

From *Performance Tools '97*, St. Malo, France, June 3–6, 1997.

## AN EFFICIENT DISK-BASED TOOL FOR SOLVING VERY LARGE MARKOV MODELS

Daniel D. Deavours and William H. Sanders\*

Center for Reliable and High-Performance Computing  
Coordinated Science Laboratory  
University of Illinois at Urbana-Champaign  
{deavours,whs}@crhc.uiuc.edu

**Abstract.** Very large Markov models often result when modeling realistic computer systems and networks. We describe a new tool for solving large Markov models on a typical engineering workstation. This tool does not require any special properties or a particular structure in the model, and it requires only slightly more memory than what is necessary to hold the solution vector itself. It uses a disk to hold the state-transition-rate matrix, a variant of block Gauss-Seidel as the iterative solution method, and an innovative implementation that involves two parallel processes: the first process retrieves portions of the iteration matrix from disk, and the second process does repeated computation on small portions of the matrix. We demonstrate its use on two realistic models: a Kanban manufacturing system and the Courier protocol stack, which have up to 10 million states and about 100 million nonzero entries. The tool can solve the models efficiently on a workstation with 128 Mbytes of memory and 4 Gbytes of disk.

### 1 Introduction

A wide variety of high-level specification techniques now exist for Markov models. These include, among others, stochastic Petri nets, stochastic process algebras, various types of block diagrams, and non-product form queuing networks. In most cases, very large Markov models result when one tries to model realistic systems using these specification techniques. The Markov models are typically quite sparse (adjacent to few nodes), but contain a large number of states. This problem is known as the “largeness problem.” Techniques that researchers have developed to deal with the largeness problem fall into two general categories: those that avoid the large state space (for example, by lumping,) and those that tolerate the large state space (for example, by recognizing that the model has a special structure and storing it in a compact form). While many largeness avoidance and tolerance techniques exist, few are applicable to models without special structure. Methods are sorely needed that permit the solution of very

---

\* This work was supported, in part, by NASA Grant NAG 1-1782.

large Markov models without requiring them to have special properties or a particular structure.

In this paper, we describe a new tool for solving Markov models with very large state spaces on a typical engineering workstation. The tool makes no assumptions about the underlying structure of the Markov process, and requires little more memory than that necessary to hold the solution vector itself. It uses a disk to hold the state-transition-rate matrix, a variant of block Gauss-Seidel as the iterative solution method, and an innovative two-process implementation that effectively overlaps retrieval of blocks of the state-transition-rate matrix from disk and computation on the retrieved blocks. The tool can solve models with ten million states and about 100 million transitions on a machine with 128 Mbytes of main memory. The state-transition-rate matrix is stored on disk in a clever manner, minimizing overhead in retrieving it from disk. In addition, the tool employs a dynamic scheme for determining the number of iterations to perform on a block before beginning on the next, which we show empirically to provide a near optimum time to convergence. Solution time is typically quick even for very large models, with only about 20% of the CPU time spent retrieving blocks from disk and 80% of the CPU resources available to perform the required computation.

In addition to describing the architecture and implementation of the tool itself, we illustrate its use on two realistic models: one of a Kanban manufacturing system [2], and another of the Courier protocol stack executing on a VME bus-based multiprocessor [14]. Both models have appeared before in the literature, and are excellent examples of models that have very large state spaces for realistic system parameter values. In particular, both models have been used to illustrate the use of recently developed Kronecker-based methods [2, 7], and the Courier protocol has been used to illustrate an approximate method based on lumping [14]. Both numerical results and solution times are presented for each model and, when possible, compared to previously obtained values and solution times. In each case we can obtain an exact solution (to the desired precision) in significantly less time than previously reported using Kronecker-based methods. It is thus our belief that if sufficient disk space is available to hold the state-transition-rate matrix, our approach is the method of choice for exact solutions.

The remainder of the paper is organized as follows. First, in Section 2, we address issues in the choice of a solution method for very large Markov models, comparing three alternatives: Kronecker-based methods (e.g., [7, 2]), “on-the-fly” methods [3], and the disk-based method that we ultimately choose. This section presents clear arguments for the desirability of disk-based methods if sufficient disk space is available. Section 3 then describes the architecture and implementation of the tool, describing solutions to issues we faced in building a practical implementation. Finally, Section 4 presents the results of the use of the tool on the two models described earlier.

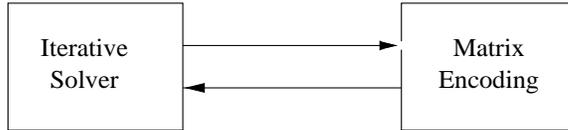


Fig. 1. Solution paradigm.

## 2 The Case for Disk-Based Methods

For our tool implementation, we considered three numerical solution techniques for tolerating large state spaces: Kronecker-based techniques, “on-the-fly” techniques, and disk-based techniques. To evaluate each method, we introduce a paradigm based on Figure 1. Here, we divide the numerical solution process into the iterative solver and the matrix encoding. The key to solving large matrices is to encode the matrix so that it takes little main memory (RAM), but still allows quick access to matrix elements. The iterative solver is thus a data consumer, and the matrix encoder is a data producer. We would like for both to be as fast as possible to obtain a solution quickly. An additional important factor is how effectively a particular iterative method uses the data it consumes. For example, certain iterative methods, such as Gauss-Seidel [12] and adaptive Gauss-Seidel [5], typically do more effective work with the same number of accesses to the matrix as Jacobi or the Power method, and hence do not require as high a data production rate to efficiently obtain a solution. We want to find a fast but general matrix encoding scheme and an effective iterative method with a low data consumption rate.

The first class of encoding schemes we consider are those of Kronecker-based methods. These methods require and make use of the fact that in certain models, particular parts of the model (called submodels) interact with one another in a limited way. One way to insure a model has this structure is to construct it according to a prescribed set of rules from smaller models, as is done, for example, in the case of stochastic automata networks [12]. If one follows these rules, one may easily express the transition rate matrix for the entire model as a function of Kronecker operators on the transition rate matrices of the submodels.

More recently there has been work on a type of model decomposition called superposed generalized stochastic Petri nets (SGSPNS) [1, 2, 4, 6, 7]. SGSPNs are essentially independent models that may be joined by synchronization of a transition. We believe [2] to be the state of the art in Kronecker operator methods, and although the more recent techniques can solve a much larger class of models than originally proposed in [4], they are still restrictive in the models that they can effectively solve.

To evaluate the speed of the Kronecker operator methods, we observe rates in which the iterative solver and matrix encoding operate. We have observed that on our computer (120 MHz Hewlett-Packard Model C110), the SOR iterative solver can consume data at a rate of about 50 Mbyte/second. From numbers published by Kemper [7], we estimate that his implementation of the Kronecker-

based method can produce data at a rate of 700 Kbyte/second on an 85 MHz Sparc 4, and we extrapolate that the rate would be about 2 Mbyte/second on our HP C110. Since both the data production and consumption require the CPU, the whole process will proceed at a rate of about 1.9 Mbyte/second. Kemper's method is also restricted to Jacobi or the Power method, which usually exhibit poor convergence characteristics, so the effectiveness of its use of generated data is low. Ciardo and Tilgner [2] present their own tool, but they do not present data in such a way that we can analyze the data generation rate. We can compare actual times to solution for their benchmark model, however, and do so in Section 4. Ciardo gives algorithms to perform Gauss-Seidel on a Kronecker representation in [1], but has not yet built a tool with which we can compare our approach.

The second class of encoding schemes we considered for implementation in this tool are "on-the-fly" methods introduced in [3]. On-the-fly methods have none of the structural restrictions of Kronecker-based methods, and they can operate on nets with general enabling predicate and state change functions, such as are present in stochastic activity networks [9, 10]. In addition, they can obtain a solution with little additional memory, or perhaps even less memory than needed by SGSPN solvers, while at the same time using Gauss-Seidel or variants. However, the prototype implementation described in [3] generates data at about 440 Kbyte/second on a HP C110. Although [3] introduces iterative methods that are usually more effective than Jacobi or the Power method in their use of data, the overall solution speed for these methods will be somewhat slower than for Kronecker-based methods, but still reasonable, given that they can be used without restrictions on the structure of a model.

The final class we considered was that of disk-based methods, where the workstation disk holds an encoding of the state-transition matrix. If we can find an iterative method that accesses data from the state-transition-rate matrix in a regular way and use a clever encoding, disks can deliver data to an iterative algorithm at a high rate (5 Mbyte/second or higher) with low CPU overhead. Furthermore, high performance disks are inexpensive relative to the cost of RAM, so we would like to find a way to utilize disks effectively. Experimental results show that if we can do both disk I/O and computation in parallel, we can perform Gauss-Seidel at a rate of 5 Mbyte/second while using the CPU only 30% of the time. Thus disk-based methods have the potential to greatly outperform Kronecker and on-the-fly methods, at the cost of providing a disk that is large enough to hold the state-transition-rate matrix of the Markov model being solved. The challenge is to find a more effective solution method that has a data consumption rate of about 5 Mbytes/second at 80% CPU utilization.

Clearly, the method of choice depends on the nature of the model being solved, and the hardware available for the solution. If the state-transition-rate matrix is too large to fit on available disk space and the model meets the requirements of Kronecker-based methods, then they should be used. If the model does not fit on disk and does not meet the requirements of Kronecker-based methods, on-the-fly methods should be used. However, SCSI disks are inexpensive relative

to RAM (in September 1996, approximately \$1400 for 4 Gbyte fast wide SCSI), so space may inexpensively be made available to store the state-transition-rate matrix. Since a single disk can provide the high data production rate only for sequential disk access, the efficiency of disk-based methods will depend on whether we can find a solution algorithm that can make effective use of the data in the sequential manner. We discuss how to do this in the following sections.

### 3 Tool Architecture and Implementation

In this section, we discuss the architecture of our tool and its implementation on an HP workstation. In particular, we discuss the basic block Gauss-Seidel (BGS) algorithm and how it maps onto a program or set of programs that run on a workstation. An important issue we solve is how to effectively do computation and disk I/O in parallel. We develop a flexible implementation with many tunable parameters that can vary widely on different hardware platforms and models.

The mathematics of BGS is generally well known (see [12], for example). We wish to solve for the steady state probability vector  $\pi$  given by  $\pi Q = 0$ . To review BGS briefly, partition the state-transition-rate matrix  $Q$  into  $N \times N$  submatrices of (roughly) the same size, labeled  $Q_{ij}$ . BGS then solves

$$\Pi_i^{(k+1)} Q_{ii} = - \left( \sum_{j=1}^{i-1} \Pi_j^{(k+1)} Q_{ji} + \sum_{j=i+1}^N \Pi_j^{(k)} Q_{ji} \right) \quad (1)$$

for  $\Pi_i$  for  $i$  ranging from 1 to  $N$ , where  $\Pi_i$  is the corresponding subvector of  $\pi$ . This is called the  $k$ -th BGS *iteration*. Solving for  $\Pi_i$  can be done by any method; our tool uses (point) Gauss-Seidel. One Gauss-Seidel iteration to solve (1) is called an *inner* iteration, and solving (1) for  $1 \leq i \leq n$  is an *outer* iteration.

The sequential algorithm for a single BGS iteration follows directly. In particular, let  $r \in \mathcal{R}^n$  be an auxiliary variable.

```

for  $i = 1$  to  $N$ 
   $r = 0$ 
  for  $j = 1$  to  $N | j \neq i$ 
     $r = r - \Pi_j Q_{ji}$ 
  Solve  $\Pi_i Q_{ii} = r$  for  $\Pi_i$ 

```

One can easily see that the access to  $Q$  is very predictable, so we may have blocks of  $Q$  ordered on disk in the same way that the program accesses them. This way the program accesses the file representing  $Q$  sequentially entirely throughout an iteration. One could then easily write an implementation of BGS and a utility to write  $Q$  appropriately to a file. What is not trivial is to build a tool that overlaps computation (the solution of  $\Pi_i Q_{ii} = r$ ) and reading from disk in a flexible, efficient way.

*Tool Architecture* Our solution to this is to have two cooperating processes, one of which schedules disk I/O, and the other of which does computation. Obviously, they must communicate and synchronize activity. We use System V interprocess communication mechanisms since they are widely available and simple to use. For synchronization, we use semaphores, and for passing messages, such as a block of  $Q$ , we use shared memory. We call the process that schedules I/O the *I/O process*, and we call the process that solves  $\Pi_i Q_{ii} = r$  the *compute process*. To minimize memory usage, we want to have as few blocks of  $Q$  in memory at one time as possible, so we must be careful how we compute the step  $r = r - \Pi_j Q_{ji}$ ,  $\forall j \neq i$ . For simplicity, we assign the task of computing  $r$  to the I/O process.

We first looked at several large matrices that were generated by GSPN models. (We looked at GSPNs because they can easily create large transition rate matrices, not because our solution technique is limited to them.) We noticed that the matrix is usually very banded; that is, for reasonable choices for  $N$ , the number of non-zero elements in the blocks  $Q_{i,j} : |i - j| > 1$  is small, if not zero. By lumping all the blocks in a column into a smaller number (three) of larger blocks, we can eliminate the overhead of reading small or empty blocks. For the  $i$ -th column, we call  $Q_{i,i}$  the *diagonal* block;  $Q_{i-1,i}$  is the *conflict* block; all other blocks are lumped into a single block that we call the *off* block. We use the term *off block* because it includes all the off diagonal blocks except the conflict block. Let D represent the diagonal block, C the conflict block, and O the off block. The following represents a matrix where  $N = 4$ .

$$\left( \begin{array}{c|c|c} \text{D} & \text{O} & \text{C} \\ \hline \text{C} & \text{D} & \text{O} \\ \hline \text{O} & \text{C} & \text{D} & \text{O} \\ \hline \text{O} & \text{O} & \text{C} & \text{D} \end{array} \right)^T$$

The reason we have a conflict block will be apparent soon.

Lumping several blocks into the off block complicates 1), but does not require any extra computation. The actual mechanics of the computation of  $\Pi Q_{\text{off},i}$  are no different than for the computation of  $\Pi_j Q_{ji}$ . For the formula  $r = \Pi Q_{\text{off},i}$ , we compute  $r = \sum_{k \neq i, i-1} \Pi_k Q_{ki}$ . We may now compute  $r$  the following way:

$$\begin{aligned} r &= -\Pi Q_{\text{off},i} \\ r &= r - \Pi_{i-1} Q_{\text{conflict},i} \end{aligned}$$

Let us denote  $r_i = \Pi_i Q_{ii}$  to distinguish between different  $r$  vectors. In order to make the computation and disk I/O in parallel, the program must solve  $\Pi_i Q_{ii} = r_i$  while at the same time compute  $r_{i+1}$ . Therefore, while the compute process is solving  $\Pi_i Q_{ii} = r_i$ , the I/O process is prefetching  $Q_{i+1,i+1}$ , and reading  $Q_{\text{off},i+1}$  and  $Q_{\text{conflict},i+1}$  to compute  $r_{i+1}$ . Notice that when computing  $r_{i+1}$ , we need the most recent value of  $\Pi_i$  to multiply by  $Q_{\text{conflict},i}$ , which introduces a data dependency. Thus, we can not completely compute  $r_{i+1}$  while in parallel computing  $\Pi_i$ . (We could also use a less recent value of  $\Pi_i$ , but that would reduce the effectiveness of BGS.)

Shared variables:  $\Pi$ ,  $Q_{\text{diag}0}$ ,  $Q_{\text{diag}1}$ ,  $r_0$ ,  $r_1$   
Semaphores:  $S_1$  locked,  $S_2$  unlocked

<p><i>Compute Process</i></p> <p>Local variable (unshared): <math>t</math>  <math>t = 0</math>  while not converged    for <math>i = 1</math> to <math>N</math>      <b>Lock</b>(<math>S_1</math>)      for <math>j = 1</math> to <math>MinIter</math>        Do GS iteration: <math>\Pi_i Q_{\text{diag}t} = r_t</math>        <math>j = MinIter + 1</math>        while <math>j \leq MaxIter</math> and          I/O process not blocked on <math>S_2</math>          Do GS iteration: <math>\Pi_i Q_{\text{diag}t} = r_t</math>          <math>j = j + 1</math>        <b>Unlock</b>(<math>S_2</math>)      <math>t = \bar{t}</math></p>	<p><i>I/O Process</i></p> <p>Local variable (unshared): <math>t</math>, <math>Q_{\text{tmp}}</math>  <math>t = 0</math>  do forever    for <math>i = 1</math> to <math>N</math>      <math>Q_{\text{diag}t} = \text{disk read}(Q_{ii})</math>      <math>Q_{\text{tmp}} = \text{disk read}(Q_{\text{off},i})</math>      <math>r_t = -\Pi Q_{\text{tmp}}</math>      <math>Q_{\text{tmp}} = \text{disk read}(Q_{\text{conflict},i})</math>      <b>Lock</b>(<math>S_2</math>)      <math>r_t = r_t - \Pi_{i-1} Q_{\text{tmp}}</math>      <b>Unlock</b>(<math>S_1</math>)      <math>t = \bar{t}</math></p>
---	---

**Fig. 2.** Compute and I/O processes for BGS algorithm.

Finally, we add synchronization to ensure that the I/O process has the most recent version of  $\Pi_i$  to compute  $r_{i+1}$ . The full algorithm we use is presented in Figure 2. We used a large, shared memory array to represent the steady state probability vector  $\Pi$ , two shared diagonal block buffers  $Q_{\text{diag}0}$  and  $Q_{\text{diag}1}$ , and two  $r$  vectors  $r_0$  and  $r_1$ . The processes share two diagonal block and  $r$  variables so that one can be used to compute (1) while the other one is being prepared for the next computation. The processes also share two locking variables,  $S_1$  and  $S_2$ , which they use to communicate and control the relative progress of the other process.

*Compute Process* We first explain the compute process. A local variable  $t$  alternates between 0 and 1, which indicates which of the two shared block and  $r$  variables the process should use. After each step,  $t$  is alternated between 0 and 1, which we denote  $t = \bar{t}$ . The function **Lock**( $S_1$ ) will lock  $S_1$  if  $S_1$  is unlocked. If  $S_1$  is already locked, it will block until  $S_1$  is unlocked (by the I/O process); then it will lock  $S_1$  and proceed. While the compute process is blocked on  $S_1$ , it uses no CPU resources.

The compute process has two parameters,  $MinIter$  and  $MaxIter$ . The compute process is guaranteed to do at least  $MinIter$  Gauss-Seidel inner iterations to approximately solve (1). Then the compute process will proceed to do up to  $MaxIter$  iterations or until the I/O process is complete with the current file I/O and is waiting for the compute process to unlock  $S_2$ , whichever comes first. This allows the compute process to do a dynamic number of Gauss-Seidel iterations, depending on how long the I/O process takes to do file I/O. We ignore the boundary conditions in the figures for simplicity. If  $i - 1 = 0$ , for example, then

we use  $N$  for  $i - 1$  instead.

The convergence criterion we use in this tool is a modification to the  $\|\pi^{(k+1)} - \pi^{(k)}\|_\infty$  criterion. In particular, we compute  $\|II_i^{(k+1)} - II_i^{(k)}\|_\infty$  for the *first* inner iteration and take the  $\max_i \|II_i^{(k+1)} - II_i^{(k)}\|_\infty$  to be the number we use to test for convergence. We use this for two reasons: the first inner iteration usually results in the greatest change of  $II_i$ , so computing the norm for all inner iterations is usually wasteful, and the computation of the norm takes a significant amount of time. We have observed experimentally that this measured norm is at least as good as the  $\|\pi^{(k+1)} - \pi^{(k)}\|_\infty$  criterion.

The dynamic nature of the iteration count is an interesting feature of this tool. If the system on which the program is running is doing other file I/O and slowing the I/O process down, the compute process may continue to proceed to do useful work. At some point, however, additional Gauss-Seidel iterations may not be useful at all, presumably after *MaxIter* inner iterations, so the process will stop doing work and block waiting for  $S_1$  to become unlocked. Choosing a good *MinIter* and *MaxIter* is difficult and requires some knowledge about the characteristics of the transition rate matrix. If we allow the compute process to be completely dynamic, some blocks may consistently get fewer inner iterations and converge more slowly than other blocks, causing the whole system to converge slowly. In Section 4, we show some experimental results of varying these parameters.

*Input/Output Process* The I/O process is straightforward. The greatest complexity comes in managing the semaphores properly. This is a case of the classical producer-consumer or bounded buffer problem, and we defer the reader to [13] or a similar text on operating systems to show the motivation and correctness of this technique. The primary purpose of the I/O process is to schedule disk reads and compute  $r_t$ . It does this by issuing a C function to read portions of the file directly into the shared block variable or the temporary block variable. Because the I/O process may execute in parallel with the compute process, the I/O process may issue read requests concurrently with the computation, and since file I/O uses little CPU (under 20%), we can effectively parallelize computation and file I/O on a modern, single-processor workstation.

This implementation of BGS uses relatively little memory. The size of the steady state probability vector  $II$  is proportional to the number of states, which is unavoidable using BGS or any other exact method. Other iteration methods, such as Jacobi, require additional vectors of the same size as  $II$ , which our program avoids. Two diagonal blocks,  $Q_{\text{diag0}}$  and  $Q_{\text{diag1}}$ , are necessary; the compute process requires one block to do the inner iteration, and the I/O process reads the next diagonal block at the same time. Two  $r$  variables are also required for the same reason. Finally, the I/O process requires a temporary variable  $Q_{\text{tmp}}$  to hold  $Q_{\text{off}}$  and  $Q_{\text{conflict}}$ . We could eliminate  $Q_{\text{tmp}}$  by instead using  $Q_{\text{diag}t}$ , but doing so would require us to reverse the order in which we read the blocks, causing us to read  $Q_{\text{diag}t}$  last. This would reduce the amount of time we could overlap computation and file I/O. We chose to maximize parallelization of computation

N	States	NZ Entries	Size (MB)	$e_1$	$e_2$	$e_3$	$e_4$	$\tau$
1	160	616	0.008	0.90742	0.67136	0.67136	0.35538	0.09258
2	4,600	28,128	0.34	1.81006	1.32851	1.32851	0.76426	0.17387
3	58,400	446,400	5.3	2.72211	1.94348	1.94348	1.52460	0.23307
4	454,475	3,979,850	47	3.64641	2.51298	2.51298	1.50325	0.27589
5	2,546,432	24,460,416	290	4.58301	3.03523	3.03523	1.81096	0.30712
6	11,261,376	115,708,992	1,367	5.53098	3.50975	3.50975	2.07460	0.33010

**Table 1.** Characteristics and reward variables for the Kanban model.

at the expense of a modest amount of memory.

## 4 Results

To better understand the algorithms presented in the previous section, we implemented them and tested the resulting tool on several large models presented in the literature. We present the models here and discuss the performance measures we took in order to better understand the issues in building and using a tool to solve large matrices, so we are not so interested here in the results of the models as much as using the models to understand the characteristics of our solver. All the solutions we present here, with the exception of one, can be solved on our HP workstation with 128 Mbyte of RAM (without using virtual memory) and 4 Gbyte of fast disk memory.

*Kanban Model* The Kanban model we present was previously used by Ciardo and Tilgner [1, 2] to illustrate a Kronecker-based approach. They chose this model because it has many of the characteristics that are ideal for superposed GSPN solution techniques, and it also does not require any mapping from the product space to the tangible reachable states. We refer to [2] for a description of the model and specification of the rates in the model. Briefly, the model is composed of four subnets. At each subnet, a token enters, spends some time, and exits or restarts with certain probabilities. Once the token leaves the first subnet, it may enter the second or third subnet, and to leave the system, the token must go through the fourth subnet. We chose to solve the model where the synchronizing transitions are timed transitions.

Table 1 shows some information about the model and the corresponding transition rate matrix. Here,  $N$  represents the maximum number of tokens that may be in a subnet at one time. There are two important variables that we may vary, the number of blocks and the number of inner iterations, that greatly affect performance. We present two experiments. First, we vary the number of blocks while keeping the number of inner iterations fixed, and second, we vary the number of inner iterations while keeping the number of blocks fixed.

For the first experiment, we use the Kanban model where  $N = 5$ . We divide the transition rate matrix into  $32 \times 3$  blocks, and perform a constant number of inner iterations. We vary the number of inner iterations from 1 to 20. The results

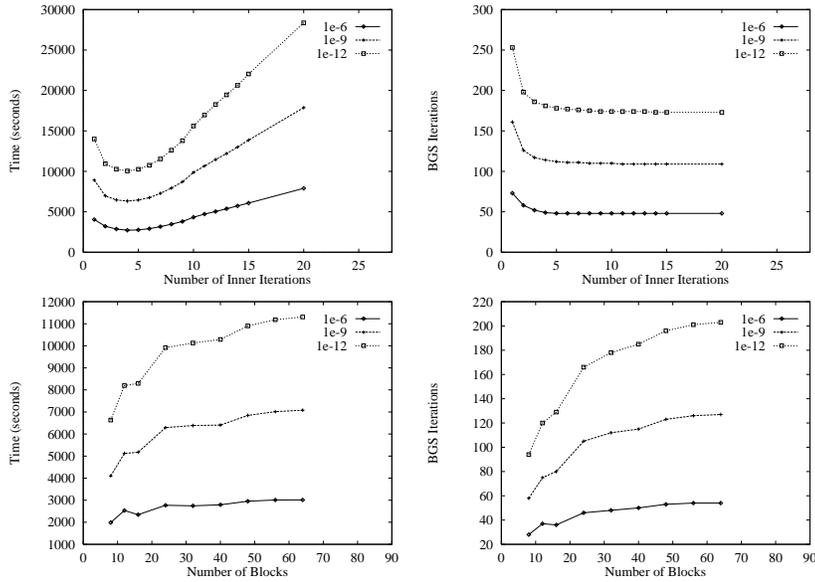
of the solution execution time and the number of BGS iterations are shown in the top two graphs in Figure 3. All the timing measurements that we present in this paper are “wall clock” times. The plots show the time to achieve three levels of accuracy based on the modified  $\|\Pi^{(k+1)} - \Pi^{(k)}\|_\infty < \{10^{-6}, 10^{-9}, 10^{-12}\}$  convergence criterion explained in Section 3.

Figure 3 shows how doing an increased number of inner iterations yields diminishing returns, so that doing more than about 7 inner iterations does not significantly help reduce the number of BGS iterations. For this model, setting *MaxIter* to 6 or 7 makes sense. It also shows that the optimal number of inner iterations with respect to execution time is 4. For fewer than four inner iterations, the compute process spends time idle and waiting for the I/O process. This leads us to choose *MinIter* to be 3 or 4.

It is interesting to note that solving this model with a dynamic number of inner iterations takes 10,436 seconds, which is more time than is required if we fix the number of inner iterations to be 3, 4, or 5 (10269, 10044, and 10252 seconds respectively). We observed that some blocks always receive 4 or fewer inner iterations, while others always receive 7 or more. This shows us several important things. First, some blocks always receive more iterations than others, and we know that the solution vector will converge only as fast as its slowest converging component. Second, we argued above that doing more than 7 inner iterations is wasteful, so allowing the number of inner iterations to be fully dynamic is wasteful since the I/O process does not read data quickly enough to keep the compute node doing useful work. Finally, if the compute process is always doing inner iterations, it checks to see if the I/O process is blocked on  $S_2$  only after completing an inner iteration. This requires the I/O process to always block on  $S_2$  and wait for the compute process to complete its inner iteration, which is wasteful since the I/O process is the slower of the two processes.

For the next experiment, we set the number of inner iterations to be 5, vary the number of blocks, and observe convergence rate, execution time, and memory usage. The bottom two plots of Figure 3 range the number of blocks from 8 to 64 and plot execution time and number of iterations respectively for the convergence criteria  $\|\Pi^{(k+1)} - \Pi^{(k)}\|_\infty < \{10^{-6}, 10^{-9}, 10^{-12}\}$ . Table 2 shows how memory usage varies with the number of blocks. Notice that between 8 and 64 blocks, the execution time is nearly double while the memory usage is about one third. We see that there is clearly a memory/speed tradeoff. Note that the solution vector for a 2.5 million state model alone takes about 20 Megabytes of memory.

Finally, as a basis for comparison, we present results given in [2] in Table 3 and compare the solution times to those of our tool. The column titled ‘Case 1’ represents the tool in [2] with mapping from the product space to tangible reachable states enabled, while ‘Case 2’ is with no mapping (an idealized case). Cases 1 and 2 are performed on a Sony NWS-5000 workstation with 90 MB of memory. We present no results for  $N = 1, 2, 3$  because the matrix was so small that the operating systems buffered the entire file in memory. In addition to computing the reward variables (see Table 1) for the Kanban model to greater accuracy than [2], we were also able to solve for the case where  $N = 6$ .



**Fig. 3.** Performance graphs of Kanban model ( $N = 5$ ).

N	Memory (MB)
8	95
12	71
16	59
24	48
32	40
40	37
48	35
56	33
64	32

**Table 2.** Number of blocks versus memory.

N	Case 1	Case 2	BGS time	No. Blocks	Memory
1	1 s	1s	-	-	-
2	13 s	2 s	-	-	-
3	310 s	2 s	-	-	-
4	4,721 s	856 s	225 s	4	28 MB
5	22,215 s	6,055 s	2,342 s	16	59 MB
6	-	-	18,563 s	128	111 MB

**Table 3.** Comparison of performance.

*Courier Protocol Model* The second model we examine is a model of the Courier protocol given in [8, 14]. This is a model of an existing software network protocol stack that has been implemented on a Sun workstation and on a VME bus-based multiprocessor. The model is well specified in [8, 14]. The GSPN models only a one-way data flow through the network. For our experiment, we are only interested in varying the window size  $N$ . The transport space and fragmentation ratio is kept at one. Varying  $N$  corresponds to initially placing  $N$  tokens in a place, and it has a substantial impact on the state space size!

	$N = 1$	$N = 2$	$N = 3$	$N = 4$	$N = 5$	$N = 6$
States	11,700	84,600	419,400	1,632,600	5,358,600	15,410,250
Nonzero	48,330	410,160	2,281,620	9,732,330	34,424,280	105,345,900
Matrix	0.6	5	28	118	414	1,264
Blocks	4	4	4	32	64	128
Generation Time (s)	5	38	218	938	3,716	11,600
Conversion Time (s)	3	24	136	581	2,076	*14,482
Solution Time (s)	2	16	143	1,138	6,040	20,742
Iterations	18	19	23	49	69	85
Memory (MB)	0.4	3.4	18	21	57	144

(\*) Time abnormally high because the computer was not dedicated.

**Table 4.** Characteristics of Courier protocol model.

Table 4 shows the characteristics of the model. The column ‘Matrix Size’ contains the size of the matrix in megabytes if the matrix were to be kept entirely in memory. One can see that this transition-rate matrix is less dense than the one for the Kanban model. For this model, we wish to show how the solution process varies as the size of the problem gets larger. We set *MaxIter* to be 6 and let  $N$  range. Table 4 summarizes these results.

There are several interesting results from this study. First, we note that for  $N < 3$ , the file system buffers significantly decrease the conversion and solution times, so they should not be considered as part of a trend. More traditional techniques would probably do as well or better for such small models. For  $N \geq 3$ , the model becomes interesting. We wrote our own GSPN state generator for these models, and it was optimized for memory (so we could generate large models), not for speed. It was also designed to be compatible with the *UltraSAN* solvers and reward variable specification. The conversion time is the time it took to convert the  $Q$ -matrix in *UltraSAN*’s format to one used by the tool, which involves taking the transpose of  $Q$  and converting from an ASCII to binary floating point representation. The conversion time shown for  $N = 6$  is the wall clock time, but it is abnormally large since it was not run in a dedicated environment. We estimate from the user time that the conversion process would take about 7,000 seconds on a dedicated system.

The data gives us a rough idea about the relative performance of each step in the solution process. The conversion process takes between half and two thirds the generation time. We believe that much of the conversion time is spent translating an ASCII representation of a real number into the computer’s internal representation. The solution times are the times to reach the convergence criterion  $\|\Pi_i^{(k+1)} - \Pi_i^{(k)}\|_\infty < 10^{-12}$  described above, and are roughly twice the generation times. This shows that the solution process does not take a disproportionate amount of time more than the state generation or conversion process.

Another interesting observation of this model is that in the case where  $N = 6$ , the transition-rate matrix generated by the GSPN state generator (a sparse textual representation of the matrix) would be larger than 2 Gbytes, which is

	$N = 1$	$N = 2$	$N = 3$	$N = 4$	$N = 5$	$N = 6$
$\lambda$	74.3467	120.372	150.794	172.011	187.413	198.919
$P_{send}$	0.01011	0.01637	0.02051	0.02334	0.02549	0.02705
$P_{recv}$	0.98141	0.96991	0.96230	0.95700	0.95315	0.95027
$P_{sess1}$	0.00848	0.01372	0.01719	0.01961	0.02137	0.02268
$P_{sess2}$	0.92610	0.88029	0.84998	0.82883	0.81345	0.80197
$P_{transp1}$	0.78558	0.65285	0.56511	0.50392	0.45950	0.42632
$P_{transp2}$	0.78871	0.65790	0.57138	0.51084	0.46673	0.43365

**Table 5.** Reward variables for Courier protocol model.

larger than the maximum allowable file size on our workstation. To solve this system, we rounded the rates to 6 decimal places. This obviously affects the accuracy of the solutions. There are obvious and simple ways to use multiple files to avoid this problem; we simply state this observation to give the reader a feel for the size of the data that the tool is manipulating. Also, of the 144 Mbytes necessary to compute the solution, 118 Mbytes of it are needed just to hold the solution vector.

In Table 5 we show several of several of the reward variables in the model as  $N$  varies from 1 to 6. The  $\lambda$  we compute here corresponds to measuring  $\lambda_{lsp}$  in the model, which corresponds to the user’s message throughput rate. The measures  $\lambda_{frg}$  can easily be computed as  $\lambda_{frg} = \lambda q_1 / q_2$ . Similarly,  $\lambda_{ack} = \lambda_{lsp} + \lambda_{frg}$ . From this, we can see how the packet throughput rate ( $\lambda$ ) increases as the window size increases. Other reward variables are explained in [8], and they correspond to the fraction of time different parts of the system are busy. We note that the values we computed here differ from those Li found by approximate techniques [8, 14]. We suspect that Li used a fragmentation ratio in his approximation techniques that is different (and unpublished) from the ratio for which he gives “exact” solutions because we were able to reproduce the exact solutions.

## 5 Conclusion

We have described a new tool for solving Markov models with very large state spaces. By devising a method to efficiently store the state-transition-rate matrix on disk, overlap computation and data transfer on a standard workstation, and utilize an iterative solver that exhibits locality in its use of data, we are able to build a tool that requires little more memory than the solution vector itself to obtain a solution. This method is completely general to any model for which one can derive a state-transition-rate matrix. As illustrated in the paper, the tool can solve models with 10 million states and 100 million non-zero entries on a machine with only 128 Mbytes of main memory. Because we make use of an innovative implementation using two processes that communicate via shared memory, we are able to keep the compute process utilizing the CPU approximately 80% of the time.

In addition to describing the tool, we have illustrated its use on two large-scale models: a Kanban manufacturing system and the Courier protocol stack executing on a VME bus-based multiprocessor. For each model, we present detailed results concerning the time and space requirements for solutions so that our tool may be compared with existing and future tools. The results show that the speed of solution is much faster than those reported for implementations based on Kronecker operators. These results show that our approach is the current method of choice for solving large Markov models if sufficient disk space is available to hold the state-transition rate matrix.

## References

1. G. Ciardo, "Advances in compositional approaches based on Kronecker algebra: Application to the study of manufacturing systems," in *Third International Workshop on Performability Modeling of Computer and Communication Systems*, Bloomington, IL, Sept. 7-8, 1996.
2. G. Ciardo and M. Tilgner, "On the use of Kronecker operators for the solution of generalized stochastic Petri nets," ICASE Report #96-35 CR-198336, NASA Langley Research Center, May 1996.
3. D. D. Deavours and W. H. Sanders, "'On-the-fly' solution techniques for stochastic Petri nets and extensions," to appear in *Petri Nets and Performance Models*, 1997.
4. S. Donatelli, "Superposed generalized stochastic Petri nets: Definition and efficient solution," in R. Valette, editor, *Application and Theory of Petri Nets 1994, Lecture Notes in computer science 815 (Proc. 15th Int. Conf. on Application and Theory of Petri Nets, Zaragoza, Spain)*, pp. 258-277, Springer-Verlag, June 1994.
5. G. Horton, "Adaptive Relaxation for the Steady-State Analysis of Markov Chains," ICASE Report #94-55 NASA CR-194944, NASA Langley Research Center, June 1994.
6. P. Kemper, "Numerical analysis of superposed GSPNs," in *Proc. Int. Workshop on Petri Nets and Performance Models (PNPM'95)*, pp. 52-61, Durham, NC, Oct. 1995. IEEE Comp. Soc. Press.
7. P. Kemper, "Numerical Analysis of Superposed GSPNs," in *IEEE Transactions on Software Engineering*, 1996, to appear.
8. Y. Li, "Solution Techniques for Stochastic Petri Nets," Ph.D. Dissertation, Department of Systems and Computer Engineering, Carleton University, Ottawa, Ontario, May 1992.
9. J. F. Meyer, A. Movaghar, and W. H. Sanders, "Stochastic activity networks: Structure, behavior, and application," In *Proc. International Workshop on Timed Petri Nets*, pp. 106-115, Torino, Italy, July 1985.
10. A. Movaghar and J. F. Meyer, "Performability modeling with stochastic activity networks," In *Proc. 1984 Real-Time Systems Symp.*, Austin, TX, December 1984.
11. W. H. Sanders, W. D. Obal II, M. A. Qureshi, F. K. Widjanarko, "The *UltraSAN* modeling environment," in *Performance Evaluation*, pp. 89-115, Vol. 24, 1995.
12. W. J. Stewart, "Introduction to the Numerical Solution of Markov Chains," Princeton University Press, 1994.
13. A. S. Tanenbaum, *Modern Operating Systems*, Prentice Hall, 1992.

14. C. M. Woodside and Y. Li, "Performance Petri Net Analysis of Communications Protocol Software by Delay-Equivalent Aggregation," in *Proc. Fourth Int. Workshop on Petri Nets and Performance Models*, pp. 64-73, Melbourne, Australia, Dec. 2-5, 1991.