

# Parallel Markov Chain Monte Carlo Simulation by Pre-Fetching

A.E. Brockwell

March 15, 2005

## Abstract

In recent years, parallel processing has become widely available to researchers. It can be applied in an obvious way in the context of Monte Carlo simulation, but techniques for “parallelizing” Markov chain Monte Carlo (MCMC) algorithms are not so obvious, apart from the natural approach of generating multiple chains in parallel. While generation of parallel chains is generally the easiest approach, in cases where burn-in is a serious problem, it is often desirable to use parallelization to speed up generation of a single chain. This paper briefly discusses some existing methods for parallelization of MCMC algorithms, and proposes a new “pre-fetching” algorithm to parallelize generation of a single chain.

**Keywords:** parallel processing, Markov chain Monte Carlo, Bayesian, inference, pre-fetching

## 1 Introduction

Over the last two decades, the increased availability of cheap computing power has dramatically altered the way statistical analyses are carried out. Many problems previously considered intractable can now be solved by intensive numerical methods. In addition, in recent years, networking has also become cheap. The vast majority of computers sold now come with built-in ethernet adaptors. Partly as a result of this, parallel computing has received new impetus, since it is possible to put together a group of networked computers for on the order of a thousand dollars per machine. Such networks, typically consisting of machines without keyboards or displays, are now widespread, and commonly referred to as “Beowulf clusters”. Furthermore, well-defined standards for communication between processors have been developed. The “Message Passing Interface” (MPI, see, e.g. Gropp

et al., 1999) is now widely accepted, and implementations in the form of C/C++/Fortran libraries are publically available. Thus parallel processing can be exploited by dividing a task into a number of sub-tasks which can be executed in parallel. Each sub-task is executed on a separate processor, and then the results are compiled, typically by a main “controlling” processor. Such a scheme can be implemented by writing code in C, C++, or Fortran, making use of the MPI library to handle inter-process communication. Alternatively, for a simpler higher-level solution, one can use the recently-developed “Snow” package for R (see <http://cran.r-project.org> and the links to packages and documentation for more details).

In statistics, perhaps the most obvious application of parallel processing is in Monte Carlo simulation, where one estimates a function

$$h_\pi = \int h(x)d\pi(x), \tag{1}$$

for some probability distribution  $\pi$ . The Monte Carlo approximation is simply

$$\hat{h} = \frac{1}{t} \sum_{j=1}^t h(X_j), \tag{2}$$

where  $\{X_j\}$  is a sequence of independent draws from the distribution  $\pi$ . Such a problem is trivial to “parallelize”. The basic principle is to subdivide the sum into  $P \geq 2$  components, and assign one processor to evaluate each component. The results for each component can then be added and normalized, either by the user, or by a “controlling” processor, to obtain the final result. When inter-processor communication time is negligible compared to the time taken to execute each sub-task, and processors are roughly the same speed (i.e., they are “balanced”), such an approach leads to an increase in speed by a factor approximately equal to  $P$ .

Markov chain Monte Carlo (MCMC) methods (see Gilks et al., 1996) are a variant of Monte Carlo schemes in which a Markov chain  $\{X_j, j = 1, 2, \dots\}$  with limiting distribution  $\pi$  is generated. It can be used for estimating posterior distributions in a wide class of models, even when the likelihood includes an unknown normalizing constant. Estimation can still be carried out using (2), but in this case, elements of the sequence  $\{X_j\}$  are not independent of each other. Furthermore, the initial value is typically not a draw from the distribution  $\pi$ . However, if the chain is constructed properly, then  $X_t \xrightarrow{d} \pi$ , and under certain conditions, the estimator  $\hat{h}$  converges to  $h_\pi$  as  $t \rightarrow \infty$ . Unfortunately, generation of a Markov chain is not well-suited to be carried out by parallel processing. The process is fundamentally sequential in nature; the distribution of  $X_{j+1}$  depends on the value of  $X_j$ , so simulation for one step is not seemingly possible until the results for the previous step have been obtained. On the other hand, it is often highly desirable to speed up MCMC simulation, particularly when convergence to the limiting distribution is slow.

Given the difficulties arising with parallelization because of the sequential nature of MCMC simulation, a natural thing to do, discussed in, for instance, Glynn and Heidelberger (1992)

and Rosenthal (2000), is simply to generate a separate Markov chain on each processor and combine the results appropriately (see, e.g. Bradford and Thomas, 1996, for an example.) This has the advantage of requiring very little effort beyond that required to program a single-processor version of the MCMC generation code. However, the drawback is that error associated with burn-in still remains in all processes. Glynn and Heidelberger (1992) and Rosenthal (2000) consider this issue in some detail, and discuss various methods for deciding how much of the initial portion of the chain to remove. (Note also that for low-dimensional problems, it may be preferable to use numerical integration methods, often referred to as “numerical quadrature” or “numerical cubature”, instead of MCMC simulation. Such methods are typically trivially parallelizable, but they tend not to perform well in high-dimensional problems.)

In certain problems, the time spent in the burn-in phase may be significant (for instance, if likelihood calculations are long, or if the convergence rate of the chain is very slow). In such cases it is often desirable to speed up generation of a single chain, rather than use multiple chains. When the state-space of the chain is high-dimensional, one possible way to do this is to divide the state-space into blocks, and then for each iteration of the Markov chain, to handle each block on a separate processor. (This is discussed, for instance, in the nice chapter-length introduction to parallel computing for Bayesian analysis given by Wilkinson (2004), and an interesting example of this kind of approach can be found in Whitley and Wilson (2004).) This approach does indeed speed up generation of a single chain, but requires additional effort, in carrying out analysis of the limiting distribution, in order to come up with appropriate blocks. This can be difficult, particularly when the conditional dependence structure in the limiting distribution is complicated. In fact, in certain cases (such as the case study given in this paper), it may be impossible. We therefore propose a new algorithm for the purpose of speeding up generation of a single chain by parallelization. The idea is to calculate multiple likelihoods ahead of time (“pre-fetching”), and only use the ones which are needed. The approach does not require any particular analysis of the limiting distribution of the chain (for instance, to sub-divide the state-space into blocks). For convenience, we also include in this paper brief discussion of two other methods of doing this, the blocking approach mentioned above, as well as an approach based on regenerative simulation. We demonstrate the potential gains in performance by applying the pre-fetching method in a time series problem.

## 2 General Considerations in Parallel Processing

Conceptually, parallel processing can be applied to almost any problem (*task*) by sub-dividing it into multiple *sub-tasks*. The execution of sub-tasks may or may not be dependent on the results of other sub-tasks. When a group of two or more sub-tasks needs to be executed, and none of the sub-tasks depends on the results of any of the others, then the sub-tasks can be executed concurrently by different processors. Assuming that no single processor is

particularly slow relative to the others, this clearly leads to an increase in speed of execution of the original task.

Before considering the specifics of parallel processing for Markov chain Monte Carlo simulation, it is important to keep in mind (at least) the following three factors in this standard approach to parallel processing.

The first is “granularity”, that is, the size of the sub-tasks. In particular, the ratio of time required to complete a sub-task to the time required to carry out inter-process communication is critically important. Typically, machines in a Beowulf cluster communicate via network connections, meaning that even simple messages between processors take on the order of milliseconds to send. On the other hand, the CPUs of the machines can typically execute, for instance, a multiplication operation in on the order of picoseconds. Therefore, if the gain in speed due to parallelization is not to be lost because of time spent communicating, it is critical that (relative to communication times), the sub-tasks each require substantial computational times.

Secondly, for statisticians using parallel processors to analyze Monte Carlo or Markov chain Monte Carlo problems, it is important to ensure that random number streams on separate processors are independent of each other. Otherwise, one can lose the benefit from parallelizing the Monte Carlo approximation - in the extreme case, if all processors generate exactly the same random number stream, then the result in the parallel approach can be exactly the same as that obtained using a single processor. A simple way of avoiding this problem is to ensure that the random number seed is set differently on each processor, although technically, it is still possible to obtain sequences on separate processors with some degree of overlap. A more sophisticated approach is to use the “Scalable Parallel Random Number Generator” (SPRNG, Mascagni and Srinivasan, 2000, see also <http://sprng.cs.fsu.edu>), which is specifically designed to generate independent streams of (pseudo-)random numbers.

Finally, individual processor speeds may be important. Many Beowulf clusters consist of machines which all run at roughly the same speed. However, in situations where some parts of the cluster are newer than others, it is not uncommon to see a range of different-speed processors. Furthermore, if the cluster is available to many users, then at any given time, a machine may be effectively slowed down if one of these users is running a computationally-intensive program. As a consequence, a single slow machine may hold up completion of the task. The solution is to employ “load-balancing” to make sure that each processor receives an amount of work proportional to its speed, and two common approaches to use here are

1. to ensure that slower processors receive smaller sub-tasks than faster processors, or
2. to adopt a queueing approach.

To use the queueing approach, one simply divides the task into a large number  $N$  of small sub-tasks. “Large” here means any number significantly greater than the number of processors

$P$ . Then the first  $P$  sub-tasks can be allocated, one to each processor. The remaining pool of  $N - P$  sub-tasks is placed in a queue. As soon as the first processor completes its sub-task, it is assigned the next sub-task, which is removed from the front of the queue. Remaining elements of the queue are subsequently assigned to new processors as the processors become available. This approach ensures that faster processors receive more sub-tasks, and provides a simple way to carry out automatic load-balancing. However, it is only effective when sub-tasks can be carried out independently and communication time requirements remain small compared to computation time required for the sub-tasks.

### 3 Parallel Generation of a Single Chain

The Metropolis-Hastings algorithm and its many variants give us straightforward schemes for obtaining (non-independent) draws  $\{\theta^{(1)}, \theta^{(2)}, \dots\}$  from a target distribution  $\pi$  defined on a state-space  $\Theta$  (Typically  $\Theta$  would be  $m$ -dimensional Euclidean space  $\mathbb{R}^m$ , and the target distribution would be the posterior distribution of a set of  $m$  real-valued parameters.) Our goal is to use parallel processors to speed up simulation of a single Metropolis-Hastings chain with limiting distribution  $\pi$ .

As Wilkinson (2004) points out, it may also be possible in certain cases to parallelize burdensome likelihood computations directly. For instance, if they involve high-dimensional matrix operations, then the ScaLAPACK library (Choi et al., 1992, also see <http://www.netlib.org/scalapack>) may be used to speed up computations. (In fact, this technique is used by Whiley and Wilson, 2004, .) Furthermore, in many cases, it is possible to speed up computation of the target density without resorting to use of parallel processors. Standard techniques for improving program efficiency include working to eliminate un-necessary network/disk accesses and redundant computations, avoiding the use of transcendental functions when possible, using efficient libraries for optimization, random number generation, etc.

When these approaches still do not provide sufficient improvement in speed, one may still be able to resort to one of the following schemes.

#### 3.1 Regeneration

For problems with low-dimensional state-space (i.e. few parameters and latent variables), regeneration (see Mykland et al., 1995; Brockwell and Kadane, 2004, for details) can be used. In the discrete state-space case, use of regeneration to parallelize generation of a single chain is conceptually straightforward. A single point  $\theta^*$  in the state-space is chosen, and the chain is started from this point. The resulting Markov chain can be divided into segments, each one beginning at  $\theta^*$  and terminating immediately before the next occurrence

of  $\theta^*$ . These segments, typically referred to as “tours”, are independent and identically distributed. Therefore each tour can be generated on a separate processor. The tours can then be patched together in a prespecified order to obtain a single long chain. In the (more common) continuous state-space case, this approach must be modified slightly, since the probability of a return to the original state  $\theta^*$  is in most cases equal to zero. However, for low-dimensional state-spaces it is not difficult to adapt the approach. Explicit algorithms are given in Brockwell and Kadane (2004). The primary limitations with this approach, however, are its lack of practical applicability to high-dimensional problems, and the introduction of an additional “re-entry” distribution, which if poorly chosen, can inhibit mixing of the chain.

## 3.2 Blocking

In MCMC problems for which the state-space is high-dimensional, it is tempting to use an update scheme where within each iteration, each processor is responsible for updating a part of the state-space. However, unless done carefully, this is *not* a valid scheme. The counterexample given in Appendix A illustrates that it does not necessarily yield a Markov chain with the correct invariant distribution. In spite of the fact that this obvious approach doesn’t yield a chain with the correct limiting distribution, under a conditional independence condition, it is still possible in many cases to make effective use of parallel processing. Suppose that the the following property holds.

**Property 3.1** *For some  $B \geq 2$ , there exists a decomposition of the state-space*

$$\Theta = \prod_{i=0}^B \Theta_i$$

*such that, for all  $\theta = (\theta_0, \dots, \theta_B) \in \Theta$ , the following equivalent conditions hold.*

1.  $\pi(\theta_1, \dots, \theta_B | \theta_0) = \pi(\theta_1 | \theta_0) \pi(\theta_2 | \theta_0) \dots \pi(\theta_B | \theta_0)$ .
2. For  $j = 1, 2, \dots, B$ ,  $\pi(\theta_j | \theta_{-j}) = \pi(\theta_j | \theta_0)$ .

This is a form of the Markov property, in which the conditional independence structure is not necessarily determined by time or spatial location. An example of a model for which such blocks can be found is the so-called generalized state-space model, discussed, along with an example of a possible block decomposition are in Appendix B. When Property 3.1 holds for some state-space  $\Theta$  and target distribution  $\pi$ , the following parallel Markov chain Monte Carlo algorithm can be used.

### Algorithm 3.1: Parallel Block-Metropolis-Hastings

**Step 1.** Choose an initial state  $\theta^{(1)} = (\theta_0^{(1)}, \dots, \theta_B^{(1)})$ . Set  $k = 1$ .

**Step 2.** Create a copy  $\theta^{(k+1)}$  of  $\theta^{(k)}$ . Set  $k \leftarrow k + 1$ .

**Step 3.** Carry out a Metropolis-Hastings update of  $\theta_0^{(k)}$  (given  $\theta_1^{(k)}, \dots, \theta_B^{(k)}$ ).

**Step 4.** On  $B$  separate processors, carry out concurrent Metropolis-Hastings updates of  $\theta_j^{(k)}$ ,  $j = 1, 2, \dots, B$ . The proposal distribution for updating  $\theta_j^{(k)}$  may depend on  $\theta_0^{(k)}$  as well as the current value of  $\theta_j^{(k)}$  but not on  $\{\theta_m^{(k)}, m \neq 0, m \neq j\}$ .

**Step 5.** Go back to Step 2.

Note that Wilkinson (2004) discusses a more general form of this blocking algorithm, which in a number of cases enables more efficient parallel block-update schemes to be implemented (in the sense that one obtains an increase in speed closer to the number of processors  $P$ ). The scheme here can be regarded as a special case of that described in Wilkinson (2004), where there are two blocks,  $T_1 = \theta_0$ , and  $T_2 = (\theta_1, \dots, \theta_B)$ .

Note also that in many cases, it is convenient to further subdivide the update for  $\theta_j^{(k)}$  into a number of Metropolis-Hastings updates of sub-components of  $\theta_j^{(k)}$ . (These, however, would all have to be carried out sequentially on the same processor.)

*Remark:* A typical update for  $\theta_j^{(k)}$  in Step 4 would involve drawing a proposal  $\theta_j^*$  from a distribution  $g_j(\cdot; \theta_j^{(k)}, \theta_0^{(k)})$ , and accepting it with probability

$$\alpha = \min \left( 1, \frac{\pi(\theta_0^{(k)}, \theta_1^{(k)}, \dots, \theta_j^*, \dots, \theta_B^{(k)})g(\theta_j^{(k)}; \theta_j^*, \theta_0^{(k)})}{\pi(\theta_0^{(k)}, \theta_1^{(k)}, \dots, \theta_j^{(k)}, \dots, \theta_B^{(k)})g(\theta_j^*; \theta_j^{(k)}, \theta_0^{(k)})} \right).$$

Since Property 3.1 is required to hold, this acceptance probability simplifies to

$$\alpha = \min \left( 1, \frac{\pi(\theta_j^* | \theta_0^{(k)})g(\theta_j^{(k)}; \theta_j^*, \theta_0^{(k)})}{\pi(\theta_j^{(k)} | \theta_0^{(k)})g(\theta_j^*; \theta_j^{(k)}, \theta_0^{(k)})} \right).$$

### 3.3 Pre-fetching

We propose a new method, which we call “pre-fetching”, for parallel generation of a Markov chains when burn-in time is significant, and the methods discussed above are not practical. (For certain problems, it may be difficult to establish Property 3.1 or the more general

property discussed in Wilkinson (2004), and even when it can be established, the particular block structure may not lead to substantial improvement in speed. In some cases, it may even be impossible to find any conditional independence structure at all. This occurs, for instance, in long-memory time series models like the ones considered in Brockwell and Chan (2004). Such an example is given in the next section.)

Suppose that a Metropolis-Hastings chain has already been developed, and that proposal generation is virtually instantaneous, as well as evaluation of prior densities, but that the main computational burden is in computing likelihoods. (This is most often the case, since proposal distributions are usually chosen to have simple density functions and be easy to sample from.)

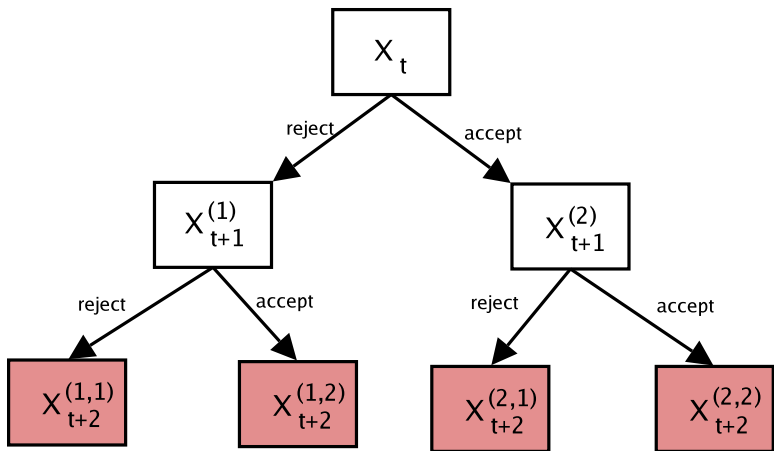


Figure 1: The possible outcomes in two iterations of a Metropolis-Hastings sampler. In two steps, there are only four unique values, those in the shaded leaf nodes, for which likelihoods must be computed. Parents have the same value as their left (“reject”) children.

For the sake of clarity, consider the possible outcomes in only two sequential iterations of the sampler, illustrated in Figure 1. Starting at time  $t$ , the chain has state  $X_t$ . At time  $t+1$ , the state  $X_{t+1}$  is either the same as  $X_t$  (if the proposal is rejected), or is equal to the proposal generated at time  $t$ . At time  $t+2$ , there are four possible outcomes, depending on both the acceptance/rejection in going from time  $t$  to  $t+1$  and that in going from time  $t+1$  to time  $t+2$ . In this context, the pre-fetching algorithm assigns one processor to evaluate the likelihood at each one of the four leaf nodes. These likelihoods determine the likelihoods in the parent nodes (a parent node has the same likelihood as its left “rejection” child). Then a single processor gathers the results, and takes two steps down the tree, drawing a uniform random variable each time and comparing it with the appropriate acceptance probability.

Formally, and in the more general case, let

$$X_{t+h}^{(I_1, I_2, \dots, I_h)}$$



denote the value of the state  $X_{t+h}$ , in the case where the proposals in going from time  $(t + j - 1)$  to  $(t + j)$  are

$$\begin{aligned} &\text{rejected, if } I_j = 1, \\ &\text{accepted, if } I_j = 2. \end{aligned}$$

(Figure 1 shows such values for  $h = 1$  and  $h = 2$ .) Also, for  $j = 0, 1, 2, \dots, h - 1$ , let

$$Z_{t+j}^{(I_1, \dots, I_j)}$$

denote the proposal generated at time  $t + j$  (for the value of the state at time  $t + j + 1$ ), contingent upon obtaining the corresponding sequence of acceptances/rejections starting at time  $t$ . Thus

$$X_{t+j}^{(I_1, \dots, I_j)} = \begin{cases} X_{t+j-1}^{(I_1, \dots, I_{j-1})}, & I_j = 1 \\ Z_{t+j-1}^{(I_1, \dots, I_{j-1})}, & I_j = 2. \end{cases} \quad (3)$$

Suppose also that the possibly time-varying proposal densities are specified by

$$g_t(x_t, z_t) = P(Z_t \in dz_t | X_t = x_t).$$

The pre-fetching algorithm for carrying out  $h$  iterations of MCMC simulation, starting with  $X_t$  and ending with  $X_{t+h}$ , requires  $2^h$  processors. (For good performance, these processors should be balanced, that is, they should be approximately the same speed.)

**Algorithm 3.2: Pre-Fetching**

**Step 1.** Compute all possible proposals and states for the  $h$  steps into the future, as follows.

1. Draw the proposal  $Z_t$  from the density  $g(x_t, z_t)$ . Determine  $X_{t+1}^{(1)}$  and  $X_{t+1}^{(2)}$  using (3).
2. Draw proposals  $Z_{t+1}^{(1)} \sim g_{t+1}(x_{t+1}^{(1)}, z_{t+1})$  and  $Z_{t+1}^{(2)} \sim g_{t+1}(x_{t+1}^{(2)}, z_{t+1})$  and then determine  $X_{t+2}^{I_1, I_2}$  for all four possible outcomes  $(I_1, I_2) \in \{1, 2\}^2$ .
3. Repeat this procedure to determine all possible proposals  $Z_{t+j}^{(I_1, \dots, I_j)}$ ,  $(I_1, \dots, I_j) \in \{1, 2\}^j$ ,  $j = 1, 2, \dots, h$ .

**Step 2.** Identify the  $2^h$  unique possible values that the states  $x_t, \dots, x_{t+h}$  can take. These are simply the values  $\{x_t \cup \{x_{t+j}^{(I_1, \dots, I_j)} : I_j = 2\}, j = 1, \dots, h\}$ . Label these  $x_i^*$ ,  $i = 1, 2, \dots, 2^h$ .

**Step 3.** Concurrently, on  $2^h$  separate processors, compute the target density  $\pi(x_i^*)$ ,  $i = 1, 2, 3, \dots, 2^h$ .

**Step 4.** For  $j = 1, 2, \dots, h$ , compute the realizations  $i_j$  of  $I_j$  using the standard Metropolis-Hastings rule

$$i_j = \begin{cases} 2, & \text{with probability } \alpha_j, \\ 1, & \text{otherwise,} \end{cases}$$

where

$$\alpha_j = \frac{\pi(z_{t+j}^{(i_1, \dots, i_{j-1})})g_t(z_{t+j}^{(i_1, \dots, i_{j-1})}, x_{t+j}^{(i_1, \dots, i_{j-1})})}{\pi(x_{t+j}^{(i_1, \dots, i_{j-1})})g_t(x_{t+j}^{(i_1, \dots, i_{j-1})}, z_{t+j}^{(i_1, \dots, i_{j-1})})}$$

**Step 5.** Set  $X_{t+j} = x_{t+j}^{(i_1, \dots, i_j)}$ , for  $j = 1, 2, \dots, h$ .

Since all computations apart from evaluation of  $\pi(\cdot)$  are assumed to be negligible, the speed-limiting component of this algorithm is Step 3. Since likelihoods are computed on separate processors, this algorithm generates  $h$  iterations of the Markov chain in the time taken to evaluate the density  $\pi$  only once.

This algorithm is not particularly efficient, since it achieves a speed-increase of only  $\log_2(P)$ , where  $P$  is the number of processors. Essentially, this is because much of the computation is wasted - only one out of  $2^h$  possible paths is actually chosen. However, in cases where none of the previously-mentioned methods are practical to implement, it provides a useful alternative. Furthermore, the algorithm is relatively straightforward to implement.

## 4 Simulation Study

Fractionally integrated autoregressive moving average models have received attention recently as potential improvements over standard autoregressive moving average models in a number of fields, including finance, meteorology, and finance, where time series appear to exhibit long-range dependence structure.

A zero-mean ARFIMA( $p, d, q$ ) process  $\{Y_t\}$  satisfies

$$\phi(B)(1 - B)^d Y_t = \vartheta(B)Z_t, \quad (4)$$

where  $B$  denotes the backshift operator,  $\phi(B) = 1 - \sum_{i=1}^p \phi_i B^i$ ,  $\vartheta(B) = 1 + \sum_{i=1}^q \vartheta_i B^i$ , the “fractional differencing parameter”  $d$  is a constant in the range  $(-0.5, 0.5)$ ,  $\{Z_t\}$  is an iid Gaussian noise sequence with mean zero and variance  $\sigma^2$ , and the roots of the polynomials  $\phi(\cdot)$  and  $\vartheta(\cdot)$  lie strictly outside the unit circle. The fractional differencing operator  $(1 - B)^d$  is interpreted in the usual manner (see, e.g. Beran, 1994) as

$$(1 - B)^d = \sum_{k=0}^{\infty} \frac{\Gamma(d + 1)}{\Gamma(k + 1)\Gamma(d - k + 1)} (-1)^k B^k.$$

To facilitate likelihood computations, the process is usually assumed to be Gaussian. When observations  $Y_1 = y_1, Y_2 = y_2, \dots, Y_N = y_n$  are made, the likelihood function is simply that of a multivariate normal,

$$l(y_1, \dots, y_N; \theta) = (2\pi)^{-N/2} \det(\Gamma)^{-1} \exp(-y^T \Gamma^{-1} y / 2),$$

where  $y = (y_1, \dots, y_N)$ , and  $\Gamma = [\gamma_{ij}]_{i,j=1,\dots,N}$  denotes the covariance matrix defined by  $\gamma_{ij} = \text{Cov}(X_i, X_j)$ . The components of  $\Gamma$  are computationally tedious to obtain, although the formula given by Sowell (1992) is very helpful. Furthermore, due to the long-memory property of  $\{Y_t\}$ , the standard (efficient) approach to likelihood evaluation for ARMA models (see, e.g. Brockwell and Davis, 1991) cannot be used. Instead, the best existing method is to carry out a Cholesky decomposition of  $\Gamma$ , or to use the Durbin-Levinson algorithm to evaluate  $l(y_1, \dots, y_N; \theta)$ . For large values of  $N$  (and indeed, one is often interested in the case where  $N$  is large when carrying out analysis of long-memory time series), this is time-consuming to compute.

To investigate the performance of the pre-fetching algorithm, we carried out a Bayesian analysis of a simulated ARFIMA(1,d,0) process with parameters  $\phi_1 = 0.5, d = 0.3, \sigma^2 = 1$ . The simulated process is shown in Figure 2. We imposed uninformative (very high variance) priors on parameter  $\phi_1, d$  and  $\log(\sigma)$ . Each step in the Metropolis-Hastings algorithm picked one of the three parameters at random, and generated a normally-distributed random-walk proposal. Variances of the random-walk step sizes were, respectively, 0.001, 0.0005, and 0.001.

## Simulated ARFIMA(1,0.3,0) Process

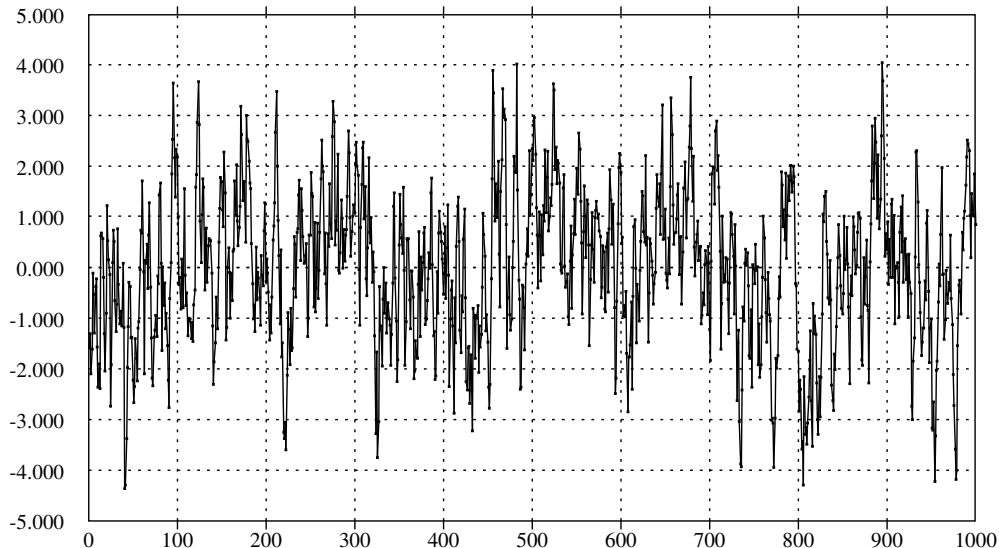


Figure 2: A simulated long-memory (ARFIMA(1,0.3,0)) process of length 1000.

The pre-fetching algorithm was implemented on a Beowulf cluster of 32 dual-CPU 1.6 GHz AMD Athlon Linux systems, connected to each other by 1 gigabit per second ethernet connections. At the time of running these tests, no other users were making use of the cluster. Programs were developed in C++, making use of the GNU scientific library, as well as the MPICH implementation (version 1.2.0) of the MPI library. In each run of the pre-fetching version of the MCMC algorithm, 10000 iterations of a Markov chain were generated, and in all cases, posterior means were indeed close to parameter values used in the simulation. Total execution time was recorded, and used to compute (Markov chain) iterations per second for the scheme. For each chosen number of processors, three runs were carried out.

On this particular cluster, average time to evaluate the likelihood was around 8 milliseconds. Observed iterations per second, for the Markov chains of length 10000, are shown in Figure 3. The solid line in the figure indicates the maximum speed we would expect to obtain, based on the assumption that processors are all running at the same speed, there are no “interruptions” on any processors, and that communication between processors is instantaneous, as well as the assumption that actual speed on each processor is the same as the average result obtained in the 1-processor runs.

Clearly, actual performance increases roughly as expected when using four processors, but starts to drop away from optimal theoretical performance as one goes to 8,16, and more processors. This is explained in part by the inherent sensitivity of the algorithm to variation in individual processing times. In particular, the time taken to carry out one  $h$  – step

update is the maximum of the times taken to evaluate the likelihoods over all  $2^h$  processors. Thus a small increase in variance of processing times can have a dramatic effect on overall performance, particularly as  $h$  grows.

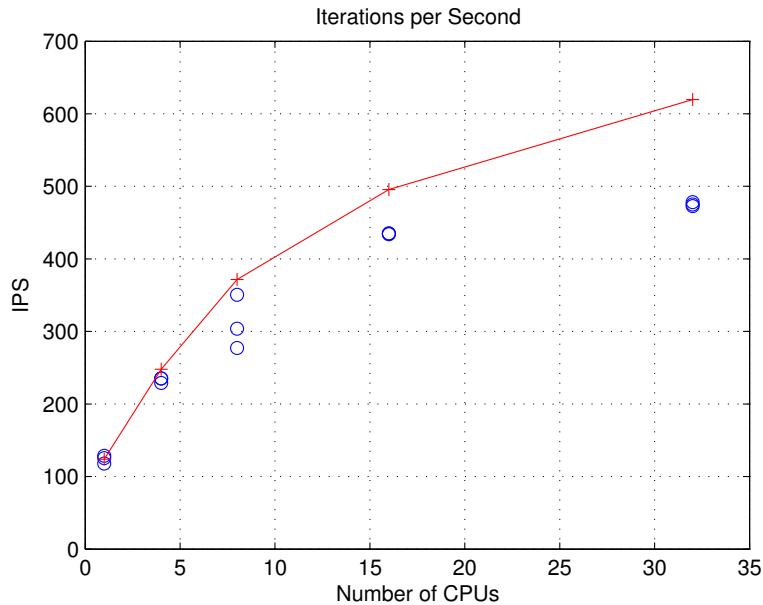


Figure 3: Iterations per second for the simulation study, as a function of number of processors. Observed rates are shown as circles for three runs within each choice of number of processors. The solid line indicates the theoretical optimal rate.

## 5 Discussion

This paper has introduced a new method for parallelizing generation of a single Markov chain in the context of MCMC simulation. As Wilkinson (2004) emphasizes, it is usually worth taking the time to improve mixing of a working single-processor MCMC algorithm, as well as optimizing program speed, before resorting to use of parallel processors. These approaches are typically easier than parallelizing generation of a chain. Furthermore, Rosenthal (2000), Glynn and Heidelberger (1992) and others have also considered the relatively simple approach of generating multiple chains in parallel, which is generally the most efficient solution when burn-in is not a serious problem. The pre-fetching method introduced here is best-suited for problems where burn-in time is significant, and other methods are not easy to implement.

The lack of robustness to processing-time variance highlighted in the simulation study suggests that, at least when going beyond the simplest 4-processor (2-times speedup) version of this algorithm, further refinements of the algorithm may be useful. By allocating spare

processors to the task, there is a range of possible ways to improve this robustness. One approach would be to allocate multiple processors to each likelihood evaluation, and simply take the first value returned. Another more complex approach would be carry out on-line evaluation of the variance of response times from processors, and allocate likelihood evaluations preferentially to the “reliable” (i.e. low-variance) machines. Furthermore, even more complicated versions of the algorithm could be developed that “travel down” the tree as soon as relevant results become available, and then cancel pending likelihood requests which as a consequence become redundant.

Another intriguing possibility, suggested to the author by an anonymous referee, arises in the case where one can guess whether or not acceptance probabilities will be “high” or “low”. In this case, the tree could be made deeper down “high” probability paths and shallower in the “low” probability paths. Theoretically, in such a case, one could exceed the log-base-two of number of processors speedup factor, since there would be a “high” probability of taking a deeper (than  $\log_2(P)$ ) path.

It is also worth noting the potential future benefits of an efficient perfect sampling algorithm, that is, an algorithm which yields a draw from exactly the distribution  $\pi$ . Since the seminal work of Propp and Wilson (1996) appeared, a lot of effort has gone into this area, but so far, no practical general-purpose scheme for typical applied Bayesian analysis problems has been developed. If it were possible to obtain these perfect samples, then the parallel chains approach would clearly be ideal. Each chain would be initialized with an independent perfect sample, this would eliminate the convergence issue, and the chains would then yield independent unbiased parameter estimates.

## 6 Acknowledgements

This work was supported in part by NSF Grant IIS-0083418. The author is also grateful to D. Wilkinson, J. Kadane, two anonymous referees, and to Brent Frye and the Pittsburgh Brain Imaging Research Center for assistance with and use of their Beowulf cluster for the case study considered in this paper.

## A An Example of Careless Blocking

The following example illustrates how asynchronous updates of separate parts of the state-space may lead to an incorrect invariant distribution for a Metropolis-Hastings chain.

**Example 1.1:** Suppose that the target distribution  $\pi$  is bivariate normal with mean  $(0, 0)^T$  and covariance matrix

$$\begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix},$$

and suppose that  $X_1 = (X_{11}, X_{12})^T \sim \pi$ .

1. The standard Gibbs sampler would generate  $X_2 = (X_{21}, X_{22})$  by drawing  $X_{21}$  from a  $N(\rho X_{12}, 1 - \rho^2)$  distribution, and then drawing  $X_{22}$  from a  $N(\rho X_{21}, 1 - \rho^2)$  distribution. This update can be written in the form

$$X_2 = \begin{bmatrix} 0 & \rho \\ 0 & \rho^2 \end{bmatrix} X_1 + \begin{bmatrix} 1 & 0 \\ \rho & 1 \end{bmatrix} \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix},$$

where  $\epsilon_1, \epsilon_2$  are iid normal random variables with mean zero and variance  $(1 - \rho^2)$ . It is easily verified that this gives  $X_2 \sim \pi$ , and thus  $\pi$  is indeed the invariant distribution of the chain.

2. The asynchronous update sampler would carry out updates of the two components by drawing from their respective full-conditional distributions concurrently. This means  $X_{21}$  would be drawn from a  $N(\rho X_{12}, 1 - \rho^2)$  distribution, and  $X_{22}$  would be drawn from an independent  $N(\rho X_{11}, 1 - \rho^2)$  distribution. The update can be written in the form

$$X_2 = \begin{bmatrix} 0 & \rho \\ \rho & 0 \end{bmatrix} X_1 + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix},$$

where  $\epsilon_1$  and  $\epsilon_2$  are as defined in the previous case. It is easily checked that this yields

$$X_2 \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho^3 \\ \rho^3 & 1 \end{bmatrix} \right).$$

Since the distribution of  $X_2$  is not the same as that of  $X_1 \sim \pi$ ,  $\pi$  cannot be the limiting distribution for the chain.

□

A simple representation of the state-updates for a bivariate distribution  $\pi$  is given in Figure 4, with the arrows indicating dependencies in updates.

## B Blocking for the Generalized State-Space Model

Consider the model

$$\begin{aligned} X_{t+1} &\sim f(\cdot; X_t, \vartheta), \quad t \geq 1 \\ Y_t &\sim g(\cdot; X_t, \vartheta), \end{aligned}$$

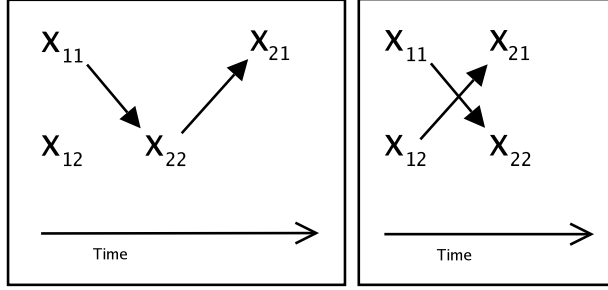


Figure 4: Left: Standard Gibbs sampling updates for a bivariate target distribution, where one component is updated according to its full-conditional distribution given the other, and then the process is repeated for the other component. Right: Asynchronous sampling updates, where one processor updates each component based on its full-conditional distribution, but updates are carried out concurrently.

where  $f(\cdot; X_t, \vartheta)$  and  $g(\cdot; X_t, \vartheta)$  are some probability density functions which depend on  $X_t$ , as well as a parameter vector  $\vartheta$ .  $\{X_t\}$  is a latent Markov chain, and  $\{Y_t\}$  is a sequence of observations whose distributions are determined by  $\{X_t\}$ . Some assumption is made about the marginal distribution  $f(X_1; \vartheta)$  - in many cases one can use the stationary distribution of  $\{X_t\}$  (assuming it exists) here. In Markov chain Monte Carlo analyses of such models (see, e.g. Carlin et al., 1992), one observes  $\{y_1, \dots, y_n\}$ , and typically defines the state-space  $\Theta$  to include all latent variables  $\{X_1, \dots, X_n\}$  as well as the parameter  $\vartheta$ . For convenience, define  $X = \{x_1, \dots, x_n\}$  and  $Y = \{y_1, \dots, y_n\}$ . The posterior distribution is then

$$\pi(\vartheta, X) = kp(\vartheta)p(X|\vartheta)p(Y|X, \vartheta),$$

where  $k$  is a normalizing constant which depends on  $y_1, \dots, y_n$ ,  $p(X|\vartheta) = f(x_1; \vartheta) \prod_{t=2}^n f(x_t; x_{t-1}, \vartheta)$ , and  $p(Y|X, \vartheta) = \prod_{t=1}^n g(y_t; x_t, \vartheta)$ . Now suppose (for the sake of giving an example) that  $n = 300$ . One could then decompose the state-space into

$$\Theta_0 = (\vartheta, X_{100}, X_{200}), \quad \Theta_1 = (X_1, \dots, X_{99}), \quad \Theta_2 = (X_{101}, \dots, X_{199}), \quad \Theta_3 = (X_{201}, \dots, X_{300}).$$

Then under the distribution  $\pi$ , Property 3.1 holds. To see this, note that

$$\begin{aligned} \pi(\theta_0, \theta_1, \theta_2, \theta_3) &= kp(\vartheta)p(\theta_0, \theta_1, \theta_2, \theta_3|\vartheta)p(Y|X, \vartheta) \\ &= kp(\vartheta)p(\theta_1, \theta_2, \theta_3|\vartheta, \theta_0)p(\theta_0|\vartheta)p(Y|X, \vartheta) \\ &= kp(\vartheta)p(\theta_1|\vartheta, \theta_0)p(\theta_2|\vartheta, \theta_0)p(\theta_3|\vartheta, \theta_0)p(\theta_0|\vartheta)p(Y|X, \vartheta). \end{aligned} \quad (5)$$

(The factorization in the last line here is valid because of the Markov property of  $\{X_1, \dots, X_n\}$ .) Next, let  $Y_0 = \{y_{100}, y_{200}\}$ ,  $Y_1 = \{y_1, \dots, y_{99}\}$ ,  $Y_2 = \{y_{101}, \dots, y_{199}\}$ , and  $Y_3 = \{y_{201}, \dots, y_{299}\}$ , and observe that

$$p(Y|X, \vartheta) = p(Y_0|\theta_0)p(Y_1|\theta_1, \vartheta)p(Y_2|\theta_2, \vartheta)p(Y_3|\theta_3, \vartheta). \quad (6)$$



It follows from equations (5) and (6) that

$$\pi(\theta_1|\theta_{-1}) = \frac{\pi(\theta_0, \theta_1, \theta_2, \theta_3)}{\int \pi(\theta_0, \theta_1, \theta_2, \theta_3)d\theta_1} = \frac{p(\theta_1|\theta_0, \vartheta)p(Y_1|\theta_1, \vartheta)}{\int p(\theta_1|\theta_0, \vartheta)p(Y_1|\theta_1, \vartheta)d\theta_1}.$$

Thus  $\pi(\theta_1|\theta_{-1})$  does not depend on  $\theta_2$  or  $\theta_3$ , so  $\pi(\theta_1|\theta_{-1}) = \pi(\theta_1|\theta_0)$ . Analogous results also hold for  $\pi(\theta_2|\theta_{-2})$  and  $\pi(\theta_3|\theta_{-3})$ . Thus Property 3.1 holds for this particular target distribution and decomposition of the state-space.

## References

- Jan Beran. *Statistics for Long-Memory Processes*. Chapman and Hall, 1994.
- R. Bradford and A. Thomas. Markov chain Monte Carlo methods for family trees using parallel processor. *Statistics and Computing*, 6:67–75, 1996.
- A.E. Brockwell and N.H. Chan. Long memory dynamic Tobit models. Technical Report 798, Carnegie Mellon University Dept. of Statistics, 2004.
- A.E. Brockwell and J.B. Kadane. Identification of regeneration times in MCMC simulation, with application to adaptive schemes. *Journal of Computational and Graphical Statistics*, To appear, 2004.
- P.J. Brockwell and R.A. Davis. *Time Series: Theory and Methods*. Springer, second edition, 1991.
- B.P. Carlin, N.G. Polson, and D.S. Stoffer. A Monte Carlo approach to nonnormal and nonlinear state-space modeling. *Journal of the American Statistical Association*, 87(418): 493–500, 1992. ISSN 0162-1459.
- J. Choi, J. Dongarra, R. Pozo, and D. Walker. ScaLAPACK: A scalable linear algebra library for distributed memory concurrent computers. In *Proceedings of the Fourth Symposium on the Frontiers of Massively Parallel Computation*, pages 120–127. IEEE Computer Society Press, 1992.
- W.R. Gilks, S. Richardson, and D.J. Spiegelhalter. *Markov Chain Monte Carlo in Practice*. Chapman & Hall/CRC, 1996.
- P.W. Glynn and P. Heidelberger. Analysis of initial transient state deletion for parallel steady-state simulations. *SIAM Journal on Scientific and Statistical Computing*, 13:904–922, 1992.
- W. Gropp, E. Lusk, and A. Skjellum. *Using MPI: Portable Parallel Programming with the Message-Passing Interface*. The MIT Press, 1999.

- Mascagni and Srinivasan. SPRNG: A scalable library for pseudorandom number generation. *ACMTMS: ACM Transactions on Mathematical Software*, 26, 2000.
- P. Mykland, L. Tierney, and B. Yu. Regeneration in Markov chain samplers. *Journal of the American Statistical Association*, 90:233–241, 1995.
- J.G. Propp and D.B. Wilson. Exact sampling with coupled Markov chains and applications to statistical mechanics. *Random Structures and Algorithms*, 9:223–252, 1996.
- J. Rosenthal. Parallel computing and Monte Carlo algorithms. *Far East J. Theor. Stat.*, 4: 207–236, 2000.
- F. Sowell. Maximum likelihood estimation of stationary univariate fractionally integrated time series models. *Journal of Econometrics*, 53:165–188, 1992.
- M. Whitley and S.P. Wilson. Parallel algorithms for Markov chain Monte Carlo methods in latent spatial gaussian models. *Statistics and Computing*, 14(3), 2004.
- D.J. Wilkinson. Parallel Bayesian computation. In *Handbook of Parallel Computing and Statistics*, chapter 18. Dekker, 2004.