm to the simulation-based class of algorithms known as Approximate Bayesian computation (ABC), which provides a likelihood-free approach to Bayesian inference. The resulting ABC algorithm is then applied to a problem in which the parameters of stochastic Lotka-Volterra predator-prey model must be estimated – a good example of a model in which the likelihood cannot be computed – to demonstrate the improvements in performance it provides.

Keywords Bayesian inference · Markov chain Monte Carlo (MCMC) · parallel tempering · anytime Monte Carlo (AMC) · Approximate Bayesian computation (ABC) · big data.

1 Introduction

Nowadays, in the age of big data, models devised to accurately represent the dynamics of data generating processes become more and more complex and high-dimensional. To deal with such problems, the widely applicable *Markov Chain Monte Carlo* (MCMC) algorithms are the best suited methods.

Consider a set of m observations $y = \{y_1, \dots, y_m\} \subset \mathbb{R}^m$ following a probability model with underlying parameters $\theta \in \Theta$ and associated *likelihood* $f(y|\theta) = f(y_1, \dots, y_m|\theta)$. The parameters θ describing the data y are unknown and considered random variables with *prior* density $p(d\theta)$. The aim of Bayesian inference is to use the information contained in the prior $p(d\theta)$ – representing prior belief – and the likelihood $f(y|\theta)$ – rep-

A. Marie d’Avigneau
University of Cambridge
E-mail: agem2@cam.ac.uk

S. Singh
University of Cambridge
E-mail: sss40@cam.ac.uk

representing evidence from the data – to obtain the *posterior* density $\pi(d\theta)$ of the parameters θ , following Equation 1, where the symbol \propto represents proportionality up to a constant. Summary statistics such as parameter estimates and credible intervals can subsequently be inferred from the posterior obtained.

$$\pi(d\theta) \propto p(d\theta)f(y|\theta) \quad (1)$$

In most cases however, the posterior π is intractable and must be approximated using computational tools such as MCMC algorithms. A commonly used and easily adaptable MCMC algorithm is the *Metropolis-Hastings* (M-H) algorithm, described in Robert and Casella [2004]. The Metropolis-Hastings algorithm starts at user-defined state θ_0 . At the n -th iteration, given current state θ_n , a new candidate $\theta' \sim q(d\theta'|\theta)$ is proposed, where q is the proposal density. Then, set $\theta_{n+1} = \theta'$ with probability

$$\begin{aligned} \alpha(\theta_n, \theta') &= \min \left\{ 1, \frac{\pi(\theta') q(\theta_n|\theta')}{\pi(\theta_n) q(\theta'|\theta_n)} \right\} \\ &= \min \left\{ 1, \frac{p(\theta') f(y|\theta') q(\theta_n|\theta')}{p(\theta_n) f(y|\theta_n) q(\theta'|\theta_n)} \right\} \end{aligned}$$

by Equation 1

otherwise retain $\theta_{n+1} = \theta_n$. Following these steps, a Markov chain is constructed whose stationary or invariant distribution is π . This enables one to sample from and hence provide an empirical approximation of the posterior.

The main drawback of MCMC algorithms is the computational cost due to the fact that they require evaluating a likelihood function $f(y|\theta)$ for the full dataset. From this, two problem may arise: the likelihood may be intractable, and it may be too computationally costly to evaluate the likelihood, either because it is too complex or because the datasets are too large. In recent years, many methods have been devised to tackle these issues and increase the efficiency of MCMC algorithms. These include divide-and-conquer methods, which reduce computing time by dividing the data into batches which are then updated in parallel before combining their result to obtain an approximation of the posterior. They also include subsampling algorithms, which speed up computations by reducing the amount of likelihood evaluations that occur at each iteration. Simulation-based methods such as Approximate Bayesian computation (ABC) have also been devised to avoid computing the likelihood altogether. Finally, the use of distributed computing has the potential to greatly speed up computations

by performing several tasks in parallel on multiple processors. A particular method that is compatible with distributed computing is parallel tempering, which allows for a more efficient exploration of the sample space by introducing exchange moves between multiple Markov chains which run in parallel, targeting different temperatures of the posterior.

It should nonetheless be noted that when dealing with distributed computing, real-time budgets arise. For example, the use of cloud computing is often financially costly and users have to deal with factors such as processor hardware, memory bandwidth, network traffic, I/O load, competing jobs on the same processors as well as potential unforeseen interruptions due to e.g. system failures, all of which affect the compute time of algorithms. Because of this, the anytime Monte Carlo framework was devised in Murray et al. [2016] to provide control over the total compute budget of Monte Carlo algorithms and ensure it is respected. This includes ensuring that all processors are simultaneously ready before they communicate instead of staying idle until the last of them has completed computations.

In this project, we combine parallel tempering with the anytime Monte Carlo framework to create Anytime Parallel Tempering Monte Carlo (APTMC) algorithm. While parallel tempering provides an increase in efficiency, the anytime framework essentially provides control over the budget of the parallel tempering algorithm and eliminates the potential bottleneck. An application of the algorithm to Approximate Bayesian computation subsequently provides a taster of its benefits in situations where the likelihood function is either unavailable or too computationally costly.

This paper is structured as follows. First, Section 2 offers a short review of the existing literature, including overviews of the parallel tempering algorithm and anytime Monte Carlo framework. Section 3 introduces the focus of this report: the Anytime Parallel Tempering Monte Carlo algorithm, illustrated in Section 5.1 with a toy example. Then, an application of the algorithm to ABC is provided in Section 4 and applied to a simple example as well as a MA(q) process parameter estimation problem before moving on to a the more advanced problem of estimating the parameters of a stochastic Lotka-Volterra model, in which the likelihood is unavailable. Finally, Section 6 concludes and presents a timeline of future work.

2 Literature review

In recent years, there have been many efforts to develop increasingly efficient or scalable MCMC algorithms, and in this chapter we will review a few examples, spending more time on the methods directly relevant to the main algorithm presented in this report.

2.1 Overview of scalable MCMC methods

An excellent review of existing methods that aim to scale up the Metropolis-Hastings algorithm for big data applications — so when the computation of the likelihood is too costly — is available in [Bardenet et al. \[2015\]](#). The authors divide the available approaches into two categories: divide-and-conquer and subsampling methods.

First of all, *divide-and-conquer* methods aim to divide the data into batches, then run the MCMC algorithm on each batch separately, usually in parallel, before combining the subposteriors obtained to form an approximation of the full posterior. Examples include works by [Neiswanger et al. \[2013\]](#), [Wang and Dunson \[2013\]](#), [Xu et al. \[2014\]](#) and [Minsker et al. \[2014\]](#). The main issues these methods must address are how to keep communication between batches minimal and how to efficiently combine the subposteriors, but they have an important advantage which is the possibility of running the algorithm on multiple processors, and thus significantly speeding up computations. This is the case of the consensus Monte Carlo algorithm presented by [Scott et al. \[2016\]](#).

On the other hand, *subsampling* methods aim to reduce the number of likelihood evaluations to speed up computations. This approach includes *pseudo-marginal* MCMC (PMCMC) methods which employ unbiased estimators of the unnormalised target distribution. This usually means that only a subsample of the data is used at each iteration, which speeds up computations. A few examples of subsampling MCMC algorithms are the Bootstrap Metropolis-Hastings algorithm by [Liang et al. \[2016\]](#) and the confidence sampler developed by [Bardenet et al. \[2014\]](#) and extended in [Bardenet et al. \[2015\]](#) and in [Kohn et al. \[2016\]](#). More examples of subsampling MCMC are available in [Quiroz et al. \[2016\]](#) and [Korattikara et al. \[2014\]](#). It is also possible to use delayed acceptance MCMC, such as the Firefly algorithm in [Maclaurin and Adams \[2014\]](#). The advantage of delayed acceptance MCMC is that it avoids computation of the likelihood if there is evidence that the proposal will be rejected, however in general it will

compute the likelihood on the full dataset otherwise, which is not ideal when the likelihood itself is too computationally costly. This is avoided in [Quiroz et al. \[2017\]](#), where delayed acceptance is combined with subsampling.

Other methods to increase the efficiency of MCMC algorithms include parallelising the Metropolis-Hastings algorithm for use on multiple processors, as described in [Calderhead \[2014\]](#) and parallel tempering, described in Section 2.2, which is particularly useful when the target distribution is multimodal. Additionally, a likelihood-free approach to MCMC in the form of Approximate Bayesian computation (ABC) is detailed in Section 4.1.

2.2 Parallel Tempering

The notion of *Parallel Tempering* (PT) was initially proposed by [Swendsen and Wang \[1986\]](#) and further developed under the name Metropolis-coupled Markov chain Monte Carlo (MC)³ by [Geyer \[1991\]](#). A parallel tempering algorithm allows for steps of various sizes to be made when exploring the parameter space and is particularly effective when the distribution we wish to sample from has multiple modes. Consider the (global) Markov chain $(X^{1:A})_{n=1}^{\infty} = (X^1, \dots, X^A)_{n=1}^{\infty}$ with initial state $(X_0^{1:A})$ and target distribution

$$\pi(\mathrm{d}x^{1:A}) \propto \prod_{\lambda=1}^A \pi_{\lambda}(\mathrm{d}x^{\lambda})$$

where the $\pi_{\lambda}(\cdot)$ are asymptotically independent marginals corresponding to the target distribution of each of A chains running in parallel at different temperatures indexed by λ . One of these chains is the *cold* chain and its target distribution $\pi_{\lambda} = \pi$ is the posterior of interest. In parallel tempering, [Geyer \[2011\]](#) identifies two types of update that can be performed at each iteration:

1. *Within-component update* or *local moves*: generally a standard Gibbs or Metropolis-Hastings update applied to each tempered chain X^{λ} in parallel. The local moves can also be performed sequentially.
2. *Between-component update*, or *exchange moves*: propose to swap the states $x \sim \pi_{\lambda}$ and $x' \sim \pi_{\lambda'}$ of one or more pairs of adjacent chains. For each pair, accept a swap with probability

$$\min \left\{ 1, \frac{\pi_{\lambda}(x')\pi_{\lambda'}(x)}{\pi_{\lambda}(x)\pi_{\lambda'}(x')} \right\} \quad (2)$$

otherwise, the chains in the pair retain their current states. With the cold chain providing more precision and the warmer chains more freedom of movement when exploring the parameter space, the combination of the two types of update allows all chains to mix much faster than any one of them would mix on its own. This provides a way to jump from mode to mode in far fewer steps than would be required under a standard Metropolis-Hastings algorithm.

An advantage of parallel tempering is that it is possible to perform the local moves in parallel on multiple processors, which speeds up computations. However, these must be synchronised before exchange moves can be performed. This means that all processors must be idle until the slowest of them has completed the local moves. This is also clearly the case if we choose to perform the local moves sequentially on a single processor. The next section introduces the anytime framework which reduces this idle time and eliminates the ensuing bottleneck.

2.3 Anytime Monte Carlo

Generally, Monte Carlo algorithms aim to simulate a pre-determined number of samples, taking a random amount of time to complete the necessary computations. The *Anytime Monte Carlo* (AMC) framework, developed in Murray et al. [2016], considers instead the situation in which we impose a real-time deadline on computations such that it is the number of samples that is random. The main appeal of the anytime framework is its application to distributed computing. Recalling parallel tempering, when one performs the local moves on multiple processors, they must be synchronised – meaning that all processors must wait for the slowest of them to complete – before local moves can be performed. By fixing a real-time budget on all processors, this waiting time is reduced and the overall algorithm is rendered more efficient.

Now, the real time taken to draw each sample may depend on the states of the Markov chain. For example, Chapter 4 deals with Approximate Bayesian computation (ABC), in which a ‘race’ takes place to determine the next sample in the Markov chain. Parameters with higher likelihood will be accepted sooner and so yield a lower computation time. More generally, the samples may not be independent of their number and a length bias with respect to computation time becomes apparent. When an empirical approximation or average over all post burn-in samples is required, this bias diminishes with time, and for a long

enough computation it may be rendered negligible. However, the bias in the final state does not diminish with time, and when this final state is important – which is the case in parallel tempering – the bias cannot be avoided by running the algorithm for longer, and other methods must be devised to correct it.

Let $(X)_{n=0}^{\infty}$ be a Markov chain with initial state X_0 , evolving on state space \mathbb{X} , with transition kernel $X_n | x_{n-1} \sim \kappa(dx|x_{n-1})$ and target distribution $\pi(dx)$. Define the *hold time* H_{n-1} as the random and positive real time required to complete the computations necessary to transition from state X_{n-1} to X_n via the kernel κ . Then let $H_{n-1} | x_{n-1} \sim \tau(dh_{n-1}|x_{n-1})$ where τ is the hold time distribution.

Given assumptions that the hold time $H > \epsilon > 0$ for minimal time ϵ , we have $\sup_{x \in \mathbb{X}} \mathbb{E}[H | x] < \infty$, and the hold time distribution τ is homogeneous in time, it is possible to construct a Markov jump process $(X, L)(t)$ with stationary distribution

$$\alpha(dx, dl) = \frac{\bar{F}_{\tau}(l|x)}{\mathbb{E}[H]} \pi(dx) dl$$

where $F_{\tau}(l|x)$ is the cdf of $\tau(dh_n|x_n)$ and $\bar{F}_{\tau}(l|x) = 1 - F_{\tau}(l|x)$, and of which the marginal is

$$\alpha(dx) = \frac{\mathbb{E}[H | x]}{\mathbb{E}[H]} \pi(dx) \quad (3)$$

The distribution α is referred to as the *anytime distribution*. When interrupted at real time t , the state of a Monte Carlo computation targeting π is distributed according to the anytime distribution α , which can essentially be seen as a length biased target distribution.

Different situations can be established to correct this bias. The main idea is to make it so expected hold time is independent of X , which leads to $\mathbb{E}[H | x] = \mathbb{E}[H]$ and hence $\alpha(dx) = \pi(dx)$, following Equation 3. While this is trivially the case for iid sampling, non-iid sampling requires that for $K \geq 0$, we simulate $K + 1$ Markov chains, each targeting π using the same transition kernel κ and hold time distribution τ . By simulating these $K + 1$ chains on the same processor in a serial schedule, we ensure that whenever the real-time deadline t is reached, states from all but one of the chains, say the $(K + 1)$ th chain, are independently distributed according to π . Since the $(K + 1)$ th chain is the currently working chain, i.e. the latest to go through the simulation process, its state

at the real-time deadline is distributed according to α . Simply discarding the state of this $(K + 1)$ th chain eliminates the length bias.

Using this multiple chain construction, it is thus possible to draw samples from π by interrupting the process at any time t . This and Section 2.2 set the basis for the focus of this paper: the Anytime Parallel Tempering Monte Carlo algorithm, described next.

3 Introduction to Anytime Parallel Tempering Monte Carlo

3.1 Overview

Consider the problem in which we wish to sample from target distribution $\pi(dx)$. In a parallel tempering framework, construct A Markov chains where each individual chain λ targets the tempered distribution

$$\pi_\lambda(dx) = \pi(dx)^{\frac{\lambda}{A}}$$

and is associated with kernel $\kappa_\lambda(dx_n | dx_{n-1})$ and hold time distribution $\tau_\lambda(dh_n | x_n)$. For each chain λ , a real-time Markov jump process $(X^\lambda, L)(t)$ can be constructed targeting the anytime distribution $\alpha_\lambda(dx, dl)$. In addition to constructing an anytime Monte Carlo algorithm, we aim to temporarily interrupt the computations on a real-time schedule of times t_1, t_2, t_3, \dots to perform exchange moves between adjacent pairs of chains before resuming.

3.2 Exchange moves

3.2.1 Exchanging on α

We propose to swap the states $(x, l) \sim \alpha_\lambda$ and $(x', l') \sim \alpha_{\lambda'}$ where the chains λ and λ' are adjacent. Similarly to Equation 2, this swap is accepted with probability

$$\frac{\alpha_\lambda(x', l') \alpha_{\lambda'}(x, l)}{\alpha_\lambda(x, l) \alpha_{\lambda'}(x', l')} \quad (4)$$

Assume that the hold time distributions are temperature-homogeneous, such that $\tau_\lambda(dh_n | x_n) = \tau(dh_n | x_n)$ for all $\lambda = 1, \dots, A$. As a direct consequence, we also have $\bar{F}_\lambda(l | x) = \bar{F}(l | x)$ for all $\lambda = 1, \dots, A$, i.e. the hold time cdfs do not change with temperature. Therefore, under this assumption, the expression in Equation 4 simplifies to the standard exchange probability for $x \sim \pi_\lambda$ and $x' \sim \pi_{\lambda'}$ given in Equation 2, and standard parallel tempering algorithms are actually implementing exchange moves on

the corresponding tempered anytime distributions.

However, when exchanging on α , the (current) final state of each Markov chain is used, meaning that there is still a length bias present, such that

$$\mathbb{E} \left\{ \frac{1}{n} \sum_{i=1}^n \zeta(X_i^A) \right\} \neq \int \zeta(x) \pi_A(dx)$$

where ζ is an arbitrary function, π_A is the desired (cold) distribution and X_A^i for $i = 1, \dots, n$ is the sampled trajectory, post burn-in, of the (cold) α_A chain. This sample bias diminishes as n increases, and would easily become negligible if the whole sampled trajectory, rather than the final state, was required here, or if exchange moves didn't occur often enough compared to local moves.

3.2.2 Exchanging on π

Here, we adapt the multi-chain construction devised to remove the bias present when sampling from A Markov chains, where each chain λ targets the distribution π_λ for $\lambda = 1, \dots, A$. Associated with each chain is MCMC kernel $\kappa_\lambda(dx_n^\lambda | dx_{n-1}^\lambda)$ and hold time distribution $\tau_\lambda(dh | x)$. The corresponding joint anytime distribution is

$$\begin{aligned} A(dx^{1:A}, dl, j) \\ = \frac{1}{A} \left(\alpha_j(dx^j, dl) \prod_{\lambda=1, \lambda \neq j}^A \pi_\lambda(dx^\lambda) \right) \frac{\mathbb{E}[H | j]}{\mathbb{E}[H]} \end{aligned}$$

for $\lambda = 1, \dots, A$. As updates are performed sequentially, j represents here the chain being updated while the other ones are not. Conditioning on x^j, j and l we obtain

$$A(dx^{1:A \setminus j} | x^j, l, j) = \prod_{\lambda=1, \lambda \neq j}^A \pi_\lambda(dx^\lambda) \quad (5)$$

Therefore, if exchange moves on the conditional $A(dx^{1:A \setminus j} | x^j, l, j)$ are performed by 'eliminating' the j -th chain to obtain the expression in Equation 5, they are being performed involving only chains distributed according to π and thus the bias is eliminated.

3.3 Implementation

3.3.1 One processor

On a single processor, the algorithm may proceed as in Algorithm 1, where in Step 3 the A chains are simulated one at a time in a serial schedule. Figure 1 provides an illustration of how the algorithm works using the example presented in Section 5.1.

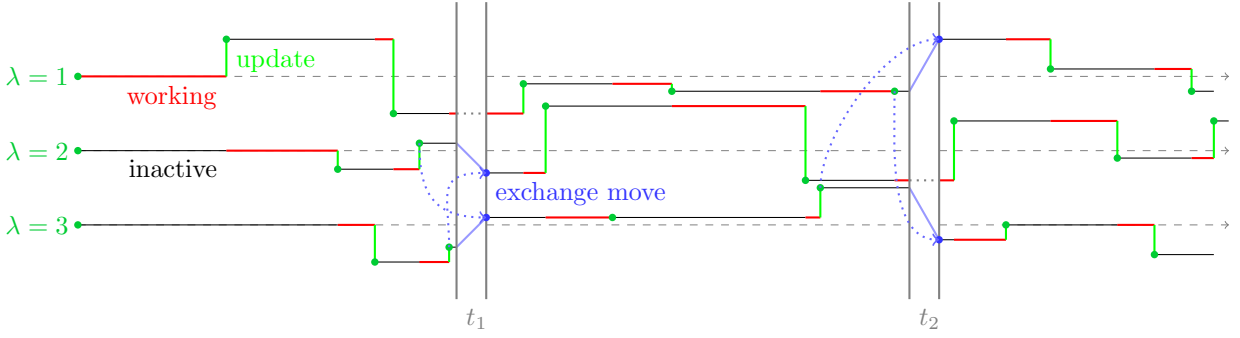


Fig. 1 Illustration of the progression of three chains in the Anytime Parallel Tempering Monte Carlo algorithm for the example in Section 5.1 on a single processor. The *green* (local move) and *blue* (exchange move) dots represent samples from the posterior being recorded as their respective local and exchange moves are completed. When exchange moves occur (at t_1 and t_2), the chain that is currently moving cannot participate in exchange moves without introducing a bias. Therefore it is ignored, and the exchange moves are performed on the remaining (inactive) chains.

Algorithm 1 Anytime Parallel Tempering Monte Carlo on one processor

- 1: Initialise real-time Markov jump process $(X^{1:A}, L, J)(0) = (x_0^{1:A}, 0, 1)$
- 2: **for** $i = 1, 2, \dots$ **do**
- 3: Simulate real-time Markov jump process $(X^{1:A}, L, J)(t)$ until real time t_i
- 4: Perform exchange steps on the conditional

$$A(dx^{1:A \setminus j} | x^j, l, j) = \prod_{\lambda=1, \lambda \neq j}^A \pi_\lambda(dx^\lambda)$$

5: **end for**

3.3.2 Multiple processors

When multiple processors are available, the chains may be run in parallel as described in Algorithm 2. Each worker uses K chains where either all chains have the same target distribution, e.g. for $W = A$ workers, worker $w = \lambda$ contains K chains targeting π_λ , or each chain has a different target distribution. For example, with $W = \frac{A}{2}$ workers, worker w could contain two chains, one with target π_w and one with target π_{2w} , or alternatively one one with target π_{2w-1} and one with target π_{2w} . Note that the multiple chain construction eliminates the intractable densities in the acceptance ratio for the exchange step when τ differs between processors.

4 Application to Approximate Bayesian Computation (ABC)

In this section we adapt the anytime parallel tempering Monte Carlo framework to Approximate Bayesian computation (ABC). The resulting algorithm will subsequently be applied to three case studies in Section 5.

Algorithm 2 Anytime Parallel Tempering Monte Carlo on multiple processors

- 1: On worker w , initialise the real-time Markov jump process $(X_w^{1:K}, L_w, J_w)(0) = (x_0^{1:K}, 0, 1)$
- 2: **for** $i = 1, 2, \dots$ **do**
- 3: On each worker, simulate the real-time Markov jump process $(X_w^{1:K}, L_w, J_w)(t)$ until real time t_i
- 4: Across all workers, perform exchange steps on the conditional

$$A(dx^{1:K \setminus j} | x^j, l, j) = \prod_{w=1}^W \prod_{k=1, k \neq j_w}^K \pi_w(dx_w^k)$$

where $\mathbf{dx}^{1:K \setminus j} = (dx_1^{1:K \setminus j_1}, \dots, dx_W^{1:K \setminus j_W})$, $\mathbf{x}^j = (x_1^{j_1}, \dots, x_W^{j_W})$, $\mathbf{l} = (l_1, \dots, l_W)$ and $\mathbf{j} = (j_1, \dots, j_W)$

5: **end for**

4.1 Overview of Approximate Bayesian Computation

The notion of ABC was developed by Tavaré et al. [1997] and Pritchard et al. [1999] and can be seen as a likelihood-free way to perform Bayesian inference, using instead simulations from the model or system of interest and comparing them to the observations available.

Let $y \in \mathbb{R}^d$ be some data with underlying unknown parameters $\theta \sim p(d\theta)$, where $p(\theta)$ denotes the prior for $\theta \in \Theta$. Suppose we are in the situation in which the likelihood $f(y|\theta)$ is either intractable or too computationally expensive, which means that MCMC cannot be performed as normal. Assuming that it is possible to sample from the density $f(\cdot|\theta)$ for all $\theta \in \Theta$, approximate the likelihood by introducing an artificial likelihood f^ε of the form

$$f^\varepsilon(y|\theta) = \text{Vol}(\varepsilon)^{-1} \int_{B_\varepsilon(y)} f(x|\theta) dx \quad (6)$$

where $B_\varepsilon(y)$ denotes a metric ball centred at y of radius $\varepsilon > 0$ and $\text{Vol}(\varepsilon)$ is its volume. The resulting approximate posterior is given by

$$p^\varepsilon(\theta | y) = \frac{p(\theta)f^\varepsilon(y | \theta)}{\int p(\vartheta)f^\varepsilon(y | \vartheta)d\vartheta}$$

The likelihood $f^\varepsilon(y | \theta)$ cannot be evaluated either, but a MCMC kernel can be constructed to obtain samples from the approximate posterior $\pi^\varepsilon(\theta, x)$ defined as

$$\pi^\varepsilon(\theta, x) = p^\varepsilon(\theta, x | y) \propto p(\theta)f(x, \theta)\mathbb{1}_\varepsilon(x)\text{Vol}(\varepsilon)^{-1}$$

where $\mathbb{1}_\varepsilon(x)$ is the indicator function for $x \in B_\varepsilon(y)$. This is referred to as *hitting* the ball $B_\varepsilon(y)$. In the MCMC kernel, one can propose $\theta' \sim q(d\theta' | \theta)$ for some proposal density q , simulate the dataset $x \sim f(dx | \theta')$ and accept θ' as a sample from the posterior if $x \in B_\varepsilon(y)$.

The *1-hit MCMC kernel*, proposed by Lee [2012] and described in Algorithm 3 introduces local moves in the form of a ‘race’: given current and proposed parameters θ and θ' , respectively simulate corresponding datasets x and x' sequentially. The state associated with the first dataset to hit the ball $B_\varepsilon(y)$ ‘wins’ and is accepted as the next sample in the Markov chain.

Algorithm 3 ABC: 1-hit MCMC kernel

```

Given current state  $(\theta_n, x_n)$ 
1: for  $i := 1, 2, \dots$  do
2:   Propose  $\theta' \sim q(d\theta' | \theta_n)$   $\triangleright$  propose a local move
3:   Compute preliminary acceptance probability  $\triangleright$  prior
   check
   
$$\alpha(\theta_n, \theta') = \min \left\{ 1, \frac{p(\theta')q(\theta_n | \theta')}{p(\theta_n)q(\theta' | \theta_n)} \right\}$$

4:   Sample  $u \sim \text{Uniform}(0, 1)$ 
5:   if  $u < \alpha(\theta_n, \theta')$  then
6:     RACE := TRUE
7:   else
8:     RACE := FALSE
9:     retain  $(\theta_{n+1}, x_{n+1}) = (\theta_n, x_n)$   $\triangleright$  automatically
   reject  $\theta'$  as it is unlikely to win
10:  end if
11:  while RACE do
12:    Simulate  $x \sim f(dx | \theta_n)$  and  $x' \sim f(dx' | \theta')$ 
13:    if  $x \in B_\varepsilon(y)$  or  $x' \in B_\varepsilon(y)$  then  $\triangleright$  stop the race
   once either  $x$  or  $x'$  hits the ball
14:      RACE := FALSE
15:    end if
16:  end while
17:  if  $x'$  falls within  $\varepsilon$  of  $y$  then  $\triangleright$  accept or reject move
18:    set  $(\theta_{n+1}, x_{n+1}) = (\theta', x')$ 
19:  else
20:    retain  $(\theta_{n+1}, x_{n+1}) = (\theta_n, x_n)$ 
21:  end if
22:   $n := n + 1$ 
23: end for
```

4.2 ABC Anytime Parallel Tempering Monte Carlo (ABC-APTMC)

We now enter an anytime Monte Carlo setting and introduce two types of exchange moves to the 1-hit MCMC kernel.

4.2.1 Fast exchange moves

Let (θ, x) and (θ', x') be the states of two chains targeting π^ε and $\pi^{\varepsilon'}$, respectively, where $\varepsilon' > \varepsilon$. Here, this is equivalent to saying θ' is the state of the ‘warmer’ chain. We already know that x' falls within ε' of the observations y , i.e. $x' \in B_{\varepsilon'}(y)$. Similarly, we also know that $x \in B_\varepsilon(y)$, and clearly that $x \in B_{\varepsilon'}(y)$. If x' also falls within ε of y , then swap the states, otherwise do not swap. The odds ratio is

$$\begin{aligned} & \frac{\pi^{\varepsilon'}(\theta, x)\pi^\varepsilon(\theta', x')}{\pi^\varepsilon(\theta, x)\pi^{\varepsilon'}(\theta', x')} \\ &= \frac{p(\theta)f(x | \theta)\text{Vol}(\varepsilon')p(\theta')f(x' | \theta')\mathbb{1}_\varepsilon(x')\text{Vol}(\varepsilon)}{p(\theta)f(x | \theta)\text{Vol}(\varepsilon)p(\theta')f(x' | \theta')\text{Vol}(\varepsilon')} \\ &= \mathbb{1}_\varepsilon(x') \end{aligned}$$

so the probability of the swap being accepted is the probability of x' also hitting the ball of radius ε centred at y . This type of exchange move is summarised in Algorithm 4. It is immediate and cheap. However, in a difficult model, where the local moves on, say, the $(K + 1)$ th chain can sometimes get stuck running for an extended time period (see Example 5.4), there is a chance the same two chains will be selected for exchange moves more than once before they have had a chance to move forward locally. Since it is entirely possible neither of the chain states will have changed after the first exchange move, all subsequent swaps of these two chains risk being identical until the $(K + 1)$ th chain finally completes its local move. Either all swaps will be rejected, leading to each new sample on the chains being the same as the last, or all swaps will be accepted, meaning the two chains will be trading the same two values back and forth. This causes an increase in autocorrelation in the chains and reduces the algorithm’s efficiency. To circumvent this issue, one can completely discard the $(K + 1)$ th chain after interruption and propose a new candidate every time it resumes local moves. This does not work if not enough time is given for the working chain to complete a single move, but anytime algorithms are generally set so that all chains have a chance to complete a local move before attempting any swap. Alternatively, when dealing with a model that doesn’t necessarily get stuck in a race for a long time, but does frequently get stuck for short periods, another

kind of exchange moves can refresh x' at each new exchange move, thus opening the possibility to the outcome of the swap being different.

4.2.2 Slow exchange moves

First of all, let us establish the detailed balance condition for the 1-hit algorithm in Lee [2012] for an exchange move. Let $0 \leq f(\theta) \leq 1$ and $0 \leq f'(\theta) \leq 1$ be the marginal likelihoods over x , i.e.

$$f(\theta) = \int \mathbb{1}_\varepsilon(x) f(dx | \theta)$$

Here, f' corresponds to the probability of hitting a larger ball, $B_{\varepsilon'}(y)$ where $\varepsilon' > \varepsilon$, when sampling from $f(dx | \theta)$. Recall $p(\theta)$ is the prior on θ and define the following probability densities

$$\pi(\theta) \propto p(\theta)f(\theta) \quad \pi'(\theta') \propto p(\theta')f'(\theta')$$

We have $\theta \sim \pi$ and $\theta' \sim \pi'$ and we set the probability of accepting the swap $(\theta, \theta') \rightarrow (\theta', \theta)$ to be $\frac{f(\theta')}{f'(\theta')}$. This is equivalent to drawing a Geometric random variable with success probability of hitting the larger ball $f'(\theta')$, and accepting if success also means hitting the smaller ball, corresponding to f . Then we have

$$\begin{aligned} & \mathbb{E}_{\theta \sim \pi, \theta' \sim \pi'} [\mathbb{P}(\text{Accept}) \cdot \varphi(\theta', \theta)] \\ &= \int \pi(\theta) \pi'(\theta') \frac{f(\theta')}{f'(\theta')} \varphi(\theta', \theta) d\theta d\theta' \\ &= Z^{-1} \int p(\theta) f(\theta) p(\theta') f'(\theta') \frac{f(\theta')}{f'(\theta')} \varphi(\theta', \theta) d\theta d\theta' \\ &= Z^{-1} \int p(\theta) f'(\theta) p(\theta') f(\theta') \frac{f(\theta)}{f'(\theta)} \varphi(\theta', \theta) d\theta d\theta' \\ &= \int \pi'(\theta) \pi(\theta') \frac{f(\theta)}{f'(\theta)} \varphi(\theta', \theta) d\theta d\theta' \\ &= \mathbb{E}_{\theta \sim \pi, \theta' \sim \pi'} [\mathbb{P}(\text{Accept}) \cdot \varphi(\theta, \theta')] \end{aligned}$$

where Z is a normalising constant. The race completes itself in an expected number of steps proportional to $f'(\theta')^{-1}$ and the acceptance probability approaches 1 as $f' \rightarrow f$. A more detailed proof is available in Appendix A.1.

4.2.3 Implementation

In practice, the acceptance probability $\frac{f(\theta')}{f'(\theta')}$ cannot be computed analytically, so a new ABC-MCMC kernel is implemented in Algorithm 5. Let $\theta \sim \pi$ and $\theta' \sim \pi'$ where θ' is the state of the ‘warmer’ chain. Then, generate $x' \sim f(x' | \theta') \mathbb{1}_{\varepsilon'}(x')$ via rejection sampling, and swap the states θ and θ' if x' also hits the smaller ball $B_\varepsilon(y)$. Given f and f' correspond to the probabilities

of hitting $B_\varepsilon(y)$ and $B_{\varepsilon'}(y)$, respectively, this is equivalent to accepting the swap with probability $\frac{f(\theta')}{f'(\theta')}$. The full implementation of the ABC Anytime Parallel Tempering Monte Carlo (ABC-APTMC) algorithm on a single processor is described in Algorithm 6. The multi-processor algorithm can similarly be modified to reflect these new exchange moves.

4.2.4 Remarks

Note that the slow exchange moves differ from the ones presented previously as they take a random amount of time to complete, as opposed to a fixed amount of time. While local moves are momentarily interrupted following a real time schedule t_1, t_2, \dots , no such interruptions are forced upon slow exchange moves, except for the hard deadline T which stops all computations. It is also important to note that if the race in the local move of, say, the $(K + 1)$ th chain takes too long to complete, none of the other K chains can progress locally. Therefore, any new samples on those ‘idle’ chains will be a result of exchange moves, and they will be essentially swapping around the same K samples until the race on the working chain finishes. This is less of an issue on multiple processors, as the chains on other workers are still able to move forward locally. All this means that, depending on where the chains were initialised, there is a risk at time T that very few effective updates will have occurred, as the race and simulation steps in Algorithms 3 and 5, respectively will have taken up most of the computation time. To suppress such a risk, solutions are the following:

- To ensure enough local updates and exchange moves occur, initialise the real-time Markov process near the centre of the distribution, e.g. simulate a few samples from the posterior using ABC rejection sampling and initialise the chains from those. This is the solution used in throughout this paper, and is enough to ensure the simulation step in Algorithm 5 never takes too long.
- Alternatively, initialise the algorithm with balls of relatively large radii ε'_0 and ε_0 to increase the chances of hitting the larger ball in Step 4 of Algorithm 5, and then reduce ε'_t and ε_t with time t to progressively increase accuracy. One can either choose to reduce the radii until a pre-defined time t_{\min} and keep them fixed thereafter, or alternatively leave them to shrink until the deadline T .

Algorithm 4 ABC: Fast exchange move between two chains

Given $\omega_n = ((\theta, x), (\theta', x'))$ where $\theta \sim \pi$, $x \sim f(dx|\theta)$ and $\theta' \sim \pi'$, $x' \sim f(dx'|\theta')$.
 \triangleright both (θ, x) and (θ', x') are outputs from Algorithm 3 for different $\varepsilon' > \varepsilon$

- 1: **if** $x' \in B_\varepsilon(y)$ **then** \triangleright accept or reject swap depending on whether x' also hits the ball of radius ε
- 2: set $\omega_{n+1} = ((\theta', x'), (\theta, x))$
- 3: **else**
- 4: retain $\omega_{n+1} = \omega_n$
- 5: **end if**
- 6: $n := n + 1$

Algorithm 5 ABC: Slow exchange move between two chains

Given $\omega_n = ((\theta, x), (\theta', x'))$ where $\theta \sim \pi$, $x \sim f(dx|\theta)$ and $\theta' \sim \pi'$, $x' \sim f(dx'|\theta')$.

- 1: $\text{SAMPLE} := \text{TRUE}$
- 2: **while** SAMPLE **do**
- 3: Simulate $z' \sim f'(dz'|\theta')$
- 4: **if** $z' \in B_{\varepsilon'}(y)$ **then** \triangleright stop sampling once z' hits the ball of radius ε'
- 5: $\text{SAMPLE} := \text{FALSE}$
- 6: **end if**
- 7: **end while**
- 8: Perform swap following Algorithm 4 using new (θ', z') to represent the warmer chain.

Algorithm 6 ABC: Anytime Parallel Tempering Monte Carlo Algorithm

- 1: Initialise the real-time Markov jump process $(\theta^{1:A}, L, J) = (\theta^{1:A}, 0, 1)$.
- 2: Set $n := 0$
- 3: **for** $i := 1, 2, \dots$ **do**
 SIMULATE THE REAL-TIME MARKOV JUMP PROCESS $(\theta, L, J)(t)$ UNTIL REAL TIME t_i
- 4: Perform local moves on (θ_n^j, x_n^j) according to Algorithm 3.
- 5: $j := j + 1$
- 6: **if** $j > \Lambda$ **then**
- 7: $j := 1$
- 8: **end if**
- 9: PERFORM EXCHANGE STEPS ON THE CONDITIONAL:

$$A(d\theta^{1:A} | \theta^j, l, j) = \prod_{\lambda=1, \lambda \neq j}^{\Lambda} \pi_\lambda(d\theta^\lambda)$$

- 9: Perform fast or slow exchange moves on $\omega_n = ((\theta_n^\lambda, x_n^\lambda), (\theta_n^{\lambda'}, x_n^{\lambda'}))$ according to Algorithm 4 or 5, respectively.
- 10: **end for**

5 Experiments

In this section, we first illustrate the workings of the algorithms presented in Section 3.3 on a simple model, in which real-time behaviour is simulated using virtual time and an artificial hold distribution. The model is also employed to demonstrate the gain in efficiency pro-

vided by the inclusion of exchange moves. Then, the ABC version of the algorithms, as presented in Section 4 is applied to three case studies. The first case is a simple model and serves to verify the workings of the ABC algorithm, including bias correction. The second case considers an ABC approach to estimating the parameters of a moving average $\text{MA}(q)$ process. It illustrates the improvements in performance introduced by the addition of ABC exchange moves as opposed to the standard ABC-MCMC algorithm. The third case considers the problem of estimating the parameters of a stochastic Lotka-Volterra predator-prey model – in which the likelihood is unavailable – and serves to evaluate the performance of the anytime parallel tempering version of the ABC-MCMC algorithm as opposed to the standard versions (with and without exchange moves) on both a single and multiple processors. The exchange moves used in the ABC case studies were the fast version, so that multiple pairs could be swapped at each iteration. The experiments were also tested with slow exchange moves and returned similar outcomes. All experiments in this paper were run on MATLAB and the code is available at <https://github.com/alixma/ABCPTMC.git>.

5.1 Toy example: Gamma mixture model

In this example we attempt to sample from an equal mixture of two Gamma distributions using the Anytime Parallel Tempering Monte Carlo (APTMC) algorithm. Define the target $\pi(dx)$ and an ‘artificial’ hold time $\tau(dh|x)$ distributions as follows:

$$X \sim \phi \text{Gamma}(k_1, \theta_1) + (1 - \phi) \text{Gamma}(k_2, \theta_2)$$

$$H|x \sim \psi \text{Gamma}\left(\frac{x^p}{\theta_1}, \theta_1\right) + (1 - \psi) \text{Gamma}\left(\frac{x^p}{\theta_2}, \theta_2\right)$$

with mixture coefficients $\phi = \frac{1}{2}$ and ψ , where $\text{Gamma}(\cdot, \cdot)$ denotes the pdf of a Gamma distribution, with shape and scale parameters (k_1, θ_1) and (k_2, θ_2) for each components, respectively, and with polynomial degree p , assuming it remains constant for both components of the mixture.

In this example, virtual time is employed instead of real time. Usually, almost nothing is known about the hold time distribution τ , and in particular not its explicit form. However, for this toy example we assume an explicit form for τ is known and simulate virtual hold times. These artificial hold times are introduced such that what in a real-time example would be the effects of polynomial computational complexity can be studied,

including constant ($p = 0$), linear ($p = 1$), quadratic ($p = 2$) and cubic ($p = 3$) complexity. Another advantage is that the anytime distribution $\alpha_A(dx)$ of the cold chain can be computed analytically and is the following mixture of two Gamma distributions

$$\alpha_A(dx) = \varphi(p) \text{Gamma}(k_1 + p, \theta_1) + (1 - \varphi(p)) \text{Gamma}(k_2 + p, \theta_2) \quad (7)$$

where

$$\varphi(p) = \frac{1}{1 + \frac{\Gamma(k_1)\Gamma(p+k_2)\theta_2^p}{\Gamma(k_2)\Gamma(p+k_1)\theta_1^p}}$$

We refer the reader to Appendix A.2 for the proof of Equation 7. In the anytime distribution, one of the components of the Gamma distribution will have an associated mixture coefficient $\varphi(p)$ or $1 - \varphi(p)$ which increases with p while the coefficient of the other component decreases proportionally. Note that for constant ($p = 0$) computational complexity, the anytime distribution is equal to the target distribution π .

5.1.1 Implementation

On a single processor, the Anytime Parallel Tempering Monte Carlo algorithm is implemented as described in Algorithm 7. We simulate $\Lambda = 8$ Markov chains, each targeting the distribution $\pi_\lambda(dx) = \pi(dx)^{\frac{\lambda}{\Lambda}}$. In Part **A** of Algorithm 7, to construct a Markov chain $(X^\lambda)_{n=0}^\infty$ with target distribution

$$\pi_\lambda(x) \propto \left[\frac{1}{2} \text{Gamma}(k_1, \theta_1) + \frac{1}{2} \text{Gamma}(k_2, \theta_2) \right]^{\frac{\lambda}{\Lambda}}$$

for $\lambda = 1, \dots, \Lambda$, simply use a *Random Walk Metropolis* update, i.e. symmetric proposal distribution $\mathcal{N}(x_n^\lambda, \sigma^2)$ with mean x_n^λ standard deviation $\sigma = 0.5$. Set $(k_1, k_2) = (3, 20)$, $(\theta_1, \theta_2) = (0.15, 0.25)$ and use $p \in \{0, 1, 2, 3\}$. The single processor algorithm is run for $T = 10^8$ units of virtual time with exchange moves alternating between occurring on all even $(1, 2), (3, 4), (5, 6)$ and all odd $(2, 3), (4, 5), (6, 7)$ pairs of inactive chains every $\delta_T = 5$ time steps. When the algorithm is running, a sample is recorded every time a local or exchange move occurs.

On multiple processors, the Anytime Parallel Tempering Monte Carlo algorithm for this example is implemented similarly. A number of $W = \Lambda = 8$ workers is used, where each worker $w = \lambda$ contains $K = 2$ chains, all targeting the same π_λ for $\lambda = 1, \dots, \Lambda$. The multiple processor algorithm is run for $T = 10^7$ units of virtual time, with exchange moves alternating between occurring on all even $(1, 2), (3, 4), (5, 6), (7, 8)$

Algorithm 7 APTMC algorithm for Gamma mixture model on a single processor

```

1: Initialise the Markov jump process  $(X^{1:\Lambda}, L, J)(0) := (x_0^{1:\Lambda}, 0, 1)$ .
2: for  $i = 1$  to  $T$  by  $\delta_T$  do
  A: SIMULATE THE MARKOV JUMP PROCESS  $(X^{1:\Lambda}, L, J)(t)$  UNTIL TIME  $i$ 
3:   while  $t < i$  do
4:     if  $l < 1$  then  $\triangleright$  make sure the  $j$ -th chain isn't already being held
5:       Draw  $h \sim \text{Gamma}\left(\frac{1}{\theta_1} \left(x_{n^j}^j\right)^p, \theta_1\right)$ 
6:        $l := l + h$ 
7:     end if
8:     Hold the  $j$ -th chain for  $l$  units of time  $\triangleright$  simulates time a chain would take to complete local moves in a real-time example
9:     Draw proposal  $x^* \sim \mathcal{N}\left(x_{n^j}^j, \sigma^2\right)$ .
10:    Compute the acceptance probability

```

$$a = \min \left\{ 1, \frac{\pi_j(x^*)}{\pi_j(x_{n^j}^j)} \right\}$$

```

11:    Set  $x_{n^j+1}^j = x^*$  with probability  $a$ , otherwise retain  $x_{n^j+1}^j = x_{n^j}^j$ .
12:     $n^j := n^j + 1$   $\triangleright$  update state index for  $j$ -th chain
13:     $j := j + 1$   $\triangleright$  move on to the next chain
14:    if  $j > \Lambda$  then reset  $j := 1$ 
15:  end while
  B: PERFORM EXCHANGE MOVES ON THE CONDITIONAL  $A(dx^{1:K \setminus j} | x^j, l, j)$ 
16:  Draw  $m$  pairs of adjacent chains uniformly based on the set  $\{1 : \Lambda\} \setminus j$ .
   $\triangleright$  e.g. discard the  $j$ -th chain and alternate between proposing to swap  $m$  even  $(1, 2), (3, 4), \dots$  and  $m$  odd  $(2, 3), (4, 5), \dots$  pairs at each iteration
17:  For each pair  $(\lambda, \lambda')$  of chains compute the acceptance probability

```

$$r(\lambda, \lambda') = \min \left\{ 1, \frac{\pi_\lambda(x_{n^{\lambda'}}^{\lambda'}) \pi_{\lambda'}(x_{n^\lambda}^\lambda)}{\pi_\lambda(x_{n^\lambda}^\lambda) \pi_{\lambda'}(x_{n^{\lambda'}}^{\lambda'})} \right\}$$

```

18:  Swap the states of the chains with probability  $r(\lambda, \lambda')$ , otherwise leave them unchanged.
19:  For each pair  $(\lambda, \lambda')$  of chains update state indices  $n^\lambda := n^\lambda + 1$  and  $n^{\lambda'} := n^{\lambda'} + 1$ 
20: end for

```

and all odd $(2, 3), (4, 5), (6, 7)$ pairs of workers every $\delta_T = 5$ steps. On each worker, the chain which was not working when calculations were interrupted is the one included in the exchange moves, as per Figure ??.

5.1.2 Verification of bias correction

To check that the single and multiple processor algorithms are successfully correcting for bias, they are also run *uncorrected*, i.e. with Part **B** adjusted to perform exchange steps on the conditional $A(dx^{1:K}, dl^{1:K})$. For

instance, Step 16 in Algorithm 7 becomes: ‘Draw m pairs of adjacent chains uniformly based on the set $\{1 : \Lambda\}$ ’ and a similar change is made for the multiple processor algorithm. This means that exchange moves are performed on α instead of π , thus causing the algorithm to yield biased results.

Since the bias is introduced by the exchange moves (when they are performed on α), we attempt to create a ‘worst case scenario’, i.e. maximise the amount of bias present when the single processor algorithm is uncorrected. The algorithm is further adjusted such that local moves are not performed on the cold chain in Part A and it is instead solely made up of samples resulting from exchange moves with the warmer chains. Additionally, the fact that exchange moves occur every $\delta_T = 5$ time steps means that a high proportion of the samples in a warmer chain come from exchange moves. The multi-processor algorithm is not run in a ‘worst case scenario’ as the virtual time required to properly verify that the corrected algorithm converges to the true posterior when $p = 3$ is too high. Local moves on the cold chain of the multi-processor algorithm are therefore allowed, and the bias caused by failing to omit the currently working chain when performing exchange moves across workers should still be apparent, if less strongly.

5.1.2.1 Results

Figure 2 shows kernel density estimates of the post burn-in cold chains resulting from runs of the single and multi-processor algorithms, uncorrected and corrected for bias. As expected, a constant ($p = 0$) computational complexity does not return any bias. While coming close but not quite completely reaching the corresponding anytime distributions – which would be the most extreme case of bias – the cold chains for the single-processor algorithm with computational complexity $p \in \{1, 2, 3\}$ have clearly converged to a shifted distribution which puts more weight the second Gamma mixture component (instead of an equal weight). Additionally, the bias becomes stronger as computational complexity p increases. A similar observation can be made for the cold chains from the multi-processor experiment – which display a milder bias due to local moves occurring on the cold chain. On the other hand, the dashed densities indicate that when the algorithms are corrected, i.e. when the currently working chain is not included in exchange moves, it successfully eliminates the bias for all $p \in \{1, 2, 3\}$ to return the correct posterior π – despite even this being the ‘worst case scenario’ in the case of the single

processor algorithm. Note that in the single processor experiment, the uncorrected density estimates never quite reach their corresponding anytime distributions because they are not solely made up of biased samples. Indeed, if chain $k = 2$ was not working before the uncorrected exchange move with the cold chain ($k = 1$) was accepted, then the new sample on the cold chain is unbiased as it’s receiving a sample from π_2 . Conversely, if chain 2 was working and the swap is performed, then the new sample on the cold chain will be from α_2 and therefore biased. By eliminating the local moves on the cold chain, we have significantly augmented the proportion of biased samples, but they still don’t make up 100% of the chain. Additionally, the stronger bias for higher computational complexities is due to more frequent exchange moves occurring on each chain, as their local moves are ‘slower’.

5.1.3 Performance evaluation

Next we verify that introducing the parallel tempering element to the anytime Monte Carlo algorithm improves performance. A standard MCMC algorithm is run for $T = 10^6$ units of virtual time and computational complexity $p \in \{0, 1, 2, 3\}$, applying the random walk Metropolis update described in Section 5.1.1. Additionally, both the single and multiple processor Anytime Parallel Tempering Monte Carlo algorithms are run again on $\Lambda = W = 8$ chains/workers, with $K = 2$ chains per worker for the multi-processor algorithm, for the same amount of virtual time and with exchange moves occurring every $\delta_T = 5$ time steps. This time, local moves are performed on the cold chain of the single processor APTMC algorithm.

To compare results, kernel density estimates of the posterior are obtained from the post burn-in cold chains for each algorithm using the `kde` function in [MATLAB \[2017\]](#), developed by [Botev et al. \[2010\]](#). It is also important to note that even though all algorithms run for the same (virtual) duration, the standard MCMC algorithm is performing local moves on a single chain uninterrupted until the deadline while the single-processor APTMC algorithm has to update $\Lambda = 8$ chains in sequence and each worker w of the multi-processor APTMC algorithm has to update $K = 2$ chains in sequence before exchange moves occur. Therefore, the algorithms are not expected to return samples of the similar sizes. For a fair performance comparison, the sample autocorrelation function (acf) is estimated first of all. When available, the acf is averaged over multiple chains to reduce variance in its estimates. Other tools employed are

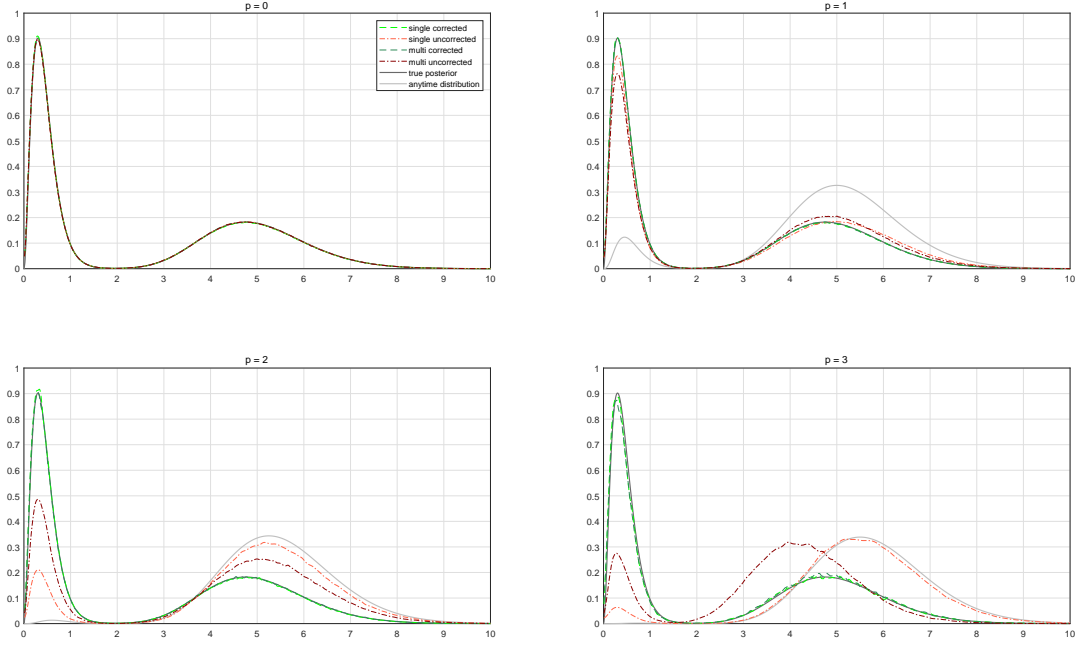


Fig. 2 Density estimates of the cold chain for bias corrected and uncorrected runs of the single- and multi-processor APTMC algorithm on various computational complexities $p \in \{0, 1, 2, 3\}$. In the single-processor case, the cold chains are made up entirely of updates resulting from exchange moves. The *dark gray* line represents the true posterior density π and the *light gray* line the anytime distribution α . The case $p = 0$ represents an instance in which, in a real-time situation, the local moves do not take a random time to complete, and therefore all densities are identical.

- *Integrated Autocorrelation Time (IAT)*, the computational inefficiency of a MCMC sampler. Defined as

$$IAT_s = 1 + 2 \sum_{\ell=1}^{\infty} \rho_s(\ell)$$

where $\rho_s(\ell)$ is the autocorrelation at the ℓ -th lag of chain s . It measures the average number of iterations required for an independent sample to be drawn, or in other words the number of correlated samples with same variance-reducing power as one independent sample. Hence, a more efficient algorithm will have lower autocorrelation values and should yield a lower IAT value. Here, the IAT is estimated using a method initially suggested in Sokal [1997] and Goodman and Weare [2010], and implemented in the Python package `emcee` by Foreman-Mackey, Daniel and Hogg, David W and Lang, Dustin and Goodman, Jonathan [2013] (Section 3). Let

$$\hat{IAT}_s = 1 + 2 \sum_{\ell=1}^M \hat{\rho}_s(\ell)$$

where M is a suitably chosen cutoff, such that noise at the higher lags is reduced. Here, the smallest M

is chosen such that $M \geq C \hat{\rho}_s(M)$ where $C \approx 6$. More information on the choice of C is available in Sokal [1997].

- *Effective Sample Size (ESS)*, the amount of information obtained from an MCMC sample. It is closely linked to the IAT by definition:

$$ESS_s = \frac{N_s}{\hat{IAT}_s}$$

where N_s is the size of the current sample s . The ESS measures the number of independent samples obtained from MCMC output.

The resulting ESS and IAT for different algorithms and computational complexities are computed and shown in Table 1. For the multi-processor APTMC algorithm, the IAT is averaged and the ESS is summed over the two output cold chains. We further note that while the proposal for local moves depends on the value of the previous state, the state of a chain if an exchange move is accepted does not, meaning that the autocorrelation in a chain containing a significant proportion of accepted samples from exchange moves will be lower. Finally, for low p , significantly more local moves occur as hold times are short while for a higher p the hold

times will be longer and hence fewer local moves are able to occur before each deadline. We therefore expect a higher proportion of samples from exchange moves for higher values of p and a more important increase in efficiency.

5.1.3.1 Performance results

In Figure 3 we observe, unsurprisingly, that the quality of the posterior estimates decreases as p increases. As a matter of fact, 10^6 units of virtual time tend to not be enough for the some of the cold chains to completely converge. The single processor APTMC algorithm overestimates the first mode and underestimates the second mode of the true posterior for $p = 2$, while none of the algorithms appear to have fully converged when $p = 3$. In general, the multi-processor APTMC returns results closest to the true cold posterior for $p = \{0, 1, 2\}$.

As for efficiency, Table 1 displays a much lower *IAT* and much higher *ESS* for both APTMC algorithms, indicating that they are much more efficient than the AMC algorithm. This is further supported by the sample autocorrelation decaying much more quickly for APTMC algorithms than for the MCMC algorithm for all p in Figure 4. In Table 1 the multi-processor APTMC yields *IAT* values that are lower than those returned by the single processor APTMC algorithm, and similarly yields effective sample sizes that are higher for $p < 3$. It is however important to note that comparing these values for $p = 3$ can be misleading, since in Figure 3, none of the algorithm outputs have converged to the desired distributions for this value of p . The efficiency of exchange moves on a single processor compared to standard MCMC, along with the reasonably accurate results displayed in Figure 3 indicate that they are a good choice for algorithms with low computational complexity if access multiple processors is not easily available. However if one has access to parallel computing tools, the multi-processor APTMC is the most efficient and accurate choice for this problem.

Next, we consider an application of the anytime parallel tempering Monte Carlo framework to a class of algorithms that are well-adapted to situations in which the likelihood is either intractable or computationally prohibitive. They are called Approximate Bayesian computation and feature a non-artificial hold time at each MCMC iteration, making them ideal candidates for adaptation to the anytime parallel tempering framework.

p	Multi-processor		Single-processor		Standard	
	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>
0	50.931	12757	70.584	1382.3	1853.4	269.71
1	46.733	5723.3	100.18	665.83	1489.0	129.52
2	54.785	1715.4	94.357	633.1	1163.5	33.748
3	210.16	281.09	77.387	735.12	876.85	9.1784

Table 1 Integrated autocorrelation time (*IAT*) and effective sample size (*ESS*) for 10^6 units long runs of the single-, multi-processor anytime Parallel Tempering and standard MCMC algorithms for different computational complexities $p \in \{0, 1, 2, 3\}$.

5.2 ABC toy example: univariate Normal distribution

To validate the results of Section 4.2, consider another simple example, initially featured in Lee [2012], and adapted here within the anytime parallel tempering Monte Carlo framework. Let Y be a Gaussian random variable, i.e. $Y \sim \mathcal{N}(y; \theta, \sigma^2)$, where the standard deviation σ is known but the mean θ is not. The ABC likelihood here is

$$f^\varepsilon(y | \theta) = \Phi\left(\frac{y + \varepsilon - \theta}{\sigma}\right) - \Phi\left(\frac{y - \varepsilon - \theta}{\sigma}\right)$$

for $\varepsilon > 0$. Using numerical integration tools in MATLAB, it is possible to obtain a good approximation of the true posterior for any ε for visualisation. Let $y = 3$ be an observation of Y and $\sigma^2 = 1$, and put the prior $p(\theta) = \mathcal{N}(\theta; 0, 5)$ on θ . In this example, the exact posterior distribution for θ can easily be shown to be $\mathcal{N}(\theta; \frac{5}{2}, \frac{5}{6})$.

When performing local moves (Algorithm 3), use a Gaussian random walk proposal with standard deviation $\xi = 0.5$. The real-time Markov jump process is run using $\Lambda = 10$ chains. The algorithm is run on a single processor for one hour or $T = 3600$ seconds in real time after a 30 second burn-in, with exchange moves occurring every $\delta_T = 5 \times 10^{-4}$ seconds (or 0.5 milliseconds). The radii of the balls $\varepsilon^{1:\Lambda}$ are defined to vary between $\varepsilon^1 = 0.1$ and $\varepsilon^\Lambda = 1.1$.

5.2.1 Verification of bias correction

First of all, we verify that bias correction must be applied for all chains to converge to the correct posterior. This is done by comparing density estimates of each of the post burn-in chains to the true corresponding posterior (obtained by numerical integration). When bias correction is not applied, meaning that the ABC-APTMC algorithm is run including the currently working chain j in Steps ??-?? of Algorithm 6, every chain converges to an erroneous distribution which

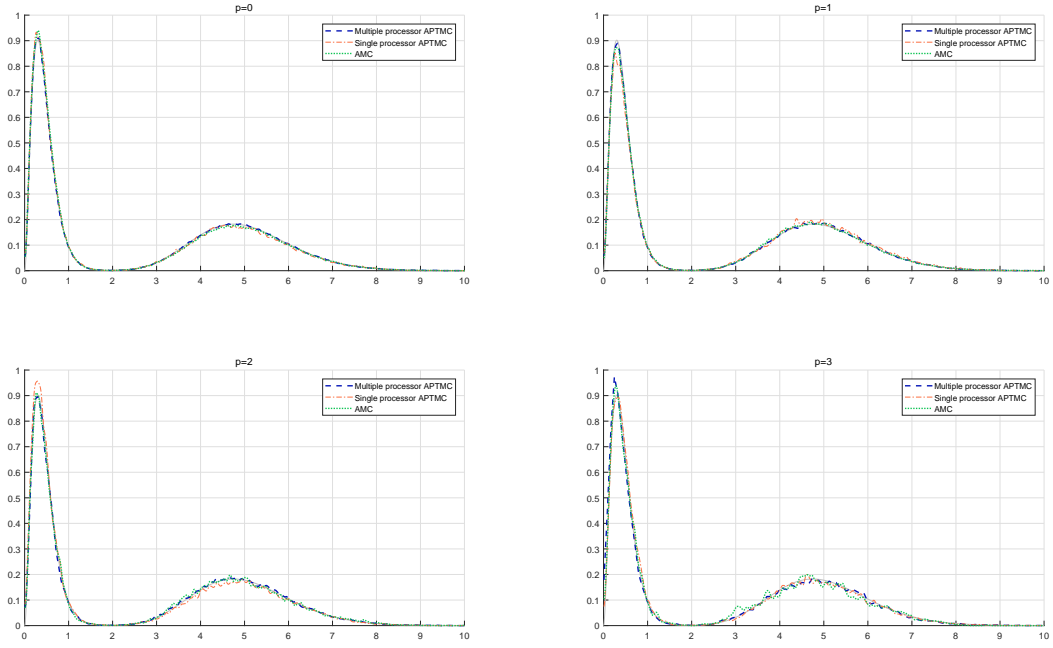


Fig. 3 Plots of kernel density estimates of the cold posterior for 10^6 units long runs of the single (*orange*) and multiple (*blue*) processor Anytime Parallel Tempering Monte Carlo (APTMC) algorithm as well as the standard (*green*) (MCMC) algorithm. Each plot corresponds to a different computational complexity $p \in \{0, 1, 2, 3\}$.

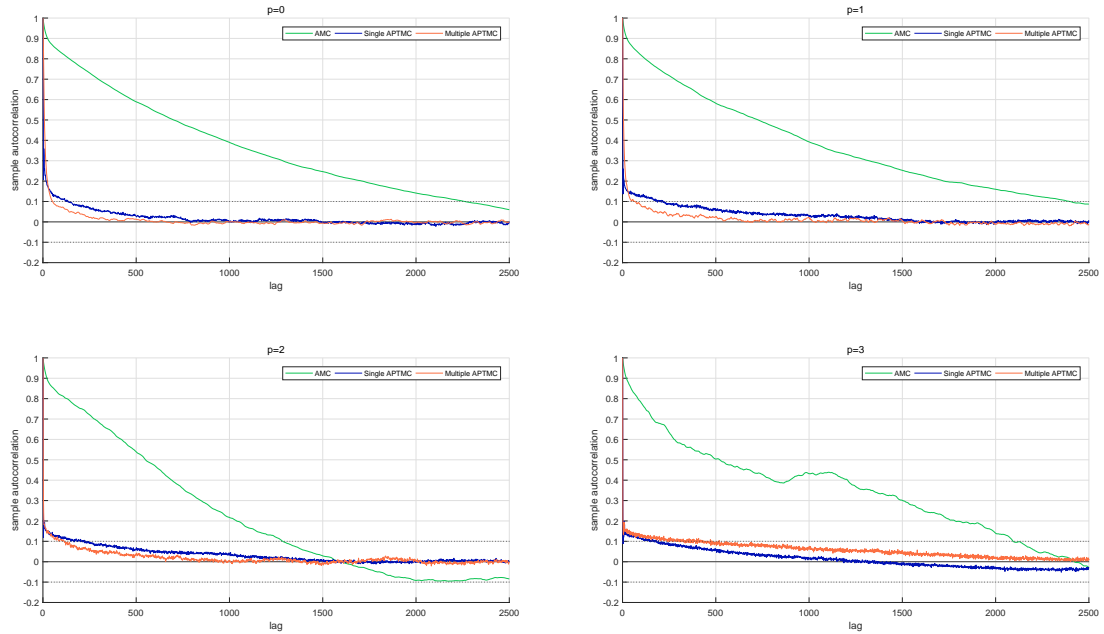


Fig. 4 Plots of the sample autocorrelation function up to lag 2500 of the post burn-in cold chain for runs of the single (*orange*) and multiple (*blue*) processor Anytime Parallel Tempering Monte Carlo (APTMC) algorithm as well as for the output of the standard Anytime Monte Carlo (AMC) algorithm (*green*). Each plot corresponds to a different computational complexity $p \in \{0, 1, 2, 3\}$.

overestimates the mode of its corresponding posterior, as is clearly visible in Figure 5. On the other hand, correcting the algorithm for such bias ensures that every chain converges to the correct corresponding posterior.

Next, we compare the performance of the ABC-APTMC algorithm to that of a standard ABC algorithm. For that, a more applied parameter estimation example is considered, for which the adoption of a likelihood-free approach is beneficial.

5.3 Moving average process

To illustrate a possible application of the ABC-APTMC algorithm, we now consider a common parameter estimation example taken from Marin et al. [2012]. A *Moving Average* process of order q , or $\text{MA}(q)$ process, is used in time series analysis to model serial autocorrelation for the stochastic process $y = (y_m)_{m \in \mathbb{N}}$ up to lag q . Consider the expression

$$y_m = u_m + \sum_{i=1}^q \theta_i u_{m-i} \quad (8)$$

where $u_m \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ for $m = 1, 2, \dots$. In this example, we aim to estimate the posterior of the parameters $\theta = (\theta_1, \dots, \theta_q)$.

5.3.1 Prior

In time series analysis, a standard invertibility condition is the following:

Condition 51 *The roots of the polynomial*

$$\mathcal{Q}(x) = 1 - \sum_{i=1}^q \theta_i x^i$$

all lie outside the unit circle in the complex plane.

We can therefore define a uniform prior over the permitted range of θ_i 's. In the case $q = 2$ this corresponds to sampling uniformly from the triangle

$$-2 < \theta_1 < 2, \quad \theta_1 + \theta_2 > -1, \quad \theta_1 - \theta_2 < 1.$$

5.3.2 Likelihood

From Hamilton [1994], the likelihood of the observed sequence $y = (y_1, \dots, y_M)$ originating from a $\text{MA}(q)$

process with parameters $\theta = (\theta_1, \dots, \theta_q)$ is available as a multivariate Gaussian of the form

$$f(y | \theta) = \frac{1}{\sqrt{(2\pi)^M |\Omega|}} \exp \left[-\frac{1}{2} y^\top \Omega^{-1} y \right]$$

where Ω is the $M \times M$ variance-covariance matrix

$$\Omega = \mathbb{E}[y^\top y] = \begin{pmatrix} \gamma_0 & \gamma_1 & \gamma_2 & \dots & \gamma_q & \dots & 0 \\ \gamma_1 & \gamma_0 & \gamma_1 & \ddots & \gamma_{q-1} & \ddots & \vdots \\ \gamma_2 & \gamma_1 & \ddots & \ddots & \ddots & \ddots & \gamma_q \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \gamma_{q-1} \\ \gamma_q & \ddots & \ddots & \ddots & \gamma_1 & \gamma_0 & \gamma_1 \\ \vdots & \ddots & \ddots & \ddots & \gamma_1 & \gamma_0 & \gamma_1 \\ 0 & \dots & \gamma_q & \dots & \gamma_2 & \gamma_1 & \gamma_0 \end{pmatrix}$$

where setting $s = |r|$ for $r = -q, \dots, q$ and $\theta_0 = 1$, the covariance for lag s is

$$\gamma_s = \sum_{i=0}^{q-s} \theta_i \theta_{i+s}$$

In practice, when M is large, the computational cost of dealing with the matrix Ω becomes prohibitive. Other means of computing the likelihood have been devised, notably in Marin and Robert [2007], but the use of ABC provides an easy, likelihood-free approach to estimate the parameters θ .

5.3.3 ABC-MCMC approach

Instead of evaluating the full likelihood at each iteration of local and exchange moves, we simulate the $\text{MA}(q)$ process $(x_m)_{m=1}^M$ and evaluate its distance to the observations $(y_m)_{m=1}^M$. There are multiple ways to evaluate such a distance: it can for example be the raw distance between the two datasets, but in general it is better to consider the distance between conveniently chosen summary statistics. More on the choice of summary statistics can be found in Marin et al. [2014]. In this project, we evaluate the quadratic distance between the first q sample autocorrelations, i.e. for $j = 1, \dots, q$ compute

$$\rho_j(y) = \frac{1}{\tau_0} \sum_{m=j+1}^M y_m y_{m-j} \quad (9)$$

where $\tau_0 = \sum_{m=1}^M y_m^2$ to obtain the vector $\rho(y) = (\rho_1(y), \dots, \rho_q(y))$ of summary statistics. Therefore, the vector $\rho(x)$ 'hitting' the ball $B_\varepsilon(\rho(y))$ of radius ε is here equivalent to $\|\rho(y) - \rho(x)\|_2 \leq \varepsilon$. The 1-hit MCMC kernel adapted for a $\text{MA}(q)$ process is detailed in Algorithm 8. Similar modifications are made to the exchange moves in Algorithm 5.

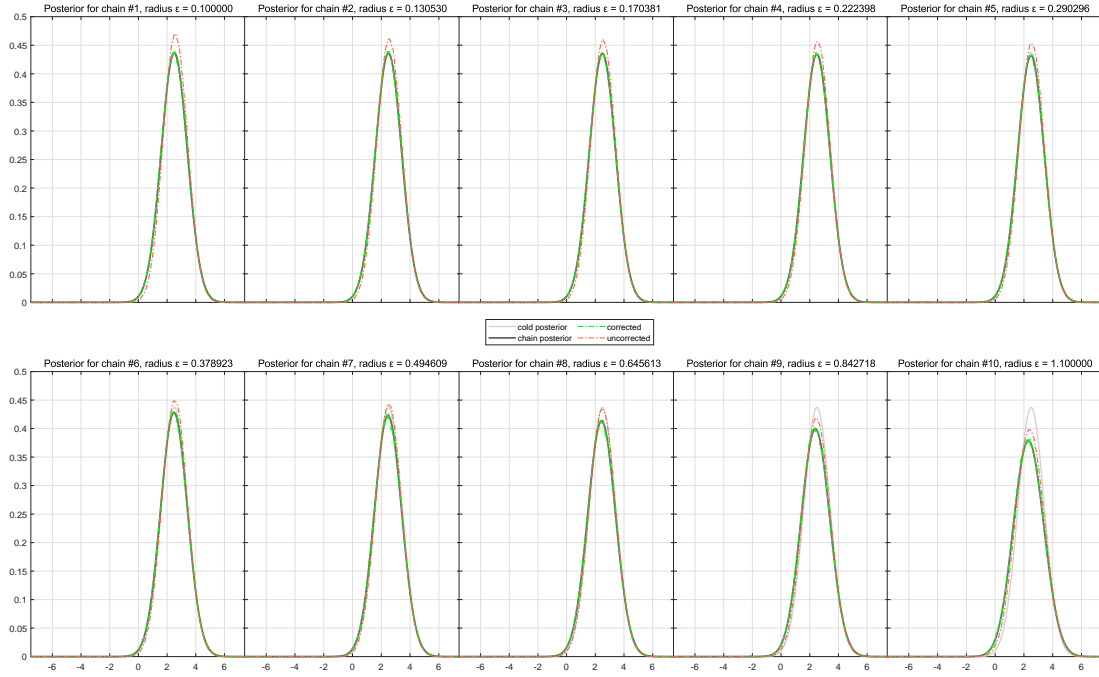


Fig. 5 Kernel density estimates of all chains for corrected and uncorrected runs of the single processor ABC-APTM algorithm. In each subplot, the *light gray* line is fixed and represents the cold posterior for reference, the *dark gray* line represents each chain's target posterior (obtained by numerical integration), the dot-dashed *green* lines are kernel density estimates of the chain's posterior returned by the corrected algorithm and agree with the dark gray line. The *orange* lines are kernel density estimates for the uncorrected algorithm, and do not agree with the gray line, as expected.

5.3.4 Methods and settings

In this example, the MA(2) process is considered for easy visualisation, i.e. $q = 2$. A sample of $M = 500$ observations $y = (y_1, \dots, y_M)$ is simulated using parameters $\theta = (\theta_1, \theta_2) = (0.6, 0.2)$. In this case, the true likelihood can be computed without too much effort, and hence the true joint and marginal posteriors can be approximated using numerical integration in MATLAB.

This example serves to illustrate a possible application of the ABC-APTM algorithm and especially the computational benefits of introducing ABC exchange moves. Therefore, we also include the standard version of the parallel tempering algorithm (ABC-PTMC), in which exchange moves are performed after a fixed amount of local moves instead of following a real-time schedule. For completeness, the ABC-APTM algorithm with no cold updates, i.e. no local moves occurring on the cold chain, is included as well. To check the algorithms yield the correct results, the ABC-APTM and standard 1-hit kernel ABC algorithm with (ABC-PTMC) and without (standard ABC) exchange moves are run eight times for $T = 10800$ seconds (or 3 hours) of real

time after an initial $t_{\min} = 1800$ second period of burn-in. The algorithms use a Gaussian random walk proposal with standard deviation $\xi = 0.25$ in the standard ABC case, and $\xi^{1:A}$ varying between $\xi^1 = 0.1$ and $\xi^A = 1$ in the ABC-APTM and ABC-PTMC cases. Additionally, the ABC-PTMC and ABC-APTM algorithms are run on $A = 10$ chains with exchange moves occurring every 15 local moves in the standard version and every 0.12 seconds in the anytime versions. These are chosen in order to ensure both algorithms spend the same median time performing local moves (see Section 5.4). The radii of the balls $\varepsilon^{1:A}$ are set to vary between $\varepsilon^1 = 0.02$ and $\varepsilon^1 = 1$. The standard 1-hit ABC kernel is run on a single chain, with the ball radius equal to ε^1 , i.e. the radius corresponding to the cold chain in the parallel tempering algorithms. We compare the efficiency of the parallel tempering ABC algorithms to that of the standard, single-chain one. For that, a sample acf plot averaged over all eight runs is drawn and the integrated autocorrelation time (IAT) and cumulative effective sample size (ESS) over all runs are computed for comparison.

5.3.4.1 Results

The scatter plots in Figure 6 indicate that the

Algorithm 8 ABC: 1-hit MCMC kernel for MA(q) process

Given current state (θ_n, x_n)

```

1: for  $i := 1, 2, \dots$  do
  SIMULATE THE REAL-TIME MARKOV JUMP PROCESS
   $(\theta, L)(t)$  UNTIL REAL TIME  $t_i$ .
2:   Propose  $\theta' \sim q(d\theta | \theta_n)$   $\triangleright$  propose a local move
3:   if  $\theta'$  satisfies Condition 51 then
4:     RACE := TRUE
5:   else
6:     RACE := FALSE  $\triangleright$  do not perform the race
7:     retain  $(\theta_{n+1}, x_{n+1}) = (\theta_n, x_n)$   $\triangleright$  automatically
      reject  $\theta'$ 
8:   end if
9:   while RACE do
10:    Simulate  $u_m$  and  $u'_m \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$  for  $m = 1, \dots, M$ 
11:    Generate the simulated stochastic processes  $x$  and
       $x'$  following Equation 8
12:    Compute summary statistics  $\rho(x)$  and  $\rho(x')$  ac-
      cording to Equation 9
13:    if  $\rho(x) \in B_\varepsilon(\rho(y))$  or  $\rho(x') \in B_\varepsilon(\rho(y))$  then
14:      RACE := FALSE  $\triangleright$  stop the race once either
       $\rho(x)$  or  $\rho(x')$  hits the ball
15:    end if
16:  end while
17:  if  $\rho(x')$  falls within  $\varepsilon$  of  $y$  first, or at the same time
      as  $\rho(x)$  then  $\triangleright$  accept or reject move
18:    set  $(\theta_{n+1}, x_{n+1}) = (\theta', x')$ 
19:  else
20:    retain  $(\theta_{n+1}, x_{n+1}) = (\theta_n, x_n)$   $\triangleright$  retain the
      initial  $x_n$  and not the  $x$  simulated during the race
21:  end if
22:   $n := n + 1$ 
23: end for

```

ABC-APTMC algorithm performs as expected, with each chain targeting a distribution increasingly close to the true posterior as ε decreases. However, as was already the case in Marin et al. [2012], the ABC approximation fails to reconstruct the posterior perfectly, even for ε as low as 0.02. Several ways to improve results are discussed in Marin et al. [2012], including decreasing ε further and considering alternative summary statistics for comparison with the observations, as well as applying corrections to the ABC output such as in Beaumont et al. [2002]. After running for 3 hours, all three algorithms return indistinguishable posteriors. What remains to be compared is the performance of each algorithm.

5.3.5 Performance evaluation

First of all, it should be noted that while all three algorithms run for the same amount of time, the standard ABC algorithm only has one chain to update while the parallel tempering algorithms must perform local moves on 10 (or 9) chains in sequence, and will hence return a cold chain with fewer samples despite the addition of exchange moves. This means that here again, a relevant comparison tool is the integrated

autocorrelation time (IAT). It is also important to verify that the parallel tempering algorithms return larger average effective sample sizes (ESS) after running for 3 hours.

The IAT values for both θ_1 and θ_2 in Table 2 are lower for all parallel tempering algorithms. As expected, the ABC-PTMC algorithm with no cold updates returns the smallest IAT values, since its cold chain is exclusively made up of samples from exchange moves. The most inefficient of the three algorithms is clearly the ABC algorithm, as it needs on average 2.7 to 4.9 times more samples to obtain the equivalent of an independent draw than it does for the parallel tempering algorithms. This is further supported by the sample acf plots in Figure 7, which display a steeper decay in sample acf for the three parallel tempering algorithms. The ESS values in Table 2 indicate that in eight 3 hour runs, the standard ABC algorithm may have output a greater overall number of samples on average, but it is so inefficient compared to the others that its resulting effective sample sizes for both θ_1 and θ_2 are still up to 6.7 times smaller on average than those of the parallel tempering algorithms.

This example serves to illustrate the improvements in performance brought by the addition of exchange moves. Indeed, the ABC-PTMC and ABC-APTMC parallel tempering algorithms display a similar performance both in the IAT values of Table 2, and in the almost indistinguishable acf decay in Figure 7. In fact, what indicates that running the algorithm within the anytime framework improved results is the much larger ESS returned by the ABC-APTMC algorithm. Even the version of the algorithm that didn't perform any local updates returns a higher ESS . We note that the multi-processor parallel tempering algorithms were not considered in this example. The next section considers a more advanced case study in which the likelihood is unavailable, and the benefits of performing exchange moves within the anytime framework are illustrated both on a single and multiple processors.

5.4 Stochastic Lotka-Volterra model

The utility of the anytime parallel tempering framework in a situation where the computation of the likelihood is impossible or prohibitive has not yet been demonstrated. Indeed, the moving average example in Section 5.3 only had $q = 2$ components and $M = 500$ data points, making the likelihood not too costly to compute. In this section, we consider the stochastic Lotka-Volterra predator-prey model (Lotka [1926], Volterra

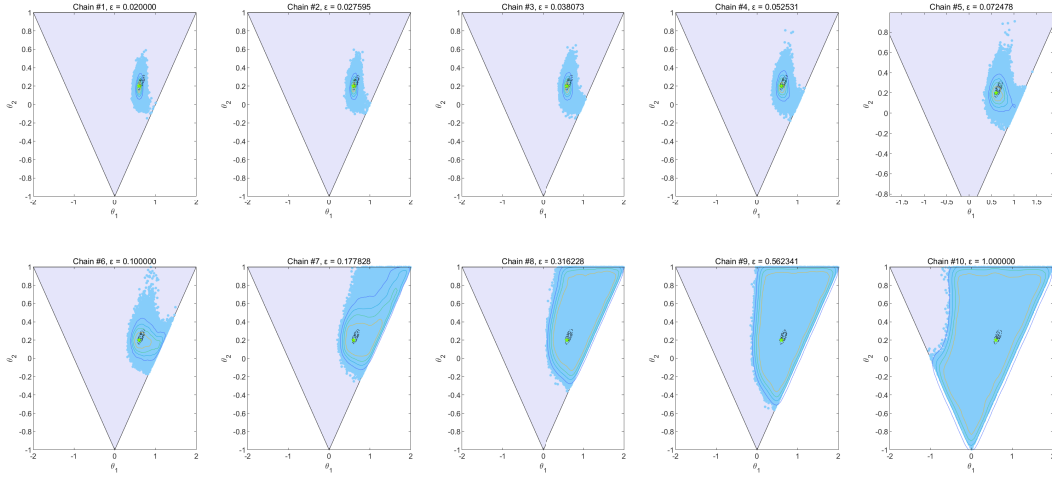


Fig. 6 Scatter plots of post burn-in chains for a run of the single processor ABC-APTMC algorithm. In each subplot, the *black* contour lines represent the level sets of the true posterior and the multicoloured contour lines a bivariate Gaussian kernel density estimate from the samples returned for each chain by the algorithm, obtained using the `gkde2` function in MATLAB. The *yellow* star represents the true value $\theta = (0.6, 0.2)$ and the triangle is the range of acceptable values of θ .

	ABC-APTMC		ABC-APTMC no cold updates		ABC-PTMC		ABC-MCMC Standard	
	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>
θ_1	7.15	205472	4.52	80372	7.14	68771	22.81	30844
θ_2	10.02	146632	4.76	76292	9.94	49388	22.30	31558

Table 2 Mean values of *IAT* and cumulative *ESS* over eight runs of the standard ABC and single processor ABC-APTMC and ABC-PTMC algorithms to estimate the posterior distributions of the parameters $\theta = (\theta_1, \theta_2)$ of a MA(2) process.

[1927]), further exploring the example considered in Lee and Łatuszyński [2014], which is itself based on an example in Chapter 6 of Wilkinson [2011]. In this case study, the posterior is intractable and some of the components of the parameters θ (namely θ_2 and θ_3) exhibit strong correlations. Let $X_{1:2}(t)$ be a bivariate, integer-valued pure jump Markov process with initial values $X_{1:2}(0) = (50, 100)$, where $X_1(t)$ represents the number of preys and $X_2(t)$ the number of predators at time t . For small time interval Δt , we describe the predator-prey dynamics in the following way

$$\begin{aligned} & \mathbb{P}\{X_{1:2}(t + \Delta t) = z_{1:2} \mid X_{1:2}(t) = x_{1:2}\} \\ &= \begin{cases} \theta_1 x_1 \Delta t + o(\Delta t), & \text{if } z_{1:2} = (x_1 + 1, x_2) \\ \theta_2 x_1 x_2 \Delta t + o(\Delta t), & \text{if } z_{1:2} = (x_1 - 1, x_2 + 1) \\ \theta_3 x_2 \Delta t + o(\Delta t), & \text{if } z_{1:2} = (x_1, x_2 - 1) \\ o(\Delta t), & \text{otherwise} \end{cases} \end{aligned}$$

In this example, the only observations available are the number of preys, i.e. X_1 at 10 discrete time points. Following theory in Wilkinson [2011] (Chapter 6), the process can be simulated and discretised using

the Gillespie [1977] algorithm, in which the inter-jump times follow an exponential distribution. The observations employed were simulated in Lee and Łatuszyński [2014] with true parameters $\theta = (1, 0.005, 0.6)$, giving $y = \{88, 165, 274, 268, 114, 46, 32, 36, 53, 92\}$ at times $\{1, \dots, 10\}$. For ABC, the ‘ball’ considered takes the following form for $\varepsilon > 0$

$$\begin{aligned} & B_\varepsilon(y) \\ &= \{X_1(t) : |\log[X_1(i)] - \log[y(i)]| \leq \varepsilon, \forall i = 1, \dots, 10\} \end{aligned} \quad (10)$$

therefore, a set of simulated $X_1(t)$ is considered as ‘hitting the ball’ if all 10 simulated data points are at most e^ε times (and at least $e^{-\varepsilon}$ times) the magnitude of the corresponding observation in y .

5.4.1 Methods and settings

In Lee and Łatuszyński [2014], the 1-hit MCMC kernel (ABC), is shown to return the most reliable results by comparison with other MCMC kernels which are not considered here. While it can be reasonably fast, it is highly inefficient as it has a very low acceptance rate,

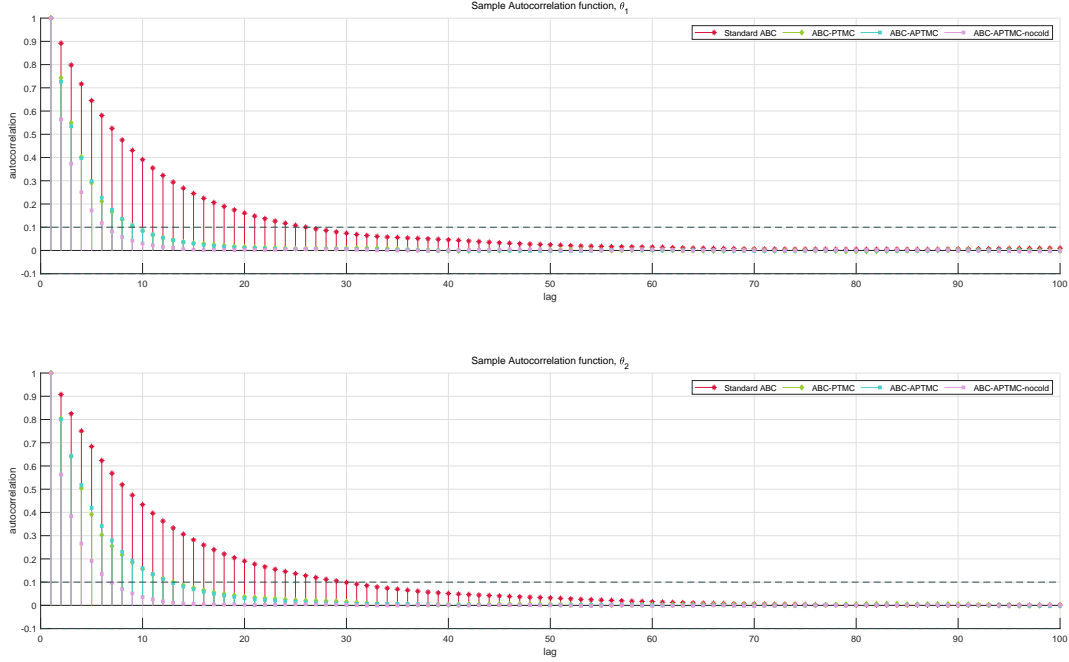


Fig. 7 Plots of the average sample autocorrelation function up to lag 100 of the post burn-in cold chain for eight 3-hour long runs of the standard ABC (red) and single processor ABC-PTMC (green) and ABC-APTMC (blue) algorithms to estimate the posterior distributions of the parameters $\theta = (\theta_1, \theta_2)$ of a MA(2) process.

and thus the autocorrelation between samples for low lags is very high. Another important issue in this particular example is that the race in the 1-hit kernel is prone to getting stuck for extended periods of time. Therefore, we aim to first of all improve performances by introducing exchange moves on a single processor (ABC-PTMC). Then – and most importantly – we further improve the algorithm by implementing both the single and multi-processor parallel tempering algorithms within the anytime framework (ABC-APTMC).

5.4.1.1 One processor

Departing slightly from the example in Lee and Łatuszyński [2014], define the prior on $\theta \in [0, \infty)^3$ for the single processor experiment to be $p(\theta) = \exp\{-\theta_1 - \theta_2 - \theta_3\}$, i.e. three independent exponential priors, all with mean 1. The proposal distribution is a truncated normal, i.e. $\theta' | \theta \sim \mathcal{TN}(\theta, \Sigma)$, $\theta' \in (0, 10)$ with mean θ and covariance $\Sigma = \text{diag}(0.25, 0.0025, 0.25)$. The truncated normal is used in order to ensure that all proposals remain non-negative. For reference, 2364 independent samples from the posterior are obtained via ABC rejection sampling with $\varepsilon = 1$ and the density estimates in Figure 6 of Lee and Łatuszyński [2014] are reproduced.

To obtain these posterior samples, 10^7 independent samples from the prior were required, yielding the very low 0.024% acceptance rate. This method of sampling from the posterior is therefore extremely inefficient, and the decision to resort to MCMC kernels in order to improve efficiency is justified.

On a single processor, the three algorithms considered are the standard 1-hit MCMC kernel (ABC), the single processor version of the algorithm with added exchange moves (ABC-PTMC-1) and the same but within the anytime framework (ABC-APTMC-1). They are run ten times for 100800 seconds (28 hours) – after 3600 seconds (1 hour) of burn-in – and their main settings are summarised in Table 3. It is important to note that the parallel tempering algorithms, having to deal with updating multiple chains sequentially, are likely to return fewer samples from the cold chain. The algorithms must therefore be properly set up such that the gain in efficiency introduced by exchange moves is not overshadowed by the greater number of chains. In this experiment, this means the parallel tempering algorithms must be run on just 6 chains, each targeting posteriors associated with balls of radii $\varepsilon^{1:6} = \{1, 1.1447, 1.3104, 1.5, 11, 15\}$ and the proposal distribution has covariance

$\Sigma^{1:6}$ where $\Sigma^\lambda = \text{diag}(\sigma^\lambda, \sigma^\lambda 10^{-2}, \sigma^\lambda)$ and $\sigma^{1:6} = \{0.008, 0.025, 0.05, 0.09, 0.25, 0.5\}$. Exchange moves are performed as described in Algorithm 5. In the anytime version, in order to determine how long the local moves should run for before exchange moves occur, the time taken for the standard algorithm to perform a fixed number δ_T of local moves is measured at each iteration, and the median over all iterations is taken. This ensures that both algorithms spend the same median time performing local moves.

5.4.1.2 Multiple processors

Then, we aim to demonstrate the gain in efficiency introduced by running the parallel tempering algorithm, not only within the anytime framework, but also on multiple processors. The algorithms considered are the single processor ABC-PTMC-1, ABC-APTMC-1 as well as their multi-processor counterparts ABC-PTMC-W and ABC-APTMC-W. This time, instead of relying on an informative, exponential prior on θ , which leads to lower acceptance rates on the warmer chains, we define a uniform prior between 0 and 3. The proposal distribution is still a truncated normal, but with tighter limits (corresponding to the prior) i.e. $\theta' | \theta \sim TN(\theta, \Sigma)$, $\theta' \in (0, 3)$. Again, 1988 independent samples from the posterior are obtained for reference. They are generated via ABC rejection sampling with $\varepsilon = 1$. To obtain these posterior samples, 10^8 independent samples from the prior were required, yielding the even lower 0.002% acceptance rate.

Since the standard ABC algorithm is not considered here, there is no need to balance gain in efficiency with the reduced number of samples on the cold chain present because of the multiple chains. Therefore, the four algorithms are run on 20 chains, each targeting posteriors associated with balls of radii ranging from $\varepsilon^1 = 1$ to $\varepsilon^{20} = 11$ and proposal distribution covariances $\Sigma^{1:20}$ where $\Sigma^k = \text{diag}(\sigma^k, \sigma^k 10^{-2}, \sigma^k)$ for chain k and where values range from $\sigma^1 = 0.008$ to $\sigma^{20} = 0.5$. These are tuned so that the acceptance rates of exchange moves between adjacent chains remain on average greater than 70%. The algorithms are run four times for 864000 seconds (24 hours) and their main settings are summarised in Table 4.

In this experiment, multiple chains running at different temperatures are present on each worker. In order to account for the approx. 1 second communication overhead when switching between local moves – running in parallel on all workers – and exchange moves –

running on the master –, exchange moves are divided into two types:

1. Exchange moves *within workers*: performed on each individual worker in parallel, between a pair of adjacent chains selected at random. No communication between workers is necessary in this case.
2. Exchange moves *between workers*: performed on the master by selecting a pair of adjacent workers at random and exchanging between the warmest eligible chain from the first worker and coldest from the second. Thus, an exchange move between two adjacent chains is effectively being performed, except this time communication between workers is required.

After (parallel) local moves are performed, each iteration alternates between *within* and *between workers* exchange moves. Communication between workers is therefore only needed once every two iterations, thus reducing the total communication overhead. Determining how long the workers in the anytime version of the algorithm should work in parallel before between-worker exchange moves occur is adapted to the new algorithm construction. The time (excluding communication overhead) taken for the ABC-PTMC-W algorithm to perform a set of parallel updates – i.e. δ_T local moves, an exchange move within workers and δ_T more local moves – is measured on each worker and at each iteration. Then, the median over all iterations is taken for each worker, and the ABC-APTMC-W algorithm is set to perform parallel updates for as long as the slowest of all workers. This ensures both that all workers have a chance to perform the δ_T local moves, and that workers containing chains which are quicker to complete their local moves do not sit idle waiting for the worker containing the slowest chains to finish. It was observed that the slowest workers were most often the ones containing the colder half of the chains.

5.4.2 Performance evaluation

All algorithms returned density estimates that were reasonably close to those obtained via rejection sampling, but all would require to run for a much longer period of time to be indistinguishable from them, which is why in this experiment we chose to focus mainly on comparing their efficiency. In order to compare the performance of the algorithms, as stated above, all algorithms compared are set to run for the same real time period. Once again the integrated autocorrelation time (*IAT*) and effective sample size (*ESS*) are computed for all algorithms.

<i>Label</i>	<i>Workers</i> <i>W</i>	<i>Chains</i> <i>K</i>	<i>Chains per</i> <i>worker</i>	<i>Exchange moves</i>	<i>Anytime</i>
ABC	1	1	1	none	No
ABC-PTMC-1	1	6	6	every 6 local moves	No
ABC-APTMC-1	1	6	6	every 2.59 seconds	Yes

Table 3 Algorithm settings for stochastic Lotka-Volterra predator-prey model on a single processor.

<i>Label</i>	<i>Workers</i> <i>W</i>	<i>Chains</i> <i>K</i>	<i>Chains per</i> <i>worker</i>	<i>Exchange moves</i>	<i>Anytime</i>
ABC-PTMC-1	1	20	20	every 20 local moves	No
ABC-APTMC-1	1	20	20	approx. every 10.3 seconds	Yes
ABC-PTMC-W	4	20	5	every 5 local moves	No
ABC-APTMC-W	4	20	5	every 5 local moves (<i>within</i> workers) and approx. every 15.3 seconds (<i>between</i> workers)	Yes

Table 4 Algorithm settings for stochastic Lotka-Volterra predator-prey model on multiple processors.

While the *IAT* and sample autocorrelation plots are good tools for comparing efficiency, they do not take into account the computational cost of running 6 chains instead of a single ones. The *ESS* on the other hand gives us how many effective samples the different algorithms can return within a fixed time frame. For example, a very fast algorithm could still return a higher *ESS* even if it has a much higher *IAT*. To illustrate how the anytime version of the parallel tempering algorithms is more computationally efficient compared to standard ABC-PTMC, the real times all algorithms take to perform local and exchange moves are measured and their timelines plotted.

5.4.2.1 One processor

Both the ABC-PTMC-1 and ABC-APTMC-1 algorithm display an improvement in performances: they return *IAT*s that are respectively 3.55 and 2.16 times lower on average than those of the standard ABC algorithm in Table 5, and display a steeper decay in sample acf in Figure 8. In the 28 hours (post burn-in) during which the algorithms ran, both parallel algorithms also yielded an increased effective sample size. As a matter of fact, introducing exchange moves every 6 local moves has multiplied the *ESS* by approximately 1.173 on average, and performing them every 2.59 seconds within the anytime framework has further increased it, returning effective sample sizes that are approximately 2.09 times as large as that of the standard ABC algorithm. We note that in this example, the ABC-PTMC-1 returned a lower *IAT* than its anytime counterpart. This is due to the low number of chains needed in this experiment. The 6 radii corresponding to each chain are further apart than they would usually be, and

because of this, the proportion of rejected exchange moves in the ABC-APTMC-1 algorithm is higher, as the construction of the anytime algorithm often requires one of the exchange moves to occur on a pair of chains that are not adjacent. On an experiment with a higher number of chains, this issue becomes negligible.

The median time spent performing local moves was the same for the ABC-PTMC-1 and ABC-APTMC-1 algorithm. However, the distribution of times spent performing local moves in the run of the ABC-APTMC-1 is more closely concentrated around the median. While many of the local moves in the run of the ABC-PTMC-1 algorithm took as little as 0.413 seconds to complete, others took over 40 seconds. Some local moves even took nearly two hours to complete. On the other hand, since a deadline was implemented in the anytime version of the algorithm, Figure 9 displays more consistent local move times.

	Standard ABC		ABC-PTMC-1		ABC-APTMC-1	
	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>
θ_1	67.993	7606.7	20.038	8529.5	36.797	13589
θ_2	118.07	4380.4	30.962	5520	49.086	10187
θ_3	145.39	3557.4	44.456	3844.5	68.419	7308.5

Table 5 Cumulative effective sample size (*ESS*) and mean integrated autocorrelation time (*IAT*) over ten 28-hour runs of the ABC, ABC-PTMC-1 and ABC-APTMC-1 algorithms to estimate the posterior distributions of the parameters $\theta = (\theta_1, \theta_2, \theta_3)$ of a stochastic Lotka-Volterra model.

5.4.2.2 Multiple processors

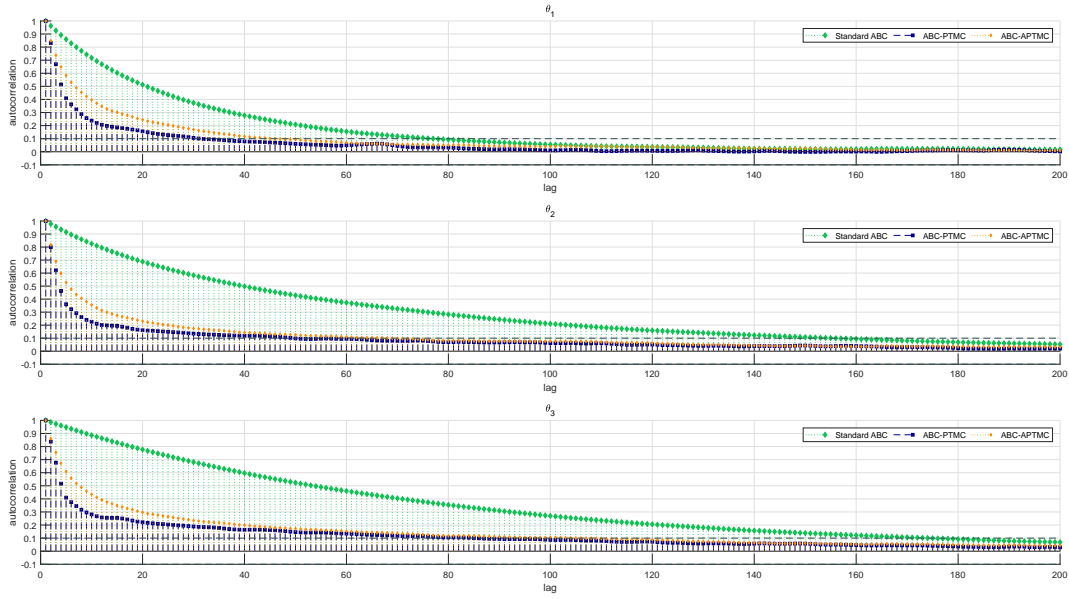


Fig. 8 Plots of the sample autocorrelation function up to lag 200 of the cold chain for runs of the standard ABC (*green*), single processor ABC-PTMC-1 (*blue*) and multi-processor ABC-PTMC-1 (*orange*) algorithms to estimate the posterior distributions of the parameters $\theta = (\theta_1, \theta_2, \theta_3)$ of a stochastic Lotka-Volterra model.

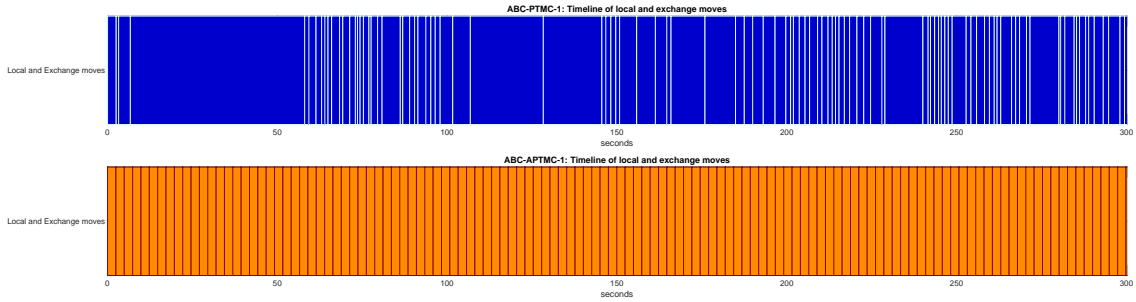


Fig. 9 Timeline of local and exchange moves for the ABC-PTMC-1 and ABC-APTMC-1 algorithms for the first 300 seconds. The exchange moves are represented by the *white* and *red* lines and the local moves by the *dark blue* and *orange* coloured blocks.

In the multi-processor case study, both the ABC-PTMC-1 and ABC-PTMC-W were set so that on each worker an exchange move occurred after all chains had been updated locally once. In theory, this should naturally give an advantage to the multi-processor version, as the various chains working in parallel are able to perform more local moves in the given time, and exchange moves occur every 5 local moves instead of every 20. In practice, if the inter-worker communication overhead occurring once every two exchange moves is too long, it may occupy too much of the allocated time and thus cause the multi-processor version of the algorithm to return fewer samples, and hence fewer effective samples. Here, this issue was avoided: the total number of samples returned by the ABC-PTMC-W algorithm is much higher for all chains (see Table 7)

and therefore the approx. 1 second communication overhead was not prohibitive. Additionally, the cold chain returned by the ABC-PTMC-W algorithm contains approximately twice as many samples originating from exchange moves as that returned by its single processor counterpart ABC-PTMC-1 (see Table 8) and has an integrated autocorrelation time that is on average 1.92 times smaller, as suggested in Table 6 and by the steeper decay in sample acf in Figure 10. Because of this, Table 6 shows that the effective sample sizes for the cold chain returned by the ABC-PTMC-W algorithm are on average 1.74 times higher than those of the single processor version.

The same observations can be made for the anytime versions of the algorithm, though note that the real

time deadline implemented means that the number of local and within-worker exchange moves occurring before a between-worker exchange move varies, and therefore it is possible for two exchange moves to occur consecutively (i.e. within workers and then between workers). The cold chain returned by the ABC-APTMC-W algorithm has an integrated autocorrelation time that is on average 1.72 times smaller than its single processor counterpart ABC-APTMC-1, as suggested in Table 6. Comparing the percentages of exchange moves present in the various chains, columns ABC-APTMC-1 and ABC-APTMC-W in Table 8 indicate that these percentages have doubled to tripled after switching to multiple processors. Even more striking is the great increase in total sample size returned for each chain. Indeed, Table 7 shows that while the number of samples on the cold chain is already 2.21 times higher, we observe by switching to multiple processors a 10-fold increase in the number of samples on warmer chains. Because of this, the effective sample sizes for the cold chain returned by the ABC-APTMC-W algorithm in Table 6 are on average 3.62 times higher than those of the single processor version.

As for the main comparison – namely anytime vs standard ABC with exchange moves – Table 6 indicates that the single processor ABC-APTMC-1 algorithm returns an effective sample size on average 2.26 times larger than the ABC-PTMC-1 algorithm. On multiple processors, the improvement is even greater, with the ABC-APTMC-W algorithm returning an effective sample size on average 4.82 times larger than the ABC-PTMC-W algorithm. Figures 11 and 12 illustrate the advantage of implementing a real-time deadline to local moves. At most local moves, the issue in which all workers sit idle waiting for the slowest to finish arises for the ABC-PTMC-W algorithm. This issue is most clearly visible on Figure 11 with Worker 2 from around 80 seconds on. On the other hand, Figure 12 clearly shows that the anytime version of the algorithm is making better use of the allocated computational resources. The anytime framework essentially ensures that all workers do not have to wait for the slowest among them to finish, allowing for more exploration of the sample space in the faster workers. Additionally, the real time deadline ensures that even if chain k on Worker w remains stuck in a race for an extended period of time, the other workers are still updating. Therefore, while the remaining four chains on Worker w wait for chain k to complete its race, they also continue to be updated at regular intervals thanks to the exchange moves with other workers. Because of this, while the integrated autocorrelation time in Table 6 for the cold chain does not seem to

be significantly lower for the ABC-APTMC-W algorithm, its sample size in Table 7 has quadrupled compared to that of the ABC-PTMC-W algorithm.

between the cold chains returned by the ABC-PTMC-W and ABC-APTMC-W algorithms, and by the roughly similar decay in sample acf in Figure 10.

The addition of ABC exchange moves in this case study proved fruitful, as the effective sample size for the parameters of the Lotka-Volterra model was increased. However this required fine tuning of the settings. The benefits of ABC parallel tempering will be stronger and more easily visible in a problem in which the parameters to be estimated have a multimodal distribution, as a single chain may get stuck in local optima while multiple tempered chains will explore more of the sample space. Nonetheless, the introduction of the anytime framework is clearly an important improvement. It ensures the various chains in ABC parallel tempering continue to be updated (via exchange moves) even when one of them is stuck performing local moves for longer than expected, and encourages the algorithm to make better use of its allocated resources, especially on multiple processors.

6 Conclusion

In an effort to increase the efficiency of MCMC algorithms, in particular for use on distributed computing, and for situations in which the likelihood is either unavailable or too computationally costly, the Anytime Parallel Tempering Monte Carlo algorithm was developed. The algorithm combines the enhanced exploration of the state space, provided by the between-chain exchange moves in parallel tempering, with control over the real-time budget and robustness to interruptions available within the anytime Monte Carlo framework.

Initially, the construction of the anytime Monte Carlo algorithm, with the inclusion of exchange moves on a single and multiple processors, was verified on a toy Gamma mixture example, and the performance improvements they brought were demonstrated by comparing the algorithm to a standard MCMC algorithm. Subsequently, the exchange moves were adapted for pairing with the 1-hit ABC-MCMC kernel, a simulation-based algorithm within the Approximate Bayesian computation (ABC) framework, which provides an attractive, likelihood-free approach to MCMC. The construction of the adapted ABC algorithm was verified using a simple univariate normal example. Then, the increased efficiency of the inclusion of exchange moves was demonstrated in comparison to that

	One processor				Multiple processors			
	ABC-PTMC-1		ABC-APTMC-1		ABC-PTMC-W		ABC-APTMC-W	
	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>	<i>IAT</i>	<i>ESS</i>
θ_1	50.444	297.94	40.041	748.08	30.992	439.93	22.153	2679.9
θ_2	37.376	339.35	36.852	726.38	22.139	580.53	22.909	2549
θ_3	65.253	238.84	53.809	509.99	26.566	482.88	31.073	1919.4

Table 6 Cumulative effective sample size (*ESS*) and mean integrated autocorrelation time (*IAT*) over four 24-hour runs of the ABC-PTMC-1, ABC-APTMC-1, ABC-PTMC-W and ABC-APTMC-W algorithms to estimate the posterior distributions of the parameters $\theta = (\theta_1, \theta_2, \theta_3)$ of a stochastic Lotka-Volterra model.

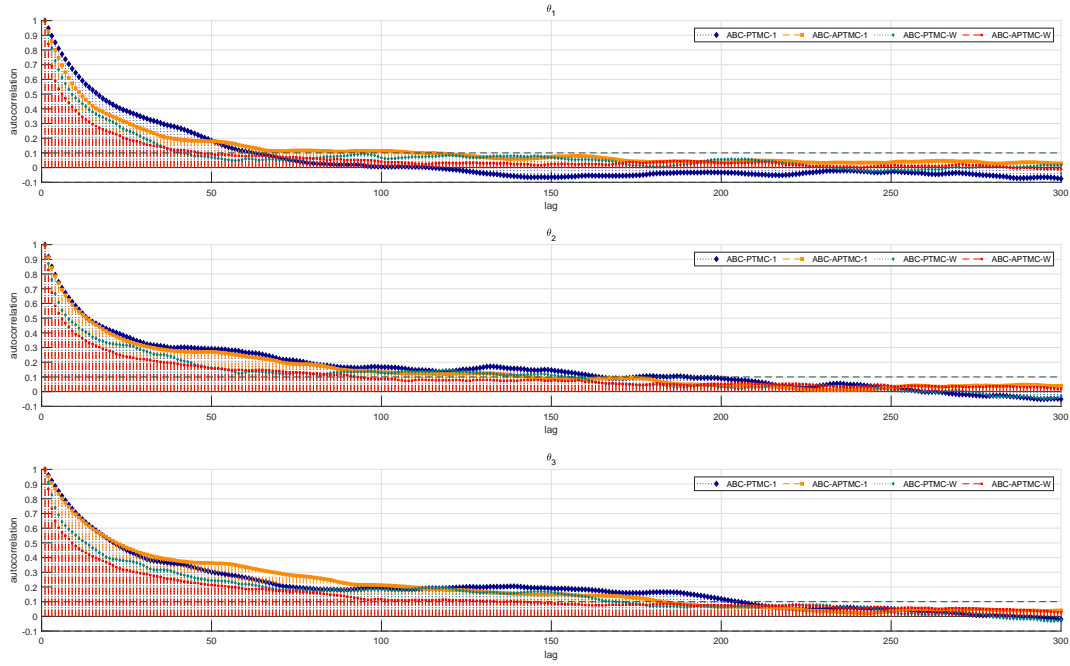


Fig. 10 Plots of the average sample autocorrelation function up to lag 290 of the cold chain for runs of the single processor ABC-PTMC-1 (*blue*) and ABC-APTMC-1 (*orange*) and multi processor ABC-PTMC-W (*teal*) and ABC-APTMC-W (*red*) algorithms to estimate the posterior distributions of the parameters $\theta = (\theta_1, \theta_2, \theta_3)$ of a stochastic Lotka-Volterra model.

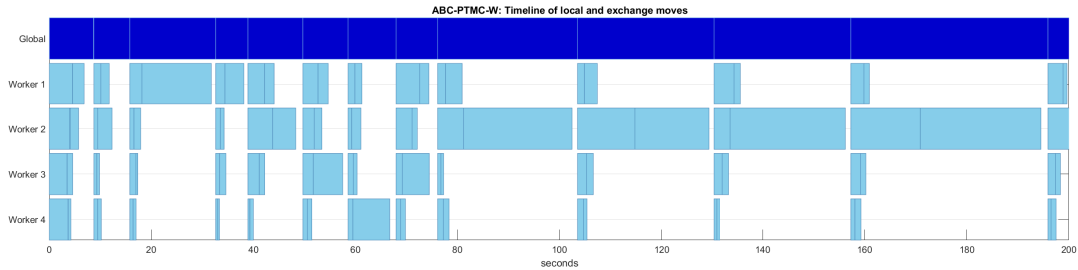


Fig. 11 Timeline of local and exchange moves for the ABC-PTMC-W algorithm for the first 200 seconds.

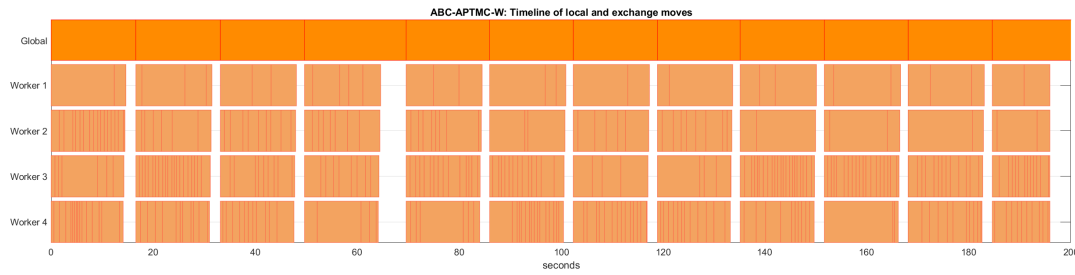


Fig. 12 Timeline of local and exchange moves for the ABC-APTMC-W algorithm for the first 200 seconds.

of a standard ABC algorithm on a parameter estimation problem involving the coefficients of a moving average $MA(q)$ process. Anytime parallel tempering ABC was finally employed to estimate the parameters of a stochastic Lotka-Volterra predator-prey model based on partial and discrete data – a problem in which the likelihood is unavailable. On a single processor, it was shown that introducing exchange moves provides an improvement in performance and an increase in the effective sample size compared to that of the standard, single chain ABC algorithm. This improvement is additionally boosted by the inclusion of the anytime framework. On multiple processors, it was shown that the anytime framework helps the parallel tempering ABC algorithm to make better use of the computational resources and thus provides an even stronger boost to effective sample size.

References

- R. Bardenet, A. Doucet, and C. Holmes. An adaptive subsampling approach for MCMC inference in large datasets. In *Proceedings of The 31st International Conference on Machine Learning*, pages 405–413, 2014.
- R. Bardenet, A. Doucet, and C. Holmes. On Markov chain Monte Carlo methods for tall data. *arXiv preprint arXiv:1505.02827*, 2015.
- M. A. Beaumont, W. Zhang, and D. J. Balding. Approximate Bayesian Computation in Population Genetics. *Genetics*, 162(4):2025–2035, 2002.
- Z. I. Botev, J. F. Grotowski, D. P. Kroese, et al. Kernel density estimation via diffusion. *The Annals of Statistics*, 38(5):2916–2957, 2010.
- B. Calderhead. A general construction for parallelizing Metropolis-Hastings algorithms. *Proceedings of the National Academy of Sciences*, 111(49):17408–17413, 2014.
- Foreman-Mackey, Daniel and Hogg, David W and Lang, Dustin and Goodman, Jonathan. emcee: The mcmc hammer. *Publications of the Astronomical Society of the Pacific*, 125(925):306, 2013.
- C. Geyer. Importance sampling, simulated tempering and umbrella sampling. *Handbook of Markov Chain Monte Carlo*, pages 295–311, 2011.
- C. J. Geyer. Markov chain Monte Carlo maximum likelihood. *Interface Foundation of North America*, 1991.
- D. T. Gillespie. Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry*, 81(25):2340–2361, 1977.
- J. Goodman and J. Weare. Ensemble samplers with affine invariance. *Communications in applied mathematics and computational science*, 5(1):65–80, 2010.
- J. D. Hamilton. *Time Series Analysis*, volume 2. Princeton university press Princeton, 1994.
- R. Kohn, M. Quiroz, M.-N. Tran, and M. Villani. Speeding up MCMC by Efficient Data Subsampling. Working Papers 2123/16205, University of Sydney Business School, Discipline of Business Analytics, 2016. URL <https://ideas.repec.org/p/syb/wpbsba/2123-16205.html>.
- A. Korattikara, Y. Chen, and M. Welling. Austerity in MCMC land: Cutting the Metropolis-Hastings budget. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 181–189, 2014.
- A. Lee. On the choice of MCMC kernels for approximate Bayesian computation with SMC samplers. In *Simulation Conference (WSC), Proceedings of the 2012 Winter*, pages 1–12. IEEE, 2012.
- A. Lee and K. Łatuszyński. Variance bounding and geometric ergodicity of markov chain monte carlo kernels for approximate bayesian computation. *Biometrika*, 101(3):655–671, 2014.
- F. Liang, J. Kim, and Q. Song. A bootstrap Metropolis-Hastings algorithm for Bayesian analysis of big data. *Technometrics*, 58(3):304–318, 2016.
- A. J. Lotka. Elements of Physical Biology. *Science Progress in the Twentieth Century (1919-1933)*, 21(82):341–343, 1926.
- D. Maclaurin and R. P. Adams. Firefly Monte Carlo: Exact MCMC with Subsets of Data. In *UAI*, pages 543–552, 2014.

- J.-M. Marin and C. P. Robert. *Bayesian core : a practical approach to computational Bayesian statistics*. Springer texts in statistics. Springer, New York, 2007. ISBN 9780387389790.
- J.-M. Marin, P. Pudlo, C. P. Robert, and R. J. Ryder. Approximate Bayesian computational methods. *Statistics and Computing*, pages 1–14, 2012.
- J.-M. Marin, N. S. Pillai, C. P. Robert, and J. Rousseau. Relevant statistics for Bayesian model choice. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(5):833–859, 2014.
- MATLAB. *version 7.10.0 (R2017a)*. The MathWorks Inc., Natick, Massachusetts, 2017.
- S. Minsker, S. Srivastava, L. Lin, and D. Dunson. Scalable and robust Bayesian inference via the median posterior. In *International Conference on Machine Learning*, pages 1656–1664, 2014.
- L. M. Murray, S. Singh, P. E. Jacob, and A. Lee. Anytime Monte Carlo. *arXiv preprint arXiv:1612.03319*, 2016. URL <https://arxiv.org/abs/1612.03319>.
- W. Neiswanger, C. Wang, and E. Xing. Asymptotically exact, embarrassingly parallel MCMC. *arXiv preprint arXiv:1311.4780*, 2013.
- J. K. Pritchard, M. T. Seielstad, A. Perez-Lezaun, and M. W. Feldman. Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Molecular biology and evolution*, 16(12):1791–1798, 1999.
- M. Quiroz, M. Villani, and R. Kohn. Exact subsampling MCMC. *arXiv preprint arXiv:1603.08232*, 2016.
- M. Quiroz, M.-N. Tran, M. Villani, and R. Kohn. Speeding up MCMC by delayed acceptance and data subsampling. *Journal of Computational and Graphical Statistics*, 2017.
- C. Robert and G. Casella. *Monte Carlo Statistical Methods*, chapter The Metropolis-Hastings Algorithm. Springer Texts in Statistics, Springer, New York, 2004. ISBN 978-1-4757-4145-2. doi: 10.1007/978-1-4757-4145-2_7.
- S. L. Scott, A. W. Blocker, F. V. Bonassi, H. A. Chipman, E. I. George, and R. E. McCulloch. Bayes and big data: The consensus Monte Carlo algorithm. *International Journal of Management Science and Engineering Management*, 11(2):78–88, 2016.
- A. Sokal. Monte Carlo Methods in Statistical Mechanics: Foundations and New Algorithms. In *Functional integration*, pages 131–192. Springer, 1997.
- R. H. Swendsen and J.-S. Wang. Replica Monte Carlo simulation of spin-glasses. *Physical Review Letters*, 57(21):2607, 1986.
- S. Tavaré, D. J. Balding, R. C. Griffiths, and P. Donnelly. Inferring coalescence times from DNA sequence data. *Genetics*, 145(2):505–518, 1997.

- V. Volterra. *Variazioni e fluttuazioni del numero d'individui in specie animali conviventi*. C. Ferrari, 1927.
- X. Wang and D. B. Dunson. Parallelizing MCMC via Weierstrass sampler. *arXiv preprint arXiv:1312.4605*, 2013.
- D. J. Wilkinson. *Stochastic Modelling for Systems Biology*. CRC press, 2011.
- M. Xu, B. Lakshminarayanan, Y. W. Teh, J. Zhu, and B. Zhang. Distributed Bayesian posterior sampling via moment sharing. In *Advances in Neural Information Processing Systems*, pages 3356–3364, 2014.

A Proofs

A.1 Detailed balance condition for ABC exchange moves

Proof Let $p(\theta)$ be the prior distribution, recall that we define $0 \leq f(\theta) \leq 1$ and $0 \leq f'(\theta) \leq 1$ such that

$$f(\theta) = \int \mathbb{1}_\varepsilon(x) f(\mathrm{d}x | \theta) \quad f'(\theta) = \int \mathbb{1}_{\varepsilon'}(x) f(\mathrm{d}x | \theta)$$

where $\mathbb{1}_\varepsilon(x)$ is the indicator function for hitting the ball $B_\varepsilon(y)$ of radius ε and we have $\varepsilon' > \varepsilon$, i.e. f' corresponds to the probability of hitting a larger ball. Now let

$$\pi(\theta) = \frac{p(\theta)f(\theta)}{c} \quad \pi'(\theta) = \frac{p(\theta)f'(\theta)}{c'}$$

where c and c' are normalising constants.

We have $\theta \sim \pi$ and $\theta' \sim \pi'$ and define φ and φ' to be such that $(\theta, \theta') \rightarrow (\varphi, \varphi')$ via the transition kernel. Then, given the set probability of accepting a swap, we have

$$\begin{cases} \mathbb{P}(\text{Accept}) = \mathbb{P}[(\varphi, \varphi') = (\theta', \theta)] = \frac{f(\theta')}{f'(\theta')} \\ \mathbb{P}(\text{Reject}) = \mathbb{P}[(\varphi, \varphi') = (\theta, \theta')] = 1 - \frac{f(\theta')}{f'(\theta')} \end{cases}$$

And we want to show that for some set A ,

$$\mathbb{P}(\varphi \in A) = \mathbb{P}(\theta \in A) \quad \mathbb{P}(\varphi' \in A) = \mathbb{P}(\theta' \in A)$$

Now, we can decompose

$$\mathbb{P}(\varphi \in A) = \mathbb{P}(\theta \in A, \text{Reject}) + \mathbb{P}(\theta' \in A, \text{Accept}) \quad (11)$$

Taking the second summand, we have

$$\begin{aligned} & \mathbb{P}(\theta' \in A, \text{Accept}) \\ &= \mathbb{E}_{\theta \sim \pi, \theta' \sim \pi'} \left[\mathbb{1}_A(\theta') \frac{f(\theta')}{f'(\theta')} \right] \\ &= \int \frac{p(\theta)f(\theta)}{c} \frac{p'(\theta')f'(\theta')}{c'} \mathbb{1}_A(\theta') \frac{f(\theta')}{f'(\theta')} \cdot \frac{f'(\theta)}{f(\theta)} \mathrm{d}\theta \mathrm{d}\theta' \\ &= \int \underbrace{\frac{p(\theta)f'(\theta)}{c'}}_{\pi'(\theta)} \underbrace{\frac{p(\theta')f(\theta')}{c}}_{\pi(\theta')} \mathbb{1}_A(\theta') \frac{f(\theta)}{f'(\theta)} \mathrm{d}\theta \mathrm{d}\theta' \\ &= \mathbb{E}_{\theta \sim \pi', \theta' \sim \pi} \left[\mathbb{1}_A(\theta') \frac{f(\theta)}{f'(\theta')} \right] \end{aligned}$$

Now relabelling $\theta := \theta'$ and $\theta' := \theta$ we obtain

$$\begin{aligned}\mathbb{P}(\theta' \in A, \text{Accept}) &= \mathbb{E}_{\theta' \sim \pi', \theta \sim \pi} \left[\mathbb{1}_A(\theta) \frac{f(\theta')}{f'(\theta)} \right] \\ &= \mathbb{P}(\theta \in A, \text{Accept})\end{aligned}$$

And hence Equation 11 becomes

$$\begin{aligned}\mathbb{P}(\varphi \in A) &= \mathbb{P}(\theta \in A, \text{Reject}) + \mathbb{P}(\theta \in A, \text{Accept}) \\ &= \mathbb{P}(\theta \in A)\end{aligned}$$

The proof that $\mathbb{P}(\varphi' \in A) = \mathbb{P}(\theta' \in A)$ is analogous.

A.2 Anytime distribution of the cold chain

Proof To obtain the anytime distribution in the Gamma mixture example in Section 5.1, compute the three components of the expression in Equation (3):

1. The density of X

$$\pi(dx) = \frac{x^{k_1-1}}{2\Gamma(k_1)\theta_1^{k_1}} e^{-\frac{x}{\theta_1}} + \frac{x^{k_2-1}}{2\Gamma(k_2)\theta_2^{k_2}} e^{-\frac{x}{\theta_2}} dx$$

where $\Gamma(\cdot)$ is the gamma function.

2. The expectation of $H \mid x$

$$\mathbb{E}[H \mid x] = \psi x^p + (1 - \psi)x^p = x^p$$

Note that the ψ factors cancel out, meaning that the anytime distribution is independent of ψ and therefore its value can be chosen to be 1 for convenience.

3. To compute $\mathbb{E}[H]$, use a property of conditional expectation and the honesty conditions of the $\text{Gamma}(k_1 + p, \theta_1)$ and $\text{Gamma}(k_2 + p, \theta_2)$ distributions:

$$\begin{aligned}\mathbb{E}[H] &= \mathbb{E}[\mathbb{E}(H \mid x)] = \mathbb{E}[x^p] \\ &= \int \frac{x^{p+k_1-1}}{2\Gamma(k_1)\theta_1^{k_1}} e^{-\frac{x}{\theta_1}} dx + \int \frac{x^{p+k_2-1}}{2\Gamma(k_2)\theta_2^{k_2}} e^{-\frac{x}{\theta_2}} dx \\ &= \frac{\Gamma(p+k_1)\theta_1^{p+k_1}}{2\Gamma(k_1)\theta_1^{k_1}} \cdot 1 + \frac{\Gamma(p+k_2)\theta_2^{p+k_2}}{2\Gamma(k_2)\theta_2^{k_2}} \cdot 1 \\ &= \frac{\Gamma(k_2)\Gamma(p+k_1)\theta_1^p + \Gamma(k_1)\Gamma(p+k_2)\theta_2^p}{2\Gamma(k_1)\Gamma(k_2)} \\ &= \frac{C}{2\Gamma(k_1)\Gamma(k_2)}\end{aligned}$$

letting $C = \Gamma(k_2)\Gamma(p+k_1)\theta_1^p + \Gamma(k_1)\Gamma(p+k_2)\theta_2^p$.

Combining the three components,

$$\begin{aligned}\alpha(dx) &= \frac{2\Gamma(k_1)\Gamma(k_2)}{C} \left(\frac{x^{p+k_1-1}}{2\Gamma(k_1)\theta_1^{k_1}} e^{-\frac{x}{\theta_1}} + \frac{x^{p+k_2-1}}{2\Gamma(k_2)\theta_2^{k_2}} e^{-\frac{x}{\theta_2}} \right) dx \\ &= \underbrace{\frac{\Gamma(k_2)\Gamma(p+k_1)\theta_1^{p+k_1}}{C\theta_1^{k_1}}}_{\varphi(p)} \underbrace{\frac{x^{p+k_1-1}}{\Gamma(p+k_1)\theta_1^{p+k_1}} e^{-\frac{x}{\theta_1}}}_{\text{Gamma}(p+k_1, \theta_1)} \\ &\quad + \underbrace{\frac{\Gamma(k_1)\Gamma(p+k_2)\theta_2^{p+k_2}}{C\theta_2^{k_2}}}_{\varphi'(p)} \underbrace{\frac{x^{p+k_2-1}}{\Gamma(p+k_2)\theta_2^{p+k_2}} e^{-\frac{x}{\theta_2}}}_{\text{Gamma}(p+k_2, \theta_2)} dx\end{aligned}$$

And now substituting back the expression C in $\varphi(p)$:

$$\begin{aligned}\varphi(p) &= \frac{\Gamma(k_2)\Gamma(p+k_1)\theta_1^p}{\Gamma(k_2)\Gamma(p+k_1)\theta_1^p + \Gamma(k_1)\Gamma(p+k_2)\theta_2^p} \\ &= \frac{1}{1 + \frac{\Gamma(k_1)\Gamma(p+k_2)\theta_2^p}{\Gamma(k_2)\Gamma(p+k_1)\theta_1^p}}\end{aligned}$$

Similarly, we can obtain $\varphi'(p) = 1 - \varphi(p)$. Therefore, the anytime distribution $\alpha(dx)$ is the following mixture of two Gamma distributions:

$$\begin{aligned}\alpha(dx) &= \varphi(p) \text{Gamma}(k_1 + p, \theta_1) \\ &\quad + (1 - \varphi(p)) \text{Gamma}(k_2 + p, \theta_2)\end{aligned}$$

B Tables

Chain k	ε^k	σ^k	ABC-PTMC-1	ABC-PTMC-W	ABC-APTMC-1	ABC-APTMC-W
1	1	0.008	2667.5	3241.8	6570.3	14488
2	1.046	0.009	2790	3564.8	6934.8	16902
3	1.094	0.011	2796.8	3567.5	6941	16837
4	1.145	0.012	2797.3	3604.5	6924.8	17264
5	1.197	0.014	2793.8	3719.8	6931.3	15710
6	1.253	0.016	2786.3	3748.8	6951.5	17157
7	1.31	0.019	2784.8	3629.3	6961.3	18947
8	1.371	0.022	2795.5	3615	6941.5	18551
9	1.434	0.025	2805.5	3608.5	6950.8	18759
10	1.5	0.029	2803.5	3711.5	6962.8	17276
11	1.661	0.034	2798.5	3799.8	6962.3	46350
12	1.84	0.039	2803.3	3693.3	6983	53716
13	2.038	0.045	2814.8	3656.5	6995.3	53289
14	2.257	0.052	2799	3658.3	7008.8	53458
15	2.5	0.06	2783.5	3796.3	7029.8	46597
16	3.362	0.092	2787.5	4054.5	7038.5	68953
17	4.522	0.14	2783.8	3936.5	7009.5	79231
18	6.082	0.214	2781	3912.5	6982.5	78917
19	8.179	0.327	2780.8	3919.5	7002.8	79038
20	11	0.5	2665.8	3598.5	6604.3	67725

Table 7 Average sample sizes per chain returned over four 24-hour runs of the ABC-PTMC-1, ABC-APTMC-1, ABC-PTMC-W, ABC-APTMC-W algorithms to estimate the posterior distributions of the parameters θ of a stochastic Lotka-Volterra model on multiple processors in Section 5.4.2. The ball radius ε^k and proposal distribution covariance $\text{diag}(\sigma^k, \sigma^k 10^{-2}, \sigma^k)$ associated with each chain k are displayed for information.

Chain k	ε^k	σ^k	ABC-PTMC-1	ABC-PTMC-W	ABC-APTMC-1	ABC-APTMC-W
1	1	0.008	4.0364	10.254	6.3818	17.032
2	1.046	0.009	8.1412	18.47	11.344	28.873
3	1.094	0.011	8.7694	18.642	11.413	28.619
4	1.145	0.012	8.8518	19.329	11.189	30.387
5	1.197	0.014	8.6446	21.793	11.27	23.531
6	1.253	0.016	8.3202	21.528	11.548	22.644
7	1.31	0.019	8.3717	18.879	11.677	29.882
8	1.371	0.022	8.8443	18.581	11.423	28.332
9	1.434	0.025	8.994	18.568	11.546	29.134
10	1.5	0.029	8.9767	20.839	11.699	23.191
11	1.661	0.034	8.8116	21.189	11.705	19.591
12	1.84	0.039	8.9863	18.829	11.968	30.582
13	2.038	0.045	9.3053	17.93	12.108	30.013
14	2.257	0.052	9.0069	18.051	12.298	30.245
15	2.5	0.06	8.628	20.982	12.581	20.011
16	3.362	0.092	8.369	19.888	12.668	18.274
17	4.522	0.14	8.1945	17.477	12.282	28.847
18	6.082	0.214	8.181	16.817	11.965	28.564
19	8.179	0.327	8.1896	17.075	12.246	28.677
20	11	0.5	4.2839	9.7452	6.8761	16.773

Table 8 Mean percentage of exchange moves per chain returned over four 24-hour runs of the ABC-PTMC-1, ABC-APTMC-1, ABC-PTMC-W, ABC-APTMC-W algorithms to estimate the posterior distributions of the parameters θ of a stochastic Lotka-Volterra model on multiple processors in Section 5.4.2. The ball radius ε^k and proposal distribution covariance $\text{diag}(\sigma^k, \sigma^k 10^{-2}, \sigma^k)$ associated with each chain k are displayed for information.

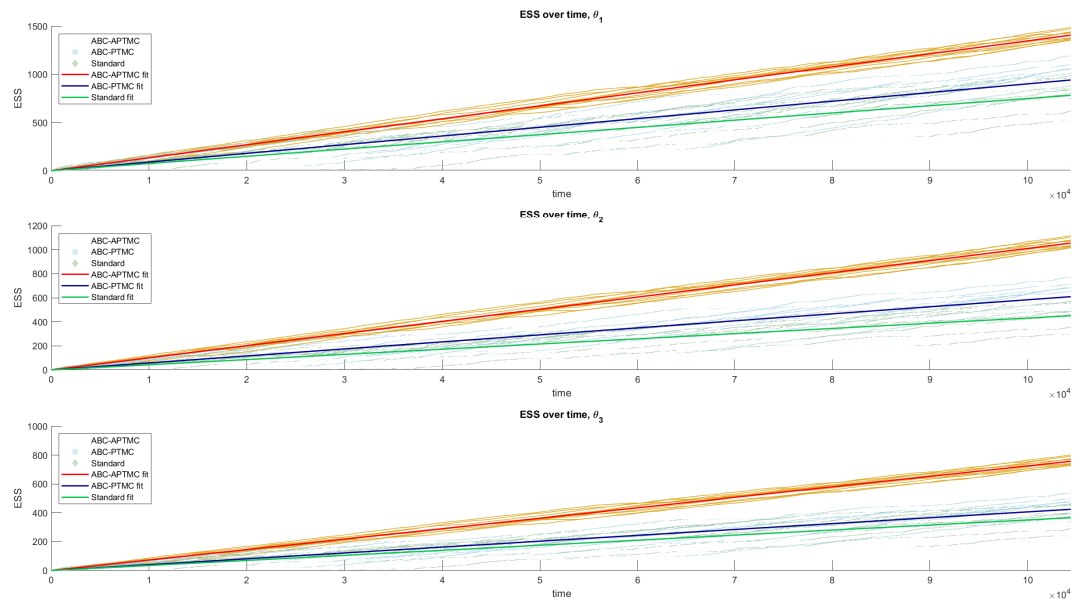


Fig. 13 Effective sample size over time for ten 28-hour runs of the standard ABC, ABC-PTMC-1 and ABC-APTMC-1 algorithms.