

# Aggregation Strategies for Large Semi-Markov Processes

Marcel C. Guenther      Nicholas J. Dingle  
Jeremy T. Bradley      William J. Knottenbelt

Department of Computing, Imperial College London  
180 Queen's Gate, London SW7 2BZ, United Kingdom  
Email: {mcg05,njd200,jb,wjk}@doc.ic.ac.uk

## Abstract

High-level semi-Markov modelling paradigms such as semi-Markov stochastic Petri nets and process algebras are used to capture realistic performance models of computer and communication systems but often have the drawback of generating huge underlying semi-Markov processes. Extraction of performance measures such as steady-state probabilities and passage-time distributions therefore relies on sparse matrix–vector operations involving very large transition matrices. The computational complexity of such operations depends on the number of rows and non-zeros in the matrix.

Previous studies have shown that exact state-by-state aggregation of semi-Markov processes can be applied to reduce the number of states. However, it is important to select the order in which states are aggregated judiciously so as to avoid a dramatic increase in matrix density caused by the creation of additional transitions between remaining states. Our paper addresses this issue by presenting the concept of state space partitioning for aggregation. Partitioning the state space entails dividing it into a number of non-intersecting subsets. In contrast to previous algorithms that perform state-by-state aggregation on the global unpartitioned graph, our new algorithm works on a local partition-by-partition basis which allows more space-efficient aggregation. We introduce different partitioning methods for this purpose. Furthermore we discuss partition sorting methods that determine the order in which partitions should be aggregated. This order has a significant impact on the connectivity of the aggregate state space, and thus the density of the transition matrix.

Aggregation of partitions can be done in one of two ways. The first way is to use exact state-by-state aggregation to aggregate each individual state within a partition. For this purpose, we present a new state ordering algorithm, which takes into account the exact number of new transitions that are created when aggregating a particular state. This technique is preferable to existing ordering methods which only approximate the number of newly created transitions. A second approach to the aggregation of partitions is *atomic partition aggregation*. Inspired by a technique used in passage-time analysis, this collapses a whole partition into a small number of semi-Markov states and transitions.

Most partitionings produced by existing graph partitioners are not suitable for use with our atomic partition aggregation techniques, and we therefore present a new deterministic partitioning method which we term *barrier partitioning*. We show that barrier partitioning is capable of splitting large semi-Markov models into two balanced partitions such that first passage-time analysis can be performed using 50% less memory than existing algorithms, without compromising speed.

# 1 Introduction

Semi-Markov processes (SMPs) are expressive tools for modelling a wide range of real-life systems. The state space explosion problem, however, hinders the analysis of large finite SMPs as it does of many stochastic and functional modelling disciplines. One approach to addressing this problem is to use aggregation techniques to remove single states or groups of states and aggregate their temporal effect into the remaining states. Many techniques exist in the Markovian domain for exact and approximate aggregation (e.g. lumpability [17], aggregation/disaggregation [12], aggregation of hierarchical models [11]) but to date analogous work on semi-Markov aggregation algorithms has been very limited. In prior work [5, 8], we presented an aggregation algorithm for semi-Markov processes which operates on each state individually. Our analysis in [8] suggests that the primary limitation of this technique is that the computational cost and memory requirements become very large as increasing numbers of states are aggregated and the transition matrices representing the SMP consequently gets less sparse.

In this paper, we present a number of novel approaches for overcoming this problem. Central to these is the concept of partitioning the state space, and we begin by considering different partitioning methods (initially inspired by those previously used for parallel sparse matrix–vector multiplication) and evaluating their suitability for our state-by-state aggregation algorithm. We demonstrate that by partitioning the state space in this way and then using the state-by-state aggregation algorithm on the separate partitions, as opposed to applying it directly to an unpartitioned state-space, we can reduce the computational cost and memory requirements of our exact aggregation approach. We also improve our exact state-by-state aggregation technique through the introduction of a new state-ordering metric which better accounts for the number of state transitions created as a result of performing the aggregation. This allows us to select states to be aggregated more effectively and hence reduce the maximum number of transitions created during aggregation.

Even this improved state-ordering metric does not resolve the central drawback of exact state-by-state aggregation: although the result of the process is an aggregated state space, the intermediate steps can actually create more state transitions (and hence require more storage and computational effort) than were present in the original unaggregated state space. Inspired by our prior work on iterative passage time analysis in SMPs [10], we therefore present *atomic partition aggregation* to overcome this limitation. This does not require each state in the partition to be aggregated in turn, but instead effectively calculates the passage time distribution across an entire partition and combines this with the state holding time distributions of relevant states outside the partition. As partitioning techniques suitable for parallel sparse matrix–vector multiplication do not produce partitions suitable for the application of atomic aggregation, we also introduce a new *barrier partitioning* strategy which is better suited. We demonstrate how this enables passage time analysis to be conducted using 50% less memory than before.

The remainder of this paper is organised as follows. Section 2 summarises background theory on the calculation of passage times in semi-Markov processes from [10], and also describes our previously-published state-by-state aggregation technique [8]. Section 3 then introduces the concept of performing aggregation on partitions of the state space, and discusses the importance of the order in which partitions are chosen to be aggregated. In Section 4 we propose a new state-selection metric for the exact state aggregation algorithm, and demonstrate how this better accounts for the number of new transitions created during aggregation. Section 5 then presents our novel atomic aggregation approach where whole partitions are aggregated by means of a passage-style analysis. Section 6 presents the barrier partitioning technique. Finally, Section 7 concludes and suggests directions for future work.

## 2 Background

### 2.1 Semi-Markov Processes

Semi-Markov Processes (SMPs) are an extension of Markov processes which allow for generally distributed sojourn times [19, 20]. Although the memoryless property no longer holds for state sojourn times, at transition instants SMPs still behave in the same way as Markov processes (that is to say, the choice of the next state is based only on the current state) and so share some of their analytical tractability.

Consider a Markov renewal process  $\{(\chi_n, T_n) : n \geq 0\}$  where  $T_n$  is the time of the  $n$ th transition ( $T_0 = 0$ ) and  $\chi_n \in \mathcal{S}$  is the state at the  $n$ th transition. Let the kernel of this process be:

$$R(n, i, j, t) = \mathbb{P}(\chi_{n+1} = j, T_{n+1} - T_n \leq t \mid \chi_n = i)$$

for  $i, j \in \mathcal{S}$ . The continuous time semi-Markov process,  $\{Z(t), t \geq 0\}$ , defined by the kernel  $R$ , is related to the Markov renewal process by:

$$Z(t) = \chi_{N(t)}$$

where  $N(t) = \max\{n : T_n \leq t\}$ , i.e. the number of state transitions that have taken place by time  $t$ . Thus  $Z(t)$  represents the state of the system at time  $t$ . We consider only time-homogeneous SMPs in which  $R(n, i, j, t)$  is independent of  $n$ :

$$\begin{aligned} R(i, j, t) &= \mathbb{P}(\chi_{n+1} = j, T_{n+1} - T_n \leq t \mid \chi_n = i) \quad \text{for any } n \geq 0 \\ &= p_{ij} H_{ij}(t) \end{aligned}$$

where  $p_{ij} = \mathbb{P}(\chi_{n+1} = j \mid \chi_n = i)$  is the state transition probability between states  $i$  and  $j$  and  $H_{ij}(t) = \mathbb{P}(T_{n+1} - T_n \leq t \mid \chi_{n+1} = j, \chi_n = i)$ , is the sojourn time distribution in state  $i$  when the next state is  $j$ . An SMP can therefore be characterised by two matrices  $\mathbf{P}$  and  $\mathbf{H}$  with elements  $p_{ij}$  and  $H_{ij}$  respectively.

### 2.2 Iterative Passage Time Algorithm

In this section we define the first passage-time random variable used throughout the paper. We also summarise from [10] an iterative algorithm for calculating first passage-time density in semi-Markov processes.

From now on, we consider a finite, irreducible, continuous-time semi-Markov process with  $N$  states  $\{1, 2, \dots, N\}$ . Recalling that  $Z(t)$  denotes the state of the SMP at time  $t$  ( $t \geq 0$ ) and that  $N(t)$  denotes the number of transitions which have occurred by time  $t$ , the first passage time from a source state  $i$  at time  $t$  into a non-empty set of target states  $\vec{j}$  is defined as:

$$P_{i\vec{j}}(t) = \inf\{u > 0 : Z(t+u) \in \vec{j}, N(t+u) > N(t), Z(t) = i\}$$

For a stationary time-homogeneous SMP,  $P_{i\vec{j}}(t)$  is independent of  $t$ :

$$P_{i\vec{j}} = \inf\{u > 0 : Z(u) \in \vec{j}, N(u) > 0, Z(0) = i\} \quad (1)$$

This formulation of the random variable  $P_{i\vec{j}}$  applies to an SMP with no immediate transitions. If such transitions are present, then the passage time can be stated as:

$$P_{i\vec{j}} = \inf\{u > 0 : N(u) \geq M_{i\vec{j}}\} \quad (2)$$

where  $M_{i\vec{j}} = \min\{m \in \mathbb{Z}^+ : \chi_m \in \vec{j} \mid \chi_0 = i\}$  is the transition marking the terminating state of the passage.

$P_{i\vec{j}}$  has an associated probability density function  $f_{i\vec{j}}(t)$ . The Laplace transform of  $f_{i\vec{j}}(t)$ ,  $L_{i\vec{j}}(s)$ , can be computed by means of a first-step analysis. That is, we consider moving from the source state  $i$  into the set of its immediate successors  $\vec{k}$  and must distinguish between those members of  $\vec{k}$  which are target states and those which are not. This calculation can be achieved by solving a set of  $N$  linear equations of the form:

$$L_{i\vec{j}}(s) = \sum_{k \notin \vec{j}} r_{ik}^*(s) L_{k\vec{j}}(s) + \sum_{k \in \vec{j}} r_{ik}^*(s) \quad : \text{ for } 1 \leq i \leq N \quad (3)$$

where  $r_{ik}^*(s)$  is the Laplace-Stieltjes transform (LST) of  $R(i, k, t)$  from Section 2.1 and is defined by:

$$r_{ik}^*(s) = \int_0^{\infty} e^{-st} dR(i, k, t) \quad (4)$$

Equation 3 has matrix-vector form  $\mathbf{Ax} = \mathbf{b}$ , where the elements of  $\mathbf{A}$  are general functions of the complex variable  $s$ . For example, when  $\vec{j} = \{1\}$ , Equation 3 yields:

$$\begin{pmatrix} 1 & -r_{12}^*(s) & \cdots & -r_{1N}^*(s) \\ 0 & 1 - r_{22}^*(s) & \cdots & -r_{2N}^*(s) \\ 0 & -r_{32}^*(s) & \cdots & -r_{3N}^*(s) \\ \vdots & \vdots & \ddots & \vdots \\ 0 & -r_{N2}^*(s) & \cdots & 1 - r_{NN}^*(s) \end{pmatrix} \begin{pmatrix} L_{1\vec{j}}(s) \\ L_{2\vec{j}}(s) \\ L_{3\vec{j}}(s) \\ \vdots \\ L_{N\vec{j}}(s) \end{pmatrix} = \begin{pmatrix} r_{11}^*(s) \\ r_{21}^*(s) \\ r_{31}^*(s) \\ \vdots \\ r_{N1}^*(s) \end{pmatrix} \quad (5)$$

We now describe an iterative algorithm for generating passage time densities that creates successively better approximations to the SMP passage time quantity  $P_{i\vec{j}}$  of Equation 1 [10]. We approximate  $P_{i\vec{j}}$  as  $P_{i\vec{j}}^{(r)}$ , for a sufficiently large value of  $r$ , which is the time for  $r$  consecutive transitions to occur starting from state  $i$  and ending in any of the states in  $\vec{j}$ . We calculate  $P_{i\vec{j}}^{(r)}$  by constructing and then numerically inverting [1, 2, 3] its Laplace transform  $L_{i\vec{j}}^{(r)}(s)$ .

Recall the semi-Markov process  $Z(t)$  of Section 2.1, where  $N(t)$  is the number of state transitions that have taken place by time  $t$ . We formally define the  $r$ th transition first passage time to be:

$$P_{i\vec{j}}^{(r)} = \inf\{u > 0 : Z(u) \in \vec{j}, 0 < N(u) \leq r, Z(0) = i\} \quad (6)$$

which is the time taken to enter a state in  $\vec{j}$  for the first time having started in state  $i$  at time 0 and having undergone up to  $r$  state transitions.

If we have immediate transitions in our SMP model (as in Equation 2) then the  $r$ th transition first passage time is:

$$P_{i\vec{j}}^{(r)} = \inf\{u > 0 : M_{i\vec{j}} \leq N(u) \leq r\}$$

This is because as the firing of an immediate transitions results in zero time being spent in the state in which it was enabled, it is not meaningful to talk about the SMP being in a particular state at a particular time. Instead, we count the transitions which have happened so that we may reason about the order in which they have occurred.

$P_{i\vec{j}}^{(r)}$  is a random variable with associated Laplace transform  $L_{i\vec{j}}^{(r)}(s)$ .  $L_{i\vec{j}}^{(r)}(s)$  is, in turn, the  $i$ th component of the vector:

$$\mathbf{L}_{\vec{j}}^{(r)}(s) = \left( L_{1\vec{j}}^{(r)}(s), L_{2\vec{j}}^{(r)}(s), \dots, L_{N\vec{j}}^{(r)}(s) \right)$$

representing the passage time for terminating in  $\vec{j}$  for each possible start state. This vector may be computed as:

$$\mathbf{L}_{\vec{j}}^{(r)}(s) = \mathbf{U} \left( \mathbf{I} + \mathbf{U}' + \mathbf{U}'^2 + \cdots + \mathbf{U}'^{(r-1)} \right) \mathbf{e}_{\vec{j}} \quad (7)$$

where  $\mathbf{U}$  is a matrix with elements  $u_{pq} = r_{pq}^*(s)$  and  $\mathbf{U}'$  is a modified version of  $\mathbf{U}$  with elements  $u'_{pq} = \delta_{p \notin \vec{j}} u_{pq}$ , where states in  $\vec{j}$  have been made absorbing. Here,  $\delta_{p \notin \vec{j}} = 1$  if  $p \notin \vec{j}$  and 0 otherwise. The initial multiplication with  $\mathbf{U}$  in Equation 7 is included so as to generate cycle times for cases such as  $L_{ii}^{(r)}(s)$  which would otherwise register as 0 if  $\mathbf{U}'$  were used instead. The column vector  $\mathbf{e}_{\vec{j}}$  has entries  $e_{k\vec{j}} = \delta_{k \in \vec{j}}$ , where  $\delta_{k \in \vec{j}} = 1$  if  $k$  is a target state ( $k \in \vec{j}$ ) and 0 otherwise.

From Equation 1 and Equation 6:

$$P_{i\vec{j}} = P_{i\vec{j}}^{(\infty)} \quad \text{and thus} \quad L_{i\vec{j}}(s) = L_{i\vec{j}}^{(\infty)}(s)$$

This can be generalised to multiple source states  $\vec{i}$  using, for example, a normalised steady-state vector  $\boldsymbol{\alpha}$  calculated from  $\boldsymbol{\pi}$ , the steady-state vector of the embedded discrete-time Markov chain (DTMC) with one-step transition probability matrix  $\mathbf{P} = [p_{ij}, 1 \leq i, j \leq N]$ , as:

$$\alpha_k = \begin{cases} \pi_k / \sum_{j \in \vec{i}} \pi_j & \text{if } k \in \vec{i} \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

The row vector with components  $\alpha_k$  is denoted by  $\boldsymbol{\alpha}$ . The formulation of  $L_{i\vec{j}}^{(r)}(s)$  is therefore:

$$\begin{aligned} L_{i\vec{j}}^{(r)}(s) &= \boldsymbol{\alpha} \mathbf{L}_{\vec{j}}^{(r)}(s) \\ &= (\boldsymbol{\alpha} \mathbf{U} + \boldsymbol{\alpha} \mathbf{U} \mathbf{U}' + \boldsymbol{\alpha} \mathbf{U} \mathbf{U}'^2 + \dots + \boldsymbol{\alpha} \mathbf{U} \mathbf{U}'^{(r-1)}) \mathbf{e}_{\vec{j}} \\ &= \sum_{k=0}^{r-1} \boldsymbol{\alpha} \mathbf{U} \mathbf{U}'^k \mathbf{e}_{\vec{j}} \end{aligned} \quad (9)$$

The sum of Equation 9 can be computed efficiently using sparse matrix–vector multiplications with a vector accumulator,  $\boldsymbol{\mu}_r = \sum_{k=0}^r \boldsymbol{\alpha} \mathbf{U}'^k$ . At each step, the accumulator (initialised as  $\boldsymbol{\mu}_0 = \boldsymbol{\alpha} \mathbf{U}$ ) is updated as  $\boldsymbol{\mu}_{r+1} = \boldsymbol{\alpha} \mathbf{U} + \boldsymbol{\mu}_r \mathbf{U}'$ .

In practice, convergence of the sum  $L_{i\vec{j}}^{(r)}(s) = \sum_{k=0}^{r-1} \boldsymbol{\alpha} \mathbf{U} \mathbf{U}'^k$  can be said to have occurred if, for a particular  $r$  and  $s$ -point:

$$|\operatorname{Re}(L_{i\vec{j}}^{(r+1)}(s) - L_{i\vec{j}}^{(r)}(s))| < \varepsilon \quad \text{and} \quad |\operatorname{Im}(L_{i\vec{j}}^{(r+1)}(s) - L_{i\vec{j}}^{(r)}(s))| < \varepsilon \quad (10)$$

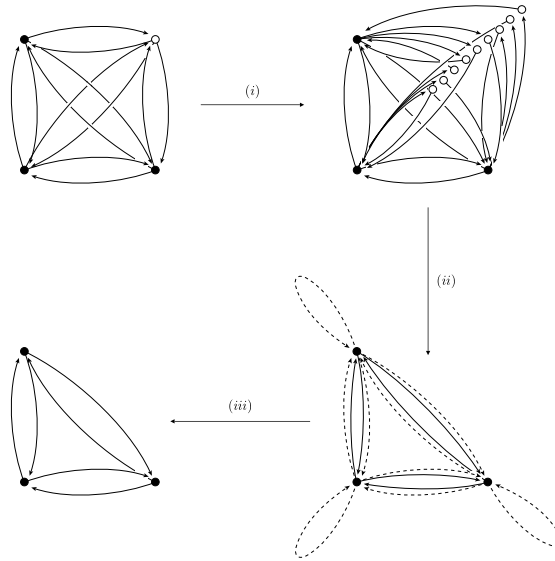
where  $\varepsilon$  is chosen to be a suitably small value, say  $\varepsilon = 10^{-16}$ .

### 2.3 Exact State Aggregation

In order to control the state space explosion which occurs when generating the state transition matrix for a semi-Markov process, we have previously developed an exact aggregation algorithm that acts on the semi-Markov state space directly [5, 8]. The aim is to apply the aggregation before performing any passage-time or transient analysis and thus reduce the calculation time required to solve the system of linear equations shown in Equation 5.

The method, illustrated in graphical terms in Figure 1, works as follows: first, a state is chosen to be aggregated. Then, from the transition graph, all paths of length two centred on that state are identified (step (i)) and aggregated into stochastically equivalent, single transitions (step (ii)). The newly-created transitions (shown dashed in Figure 1), which duplicate the route of existing transitions, are combined with the existing transitions. Finally, cyclic transitions are eliminated (step (iii)).

The result is to remove the chosen state and thus reduce the order of the transition matrix by one. Repeated application of this algorithm on different states will reduce the SMP to an arbitrary size



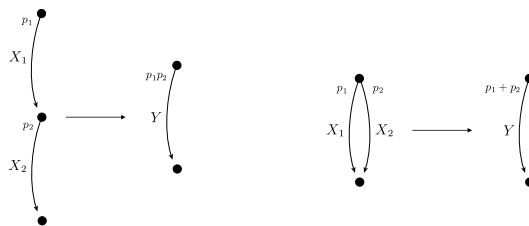
**Fig. 1.** Reducing a complete 4 state graph to a complete 3 state graph.

( $\geq 2$  states), while still preserving the exact passage-time distributions between all pairs of the remaining states. This style of aggregation is not possible in a Markovian context as aggregation operations of this type do not have a closed form in the Markov domain (i.e. the convolution of two Markovian delays is not itself Markovian).

There are three basic reduction steps for aggregating a single state of an SMP. These deal with convolutions, branching and cycles as follows:

### Sequential Reduction

In Figure 2(a),  $Y = X_1 + X_2$  is a convolution and therefore in Laplace form  $L_Y(s) = L_{X_1}(s)L_{X_2}(s)$ . In order to extract the path from an SMP we have to take into account the probabilities  $p_1$  and  $p_2$  of the first transition and second transitions of the path being selected. This gives us the overall path probability of  $p_1p_2$ .



(a) Sequential transitions.

(b) Branching transitions.

**Fig. 2.** Aggregating transitions in an SMP.

### Branch reduction

In Figure 2(b), we can sum the respective probabilities to get the overall selection probability for the aggregate path. Thus the aggregate probability for the branch is  $p_1 + p_2$ . Our

aggregate distribution,  $Y$ , is given by:

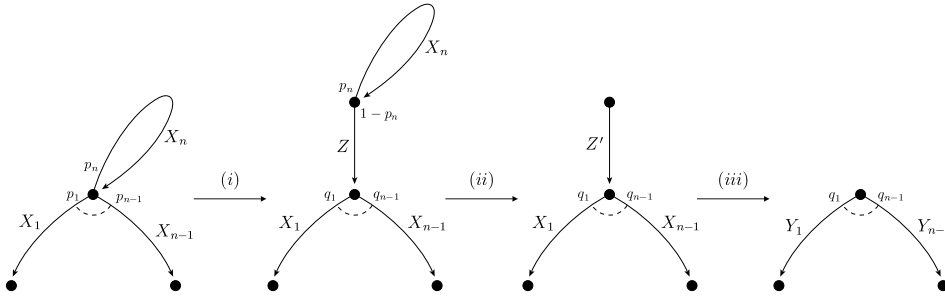
$$L_Y(s) = \frac{p_1}{p_1 + p_2} L_{X_1}(s) + \frac{p_2}{p_1 + p_2} L_{X_2}(s)$$

so that for both aggregate and unaggregated forms the total sojourn-time distribution has Laplace transform  $p_1 L_{X_1}(s) + p_2 L_{X_2}(s)$ .

### Cycle Reduction

When there is a state with at least one out-transition and a transition to itself, as shown in Figure 3, we can remove the cycle by making its stochastic effect part of the out-going transitions.

Consider a state transition system as being in the first stage of Figure 3, with  $(n - 1)$  out-transitions and probability  $p_i$  of departure along edge  $i$ . Each out-transition has an associated sojourn  $X_i$ ; the cycle probability is  $p_n$  with sojourn  $X_n$ .



**Fig. 3.** The three-step removal of a cycle from an SMP.

The first step, *(i)*, is to isolate the cycle and treat it separately from the branching out-transitions. We do this by rewriting the system to include an instantaneous delay and extra state immediately after the cycle,  $Z \sim \delta(0)$ ; the introduction of an extra state is only to aid our visualisation of the problem and is not necessary (or indeed performed) in the actual aggregation algorithm. Clearly the instantaneous transition will be selected with probability  $(1 - p_n)$ . We now have to renormalise the  $p_i$  probabilities on the branching state to become  $q_i = p_i / (1 - p_n)$ .

In step *(ii)* of Figure 3, we aggregate the delay of the cycle into the instantaneous transition creating a new transition with distribution  $Z'$ . By treating the system as a geometric sum of the random variable  $X_n$ , we can write:

$$L_{Z'}(s) = \frac{1 - p_n}{1 - p_n L_{X_n}(s)}$$

In stage *(iii)* of the process, the  $Z'$  delay can be sequentially convolved with the  $X_i$  sojourns to give us our final system.

In summary, we have reduced an  $n$ -out-transition state where one of the transitions was a cycle to an  $(n - 1)$ -out-transition state with no cycle such that:

$$q_i = \frac{p_i}{1 - p_n}$$

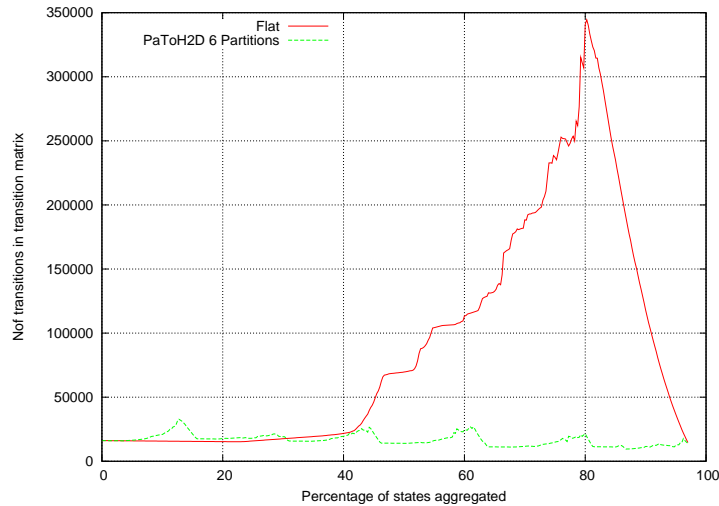
and:

$$L_{Y_i}(z) = \frac{1 - p_n}{1 - p_n L_{X_n}(z)} L_{X_i}(z)$$

## 2.4 Case Study Semi-Markov Models

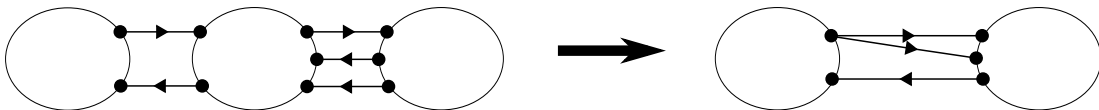
Throughout this paper we use three semi-Markov models as running examples. The Courier model [21] represents the ISO Application, Session and Transport layers of the Courier sliding-window communication protocol. The Voting model is a model of a distributed voting system with voters, failure-prone voting booths and failure-prone central servers [6, 9, 10]. The Web-server model represents a web content authoring system, and contains a number of clients, authors web servers and a write buffer [9, 10]. All three models were originally represented in a high-level Semi-Markov Stochastic Petri Net (SM-SPN) [7] form, from which semi-Markov processes of varying sizes can easily be generated. Further detail can be found in [15].

## 3 Partition Aggregation



**Fig. 4.** The effect of partition aggregation compared to flat aggregation of the 4 050 state Voting model.

Figure 4 illustrates the main problem encountered in applying the exact state-by-state aggregation algorithm outlined in Section 2.3 sequentially across the “flat” state space of an SMP with 4 050 states. The transition matrix initially contains approximately 15 000 non-zeros, but by the time that approximately 80% of the states have been aggregated (c. 3 200 states) the number of non-zeros in the transition matrix has increased to nearly 350 000 even though the dimensions of the matrix have been reduced dramatically. This is an important since it is the number of non-zeros that determine the storage requirements and run-time performance of our performance analysis algorithms, however, rather than just the dimensions of the matrix.



**Fig. 5.** Partition aggregation.

To avoid this dramatic peak in non-zeros, we propose partition aggregation. As shown in Figure 5, the state space of the SMP is divided into a number of partitions and the states within each of these are aggregated together, leaving only the transitions between the states on the boundaries

of each partition. The result of this can be seen in the lower curve in Figure 4; the peak in the number of non-zeros now occurs for each partition, but each peak is an order of magnitude smaller than the peak in non-zeros which occurs when aggregating the entire state space sequentially.

### 3.1 Partitioning Techniques

Central to this new aggregation technique is the ability to partition the SMP’s state space effectively. We divide  $n$  non-source and non-target states into  $k$  partitions, such that  $k|n$ . Inspired by our experiences in parallelising sparse matrix–vector multiplication, we consider the following three partitioning techniques:

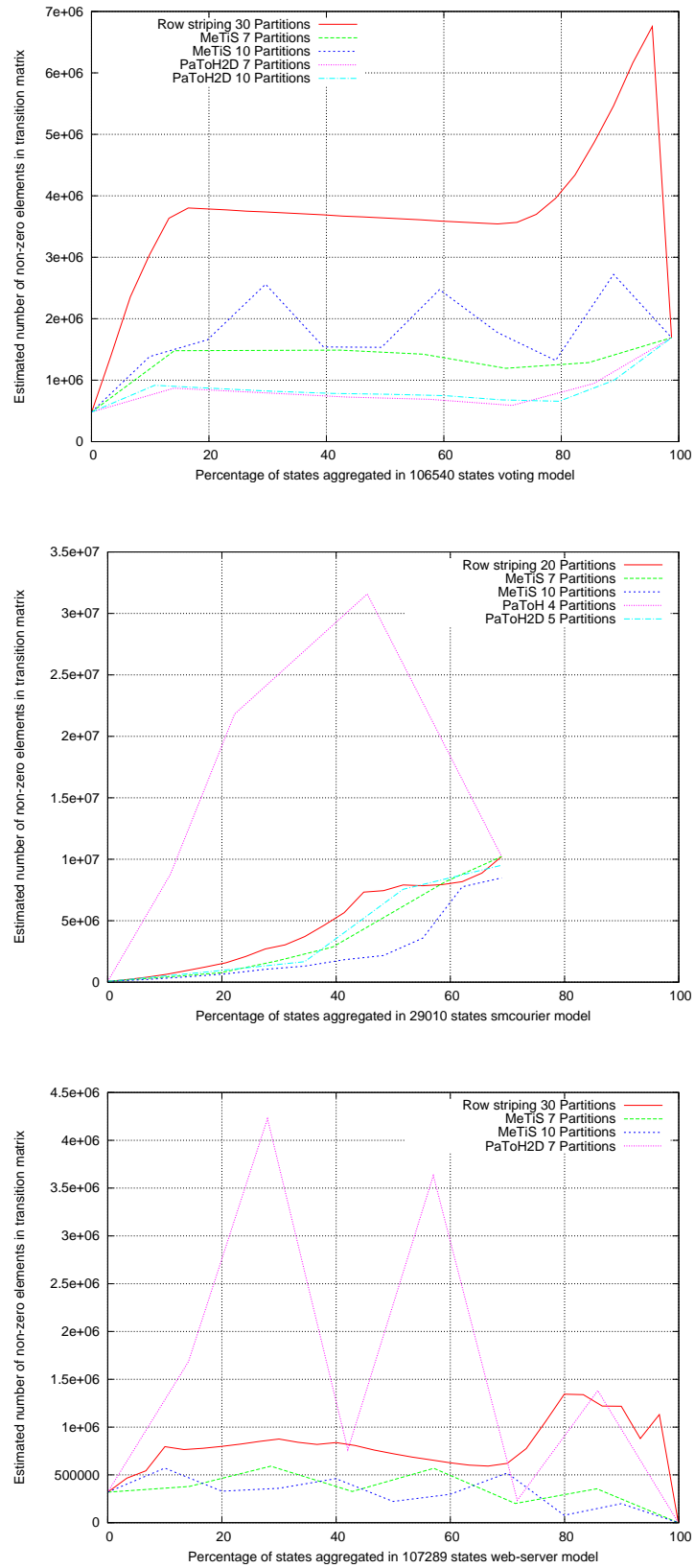
**Row striping** The simplest partitioning strategy is to divide the matrix into blocks of contiguous rows such that each block contains approximately the same number of non-zeros. For  $k$  partitions and  $n$  matrix rows, the first partition contains the first  $n/k$  matrix rows, the second is assigned the next  $n/k$  rows and so on. This scheme has the advantage of being very easy to compute and also of achieving good load balance.

**Graph partitioner** In a row-stripped decomposition, the the  $n \times n$  sparse transition matrix  $\mathbf{P}$  of an SMP can be represented as an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  where each row  $i$  ( $1 \leq i \leq n$ ) in the matrix corresponds to vertex  $v_i \in \mathcal{V}$  in the graph. The corresponding weight  $w_i$  of vertex  $v_i$  is the total number of non-zeros in row  $i$ . For the edge-set  $\mathcal{E}$ , edge  $e_{ij}$  connects vertices  $v_i$  and  $v_j$  with weight  $w_{ij} = 1$  if either one of  $p_{ij} > 0$  or  $p_{ji} > 0$ , and with weight  $w_{ij} = 2$  if both  $p_{ij} > 0$  and  $p_{ji} > 0$  [13]. Graph partitioners try to minimise the number of edges which span two partitions (these are said to be *cut*) while balancing the number of non-zero elements in each partition. We use the MeTiS sequential  $k$ -way graph partitioning library [16].

**Hypergraph partitioner** A hypergraph  $\mathcal{H} = (\mathcal{V}, \mathcal{N})$  is defined by a set of vertices  $\mathcal{V}$  and a set of nets (or hyperedges)  $\mathcal{N}$ , where each net is a subset of the vertex set  $\mathcal{V}$  [4]. A hypergraph is therefore a generalised graph data structure in which edges can connect arbitrary non-empty subsets of vertices. In the context of a row-wise decomposition of a sparse matrix, matrix row  $i$  ( $1 \leq i \leq n$ ) is represented by a vertex  $v_i \in \mathcal{V}$  while column  $j$  ( $1 \leq j \leq n$ ) is represented by net  $N_j \in \mathcal{N}$  [13]. The vertices contained within net  $N_j$  correspond to the row numbers of the non-zero elements within column  $j$ , i.e.  $v_i \in N_j$  if and only if  $p_{ij} \neq 0$ . Weights are assigned to vertices in the same manner as to the vertices of a graph. The weight of all nets is one, with a net’s contribution to the hyperedge cut being defined as one less than the number of different partitions spanned by that net. The overall objective of a hypergraph partitioning is to minimise the hyperedge cut while maintaining a balance criterion. In this paper we use the PaToH library [14] to perform hypergraph partitioning.

We distinguish between 1D hypergraph partitioning, where the hypernets either represent the successor states of each state (rows) or the predecessor states of each state (columns) and the 2D approach, where we use both successor and predecessor hypernets. Note that our definition of 2D hypergraph partitioning differs slightly from the definition commonly found in literature, where each non-zero matrix element becomes a vertex in the 2D hypergraph. In our case 2D simply implies that we use information from both rows and columns of the SMP transition matrix to construct hypernets.

We now investigate how the choice of the partitioner affects the number of non-zeros created in the transition matrix during exact state-by-state aggregation of partitions. Recall that our idea is to partition the state space of the SMP and run the exact state aggregation algorithm on each partition separately, and thus avoid the dramatic increase in non-zeros observed (and hence the amount of memory required) when aggregating the unpartitioned state space.



**Fig. 6.** The number of transitions in the transition matrix of different models during aggregation when partitioned with the three partitioners.

Figure 6 compares the number of non-zeros in the transition matrices of the three semi-Markov models when their state spaces are partitioned using these three techniques and then aggregated. We conclude that PaToH, which only uses the rows of the matrix as hypernets for partitioning, gives the worst results of all partitioners we tested as it leads to the largest number of non-zeros being created. For the Courier model (see Figure 6(b)), PaToH yields the worst matrix fill-in, while for the larger Voting and Web-server models, it either took too long to complete or exhausted the available memory on the test machine. The naïve row striping yielded good results in the Web-server and Courier model, but in the slightly more dense Voting model it performed much worse than MeTiS and PaToH2D. In general, we conclude that is very difficult to reliably use any one of these techniques to produce the best partitions; the choice of best partitioner varies depending on the model and the number of partitions required. This inspires our alternative barrier partitioning approach discussed in Section 6 below.

### 3.2 Partition Ordering

Our prior work on exact state aggregation [8] has shown the importance of choosing the order in which states should be aggregated carefully, and the same also applies to selecting the order in which partitions should be aggregated. Inspired by the state selection criteria in [8], we now compare two potential methods for partition order selection.

**Fewest-paths-first (FPF) partition sort** Suppose a partition has  $m$  predecessor states, i.e. states that lie outside the partition but have outgoing transitions to states in the partition, and  $n$  successor states, i.e. states that lie outside the partition and have incoming transitions from states in the partition. The number of transitions from the predecessor to the successor states in the SMP transition matrix after the aggregation of the partition is  $mn$  if all  $m$  predecessor states can reach all  $n$  successor states via paths through the partition. The FPF-value of a partition is:

$$mn - \textit{outgoing transitions}$$

where *outgoing transitions* is the total number of outgoing transitions from states in the partition. To choose a partition for aggregation using FPF sort we simply greedily select the one with the lowest FPF-value.

**Enhanced-fewest-paths-first (EFPPF) partition sort** Despite a being a good estimator for the total number of new transitions created after the aggregation of a partition, the FPF-value does not take into account the number of *incoming transitions* from the predecessor states of the partition. Further it does not count the *existing transitions* between the predecessor and successor states of the partition. The total number of new transitions after the aggregation can thus be estimated more accurately using *enhanced-fewest-paths-first (EFPPF) sort*. The EFPPF-value is:

$$mn - \textit{outgoing transitions} - \textit{incoming transitions} - \textit{existing transitions}$$

Note that the EFPPF-value of a partition is only an upper bound for the total number of new transitions in the transition matrix after the aggregation of a partition. This is because there may not be a path from every predecessor state to every successor states with all intermediate states of the path being partition internal states. Even for small values of  $m$  and  $n$  this may cause significant differences between the estimated and the actual number of partitionwise transitions.

Even though it is more expensive to calculate, our experiments have shown that EFPPF partition sort usually gives better results than FPF or picking the partitions in a random order. Figure 7 shows one situation where this is the case, specifically for a 5-way partitioning of the 10 300 state Voting model. For this reason we confine ourselves to considering only EFPPF partition sorting in the following sections.

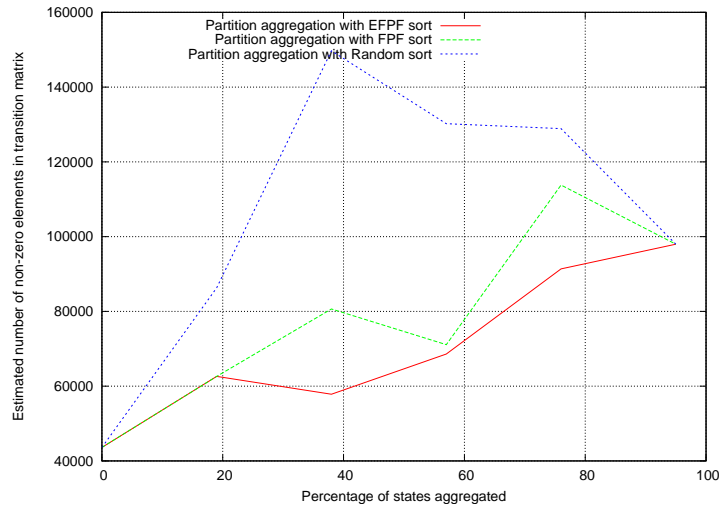


Fig. 7. Comparing EFPF partition sort with FPF partition sort.

## 4 Improved Exact State Aggregation

In this section we present an improvement to the exact state-by-state aggregation technique described in Section 2.3 and [8] to aggregate partitions of states. In [8] various state sorting techniques were introduced and tested; we describe the most successful here (Fewest-Paths-First aggregation) before presenting a new state ordering method (which we term Exact-Fewest-Paths-First).

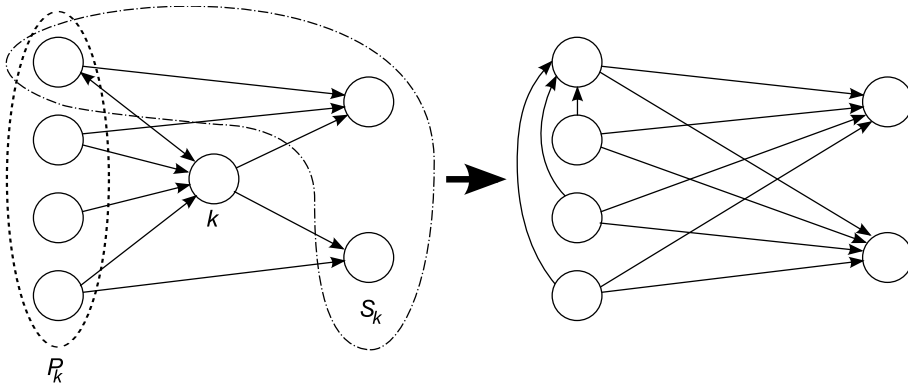
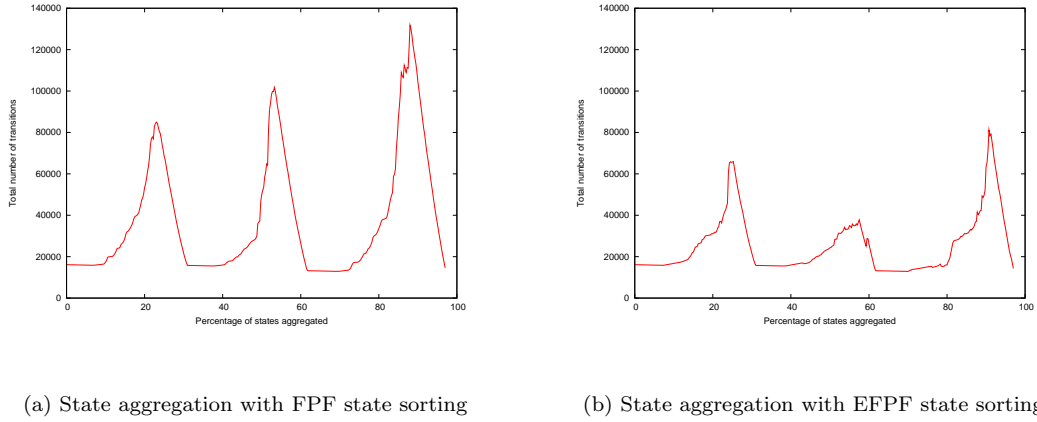


Fig. 8. The transition structure of an SMP before (left) and after (right) the aggregation of state  $k$ , where  $P_k$  is the set of predecessor states of state  $k$  and  $S_k$  the set of successor states of state  $k$ . Note that the self cycle of the state that lies in both  $P_k$  and  $S_k$  has been removed after the aggregation of state  $k$ .

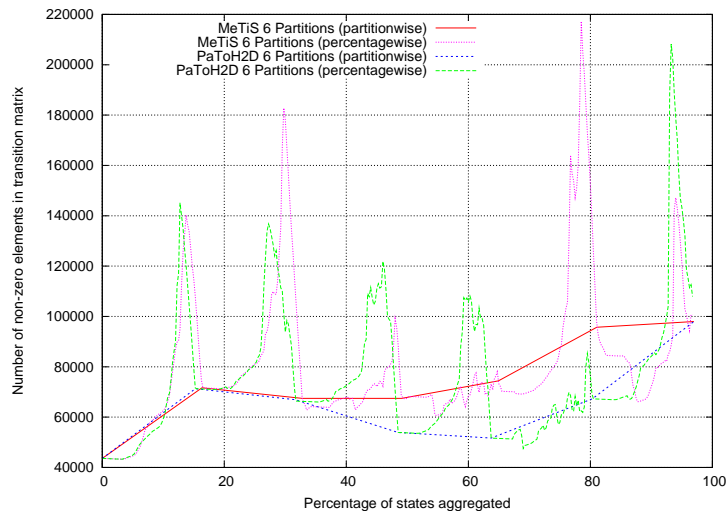
**Fewest-Paths-First (FPF) state aggregation** Given each state has  $m$  predecessor states and  $n$  states, FPF chooses the state with lowest  $mn$ -value first. This is designed to minimise computation, as  $O(mn)$  convolution and branching aggregations are required to eliminate a state. The downside of FPF is that it does not take into account existing transitions between predecessor and successor states and so can overestimate the number of transitions produced by aggregating a state. Figure 8 illustrates this problem with reference to the aggregation of a state  $k$ . The FPF algorithm of [8] would calculate a cost of  $4 \cdot 3 = 12$  for aggregating state  $k$ , while the actual number of newly created transitions is only  $4 \cdot 3 - 4 - 3 - 4 = 1$ .



**Fig. 9.** Voting model with 4050 states and 3 partitions. Partitions were sorted using EFPF partition sort.

**Exact-Fewest-Paths-First (EFPF) state aggregation** To overcome the inaccuracy of the FPF metric, we present *Exact-Fewest-Paths-First (EFPF) aggregation*. Suppose a state  $k$  has  $m$  predecessors,  $n$  successors and  $i \in \{0, 1\}$  self-loops. Moreover assume that there are  $t$  existing transitions between the successor and predecessor states, not including the transitions starting or ending in state  $k$ . The latter restriction is important as a state with a self-loop is its own predecessor and successor state. The EFPF-value of state  $k$  is  $(m - i)(n - i) - m - n - t$ . Note that we do not count self-loops, which are created when the set of predecessor states intersects with the set of successor states, as new transitions. This is because all these loops can be removed after each aggregation. Figure 9 shows example results where EFPF aggregation outperforms the FPF technique and leads to fewer extra non-zeros being created.

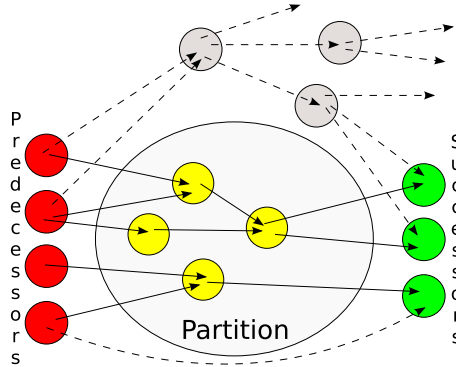
## 5 Atomic Partition Aggregation



**Fig. 10.** State-by-state partition aggregation on 10300 states Voting model.

Compared to flat state-by-state aggregation, the partition-by-partition aggregation approach re-

duces the transition matrix fill-in drastically. However, there is still the problem that the maximum number of transitions generated during the aggregation of a partition is much higher than the final number of transitions in the aggregated state space (see Figure 10). Indeed, there is also the problem that the final number of non-zeros in the aggregated state space can be higher than in the initial unaggregated one. Such density peaks are undesirable because it requires a significant amount of memory to store all temporary transitions, and the fill-in also slows down the aggregation of states as we need to perform more sequential and branching aggregation operations to remove states when the sub-matrix of a partition becomes dense. This observation prompted us to investigate an approach inspired by first passage time analysis which avoids these peaks by aggregating an entire partition in one go. We term this *atomic aggregation*.



**Fig. 11.** Atomic aggregation.

The general concept is illustrated in Figure 11. First we compute the passage time from each predecessor state  $p$  to every successor state  $s$  including only paths whose intermediate states lie entirely in the partition (denoted by the solid arcs in Figure 11). In a second step we aggregate the passage time and the probability of these internal transitions with the passage time and probability of the existing one-step transition from  $p$  to  $s$  (denoted by the lower dashed arc in Figure 11), if such a transition exists, using the branch aggregation technique from Section 2.3. If this one-step transition from  $p$  to  $s$  does not exist then the transition we computed in the first step becomes the new transition from  $p$  to  $s$ .

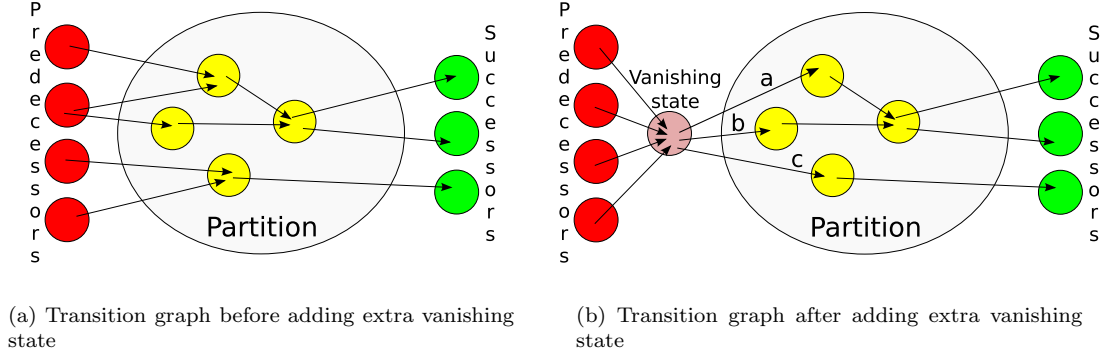
We only consider outgoing transitions from predecessor states of the partition to partition internal states. All other outgoing transitions of the predecessor states are ignored. We do not normalise the transition probabilities of outgoing transitions from the predecessor states as the sum of probabilities of the transitions of a predecessor state to each of the successor states after the aggregation of the partition is the same as the sum of probabilities of the transitions from the predecessor state to the partition internal states and successor states before aggregation. This can be formally justified by the flow conservation law, as we ensure that there are no final strongly-connected components of states within the partition.

Even though this appears to be a good strategy for aggregating an entire partition at once, it has one major disadvantage. Assume a partition has  $m$  predecessor,  $n$  successor and  $i$  partition internal states. In order to calculate the transition from every predecessor to every successor state using partition internal paths only, we have to solve  $m$  sets of  $i + n$  equations.

## 5.1 Modified Atomic Aggregation

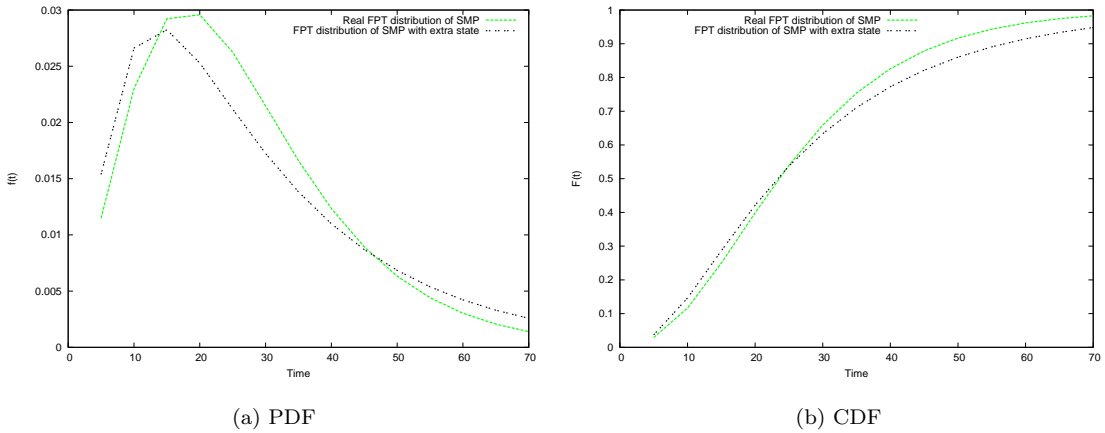
The main problem with atomic aggregation is that the number of linear equations to be solved to aggregate a partition depends on the number of predecessor and successor states of that partition, and that it may not be possible to find a partition of an SMP's state space that keeps the number

of such states low. To overcome this, we investigate inserting extra states into the SMP to try to ensure that partitions have only one predecessor or successor state. Adding extra states was inspired by the application of hidden nodes in Bayesian inference [18].



**Fig. 12.** Insertion of an extra vanishing state to improve atomic aggregation.

The general approach is shown in Figure 12. Through the extra state all four predecessor states have become connected to all partition entry states and can thereby reach each of the successor states of the partition. The number of linear equations required to aggregate the partition is therefore lower, but we have changed the structure of the SMP and so introduced error into any performance measures calculated upon it.



**Fig. 13.** Effect on the first passage time density and distribution of adding an extra state to the Courier model with 29 010 states.

To illustrate the error in the first-passage time distribution introduced by adding an extra state to the transition matrix, we compare the results from the unmodified model with results from the same model with an extra predecessor state. Figure 13 shows the resulting nature of the approximation to the first-passage time distribution of the original SMP when analysing the modified graph. The Kolmogorov-Smirnov statistic for the two distributions (the maximum absolute difference between the two) is 0.0573 (4 d.p.), but nevertheless the resulting pdf and cdf appear to be good approximations to the real passage time density and distribution respectively. In a second experiment we tested the impact of adding an extra predecessor state in the 107 289 state Web-server model. In this example we achieved a better approximation with a Kolmogorov-Smirnov

statistic for the two distributions of 0.0002 (4 d.p.).

Note, however, that the runtime of the passage time analyser in both cases was twice as long for the model with the added state as for the unmodified SMP. It was possible, however, sometimes to achieve a speed-up. The algorithm was tested on a Intel Duo Core 1.8 Ghz processor with 1Gbyte RAM. For the 106 540 state Voting model the total time taken to do atomic aggregation and the subsequent passage time analysis for 165 Laplace transform samples with convergence precision  $10^{-16}$  was 306 seconds. The total number of complex number multiplications was 2 553 489 711. In contrast, it took 398 seconds and 3 709 928 347 complex number multiplications to do the same passage time calculation on the initial SMP graph without aggregation.

## 6 Barrier Partitioning

Atomic aggregation requires us to find partitions that have a low number of predecessor or successor states. As partitioners such as PaToH and MeTiS are not guaranteed to find such partitions, we need to investigate further partitioning methods for transition graphs of large semi-Markov models. In this section we introduce a new partitioning method called *barrier partitioning*.

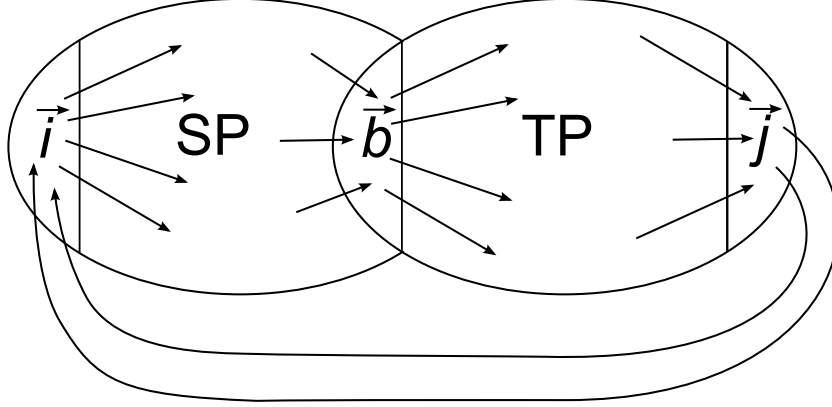
To perform first-passage time analysis on a SMP with  $n$  states we need to solve  $n$  linear equations to obtain  $L_{i\vec{j}}(s)$  (see Section 2.2). We observe that first-passage time analysis can be done forward, i.e. from each source state to the set of target states, as well as in reverse, i.e. from the set of target states to the individual source states, by transposing the SMP transition matrix and swapping source and target states. Such reverse passage time calculation works well in Laplace space since complex multiplication is an associative operation. The barrier partitioning method exploits this duality between the forward and reverse calculation of the first-passage time distribution and allows us to split the first-passage time calculation into two separate calculations. The combined cost of doing the two separate calculations is the same as the cost of the original first-passage time calculation, but with the advantage that each of the two separate calculations requires only half the amount of memory as the complete original.

**Definition 1.** Assume we have an SMP with a set of start states  $\vec{i}$  and a set of target states  $\vec{j}$ . If any state is a source and a target state at the same time it can be split up into a target and source state, by adding an immediate transition from the new target to the new source state, without changing any measures of the SMP model represented by the new graph. We divide the state space into two partitions  $SP$  and  $TP$ .  $SP$  contains all source states and a proportion of the intermediate states such that any outgoing transitions from  $SP$  to  $TP$  go into a set of barrier states  $\vec{b}$  in  $TP$ . Furthermore the only outgoing transitions from states in  $TP$  to states in  $SP$  are from target states  $\vec{j}$  to source states  $\vec{i}$ . Thus once a path has entered  $TP$  it can only ever go back to  $SP$  by going through states  $\vec{j}$ . Note that  $\vec{b}$  and  $\vec{j}$  may intersect. The resulting partitioning is a *barrier partitioning*. See Figure 14 for a graphic representation.

**Proposition 1.** Assume that we can divide the state space  $\mathcal{S}$  of a connected SMP graph into two partitions such that the resulting partitioning is a barrier partitioning. Clearly we have  $\vec{i} \cap \vec{j} = \emptyset$ ,  $SP \cup TP = \mathcal{S}$ . We denote the set of source states as  $\vec{i}$ , the set of barrier states as  $\vec{b}$  and the set of target states as  $\vec{j}$ . The result of first-passage time calculation from a source state  $i$  to the set of target states  $\vec{j}$  is same as the result obtained by doing a first-passage time calculation from  $i$  to the set of barrier states  $\vec{b}$ , convolved with the first-passage time calculation from the set of barrier states  $\vec{b}$  to the set of target states  $\vec{j}$ . In the Laplace domain this translates to:

$$L_{i\vec{j}}(s) = \sum_{b \in \vec{b}} L_{ib}^R(s) L_{b\vec{j}}(s)$$

where  $L_{ib}^R(s)$  denotes a restricted first-passage time distribution from state  $i$  to state  $b \in \vec{b}$ , where all states in  $\vec{b}$  are made absorbing for the calculation of  $L_{ib}^R(s)$ . This ensures that we only consider



**Fig. 14.** Barrier partitioning.

paths of the form  $i \rightarrow k_1 \rightarrow \dots \rightarrow k_m \rightarrow b$ , with  $k_x \in SP$ . In other words we do not consider paths through TP for the calculation of  $L_{ib}^R(s)$ .

*Proof.* By Equation 3 we have:

$$L_{i\vec{j}}(s) = \sum_{k \in (SP \cup TP) \setminus \vec{j}} r_{ik}^*(s) L_{k\vec{j}}(s) + \sum_{k \in \vec{j}} r_{ik}^*(s)$$

hence:

$$L_{i\vec{j}}(s) = \sum_{k \in (SP \cup TP)} r_{ik}^*(s) L_{k\vec{j}}(s)$$

where  $L_{k\vec{j}}(s)$  is equal to 1 if  $k \in \vec{j} \cap \vec{b}$ . We can rewrite  $k \in SP \cup TP$  since  $k \in SP \cup \vec{b}$  as there is no transition from any state in  $SP$  to any state in  $TP \setminus \vec{b}$  by construction of the barrier.

$$\begin{aligned} L_{i\vec{j}}(s) &= \sum_{k \in (SP \cup \vec{b})} r_{ik}^*(s) L_{k\vec{j}}(s) \\ &= \sum_{b \in \vec{b}} r_{ib}^*(s) L_{b\vec{j}}(s) + \sum_{k \in SP} r_{ik}^*(s) L_{k\vec{j}}(s) \end{aligned}$$

also by construction of the barrier partitioning and the fact that target states are absorbing states we know that once we have entered  $TP$  (i.e. reached a state in  $\vec{b}$ ) we cannot find a path back to a state in  $SP$ . Hence:

$$\begin{aligned} L_{i\vec{j}}(s) &= \sum_{b \in \vec{b}} r_{ib}^*(s) L_{b\vec{j}}(s) + \sum_{k \in SP} r_{ik}^*(s) \sum_{b \in \vec{b}} L_{kb}^R(s) L_{b\vec{j}}(s) \\ &= \sum_{b \in \vec{b}} r_{ib}^*(s) L_{b\vec{j}}(s) + \sum_{b \in \vec{b}} \sum_{k \in SP} r_{ik}^*(s) L_{kb}^R(s) L_{b\vec{j}}(s) \\ &= \sum_{b \in \vec{b}} \left( r_{ib}^*(s) L_{b\vec{j}}(s) + \sum_{k \in SP} r_{ik}^*(s) L_{kb}^R(s) L_{b\vec{j}}(s) \right) \\ &= \sum_{b \in \vec{b}} \left[ \left( \sum_{k \in SP} r_{ik}^*(s) L_{kb}^R(s) + r_{ib}^*(s) \right) L_{b\vec{j}}(s) \right] \end{aligned}$$

by definition  $\sum_{k \in SP} r_{ik}^*(s) L_{kb}^R(s) + r_{ib}^*(s)$  is the restricted first-passage time from state  $i$  to barrier state  $b$ . Therefore:

$$L_{i\vec{j}} = \sum_{b \in \vec{b}} L_{ib}^R(s) L_{b\vec{j}}(s)$$

■

**Corollary 1.1.** Let  $L_{\vec{i}\vec{b}}^R(s) = \{L_{\vec{i}b_1}^R(s), \dots, L_{\vec{i}b_l}^R(s)\}$ , where  $L_{\vec{i}b_m}^R(s) = \{\alpha_1 L_{i_1 b_m}^R(s) + \dots + \alpha_l L_{i_l b_m}^R(s)\}$  and  $L_{\vec{b}\vec{j}}(s) = \{L_{b_1 \vec{j}}(s), \dots, L_{b_l \vec{j}}(s)\}$  then in steady-state we have  $L_{\vec{i}\vec{j}}(s) = \sum_{b \in \vec{b}} L_{\vec{i}\vec{b}}^R(s) L_{b\vec{j}}(s) = L_{\vec{i}\vec{b}}^R(s) \cdot L_{\vec{b}\vec{j}}(s)$

*Proof.* Let  $\alpha_1, \alpha_2, \dots, \alpha_l$  be the normalised steady-state probabilities of the source states  $\vec{i} = (i_1, i_2, \dots, i_l)$  as defined in Equation 8. By Equation 9 we have:

$$\begin{aligned} L_{\vec{i}\vec{b}} &= \alpha_1 L_{i_1 \vec{j}}(s) + \alpha_2 L_{i_2 \vec{j}}(s) + \dots + \alpha_l L_{i_l \vec{j}}(s) \\ &= \sum_{b \in \vec{b}} \left( \alpha_1 \left( L_{i_1 b}^R(s) L_{b\vec{j}}(s) \right) + \dots + \alpha_l \left( L_{i_l b}^R(s) L_{b\vec{j}}(s) \right) \right) \\ &= \sum_{b \in \vec{b}} \left( \alpha_1 L_{i_1 b}^R(s) + \dots + \alpha_l L_{i_l b}^R(s) \right) L_{b\vec{j}}(s) \end{aligned}$$

■

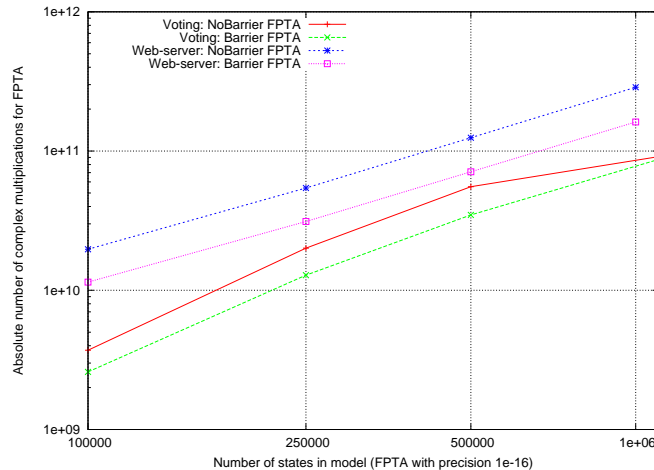
## 6.1 Barrier Partitioning In Practice

To compute the first-passage time distribution of a model whose state space has been split into partitions  $SP$  and  $TP$ , we start by calculating  $L_{\vec{i}\vec{b}}(s)$  using iterative first-passage time calculation. For this the source states remain unmodified, but the barrier states become absorbing target states. Also as this calculation is part of the final first-passage time calculation we need to weight the source states by their normalised steady state probabilities. Having calculated  $L_{\vec{i}\vec{b}}(s)$  we use it as our  $\mu_0$  (see Section 2.2) in the subsequent first-passage time calculation from the set of barrier states  $\vec{b}$  to the set of target states  $\vec{j}$ .

This technique reduces the amount of memory that we need for a first-passage time calculation as we only have to keep either the sub-matrix of the source partition or the target partition in memory at any point in time. Another advantage of barrier partitionings is that we can easily find barrier partitions in large models at low cost. Firstly, since we are doing first-passage time analysis we can discard the outgoing transitions from all target states. Secondly, we explore the entire state space using breadth-first search, with all source states being at the root level of the search. We store the resulting order in an array. To find a barrier partitioning we first add all non-target states among the first  $m$  states in the array to our source partition. Note that  $m$  has to be larger than the number of source states in the SMP. We then create a list of all predecessor states of the resulting partition. In the next step we add all predecessor states in the list to the source partition and recompute the list of predecessor states. We repeat this until we have found a source partition with no predecessor states. Since we discarded all outgoing edges of the target states, this method must give us a barrier partitioning. In the worst case this partitioning has all source and intermediate states in  $SP$  and  $TP$  only contains the set of target states.

In both the Voting and the Web-server model it is possible to split the state space such that each partition contains roughly 50% of the total number of transitions. Even more surprisingly, we

easily found balanced partitions (those where  $SP$  and  $TP$  contain a similar number of transitions) for large versions of the Voting and Web-server models with several million transitions. In addition our barrier partitioning algorithm is very fast. The computation of a balanced barrier partitioning for the 1.1 million state Voting model takes less than 10 seconds on an Intel Duo Core machine with two 1.8 Ghz processors and 1 Gbyte of RAM. The computation of a 2-way partitioning with PaToH2D takes about 60 seconds on the same machine, but the resulting partitioning is not suitable for atomic aggregation as both partitions have large numbers of predecessor and successor states.



**Fig. 15.** Log-log comparison of the absolute number of multiplications required under different aggregation strategies in the Voting and Web-server models.

The log-log plot in Figure 15 compares the number of complex multiplications needed for our different aggregation methods to calculate 165 Laplace transform samples for the estimation of the mode of the distribution and samples for 4 other  $t$ -points of the distribution. It is interesting to observe that the Barrier method generally seems to require fewer complex multiplications than the NoBarrier method in both models. However, the steep increase in the number of complex multiplications needed by the Barrier method between the 500 000 and the 1 100 000 state Voting model might highlight a trend that the Barrier method needs more complex multiplications than the NoBarrier method in large SMPs. Further investigation on larger models is necessary to see if this is a general trend or if it is simply due to the nature of the Voting model.

## 7 Conclusion

In this paper we have presented a number of improved aggregation techniques for SMPs. We have shown how dividing an SMP's state space into a number of loosely-connected partitions reduces the maximum number of transitions generated during the application of our state-by-state exact aggregation algorithm, and devised two partition-ordering metrics (analogous to the state-ordering metrics of the exact aggregation algorithm) to determine the order in which partitions should be aggregated. Of these, we concluded that our EFPF method gave better results than the FPF method. We also evaluated a similar EFPF metric to decide the order in which states should be aggregated within a partition, and again concluded that it produced better results than the previously-employed FPF metric.

Even with the partition aggregation approach with improved state and partitioning ordering metrics, however, we could not complete escape the fact that many additional temporary transitions were being created during aggregation. This inspired us to propose a scheme, based on first passage time analysis, for the atomic aggregation of partitions. Provided we find a suitable partition, atomic partition aggregation is more efficient than state-by-state aggregation of partitions. Like state-by-state aggregation, it may not always yield a speed up in computation time of the passage time analysis, but it can always be used to save memory as we only need to store the sub-matrix of the partition under consideration.

The biggest problem with atomic barrier partitioning is that we may not be able to find suitable partitions using existing state-space partitioning techniques. Introducing additional vanishing states alleviates this somewhat, but results in errors being introduced into the final calculated first passage time distributions. We therefore developed barrier partitioning, which deterministically partitions the SMP's state space into two partitions and thus allows first passage time analysis to be conducted using 50% of the memory required for the unaggregated SMP.

For the future, it would be interesting to investigate if graph and hypergraph partitioners can be modified to produce better partitionings for atomic aggregation. This could potentially be done by finding more suitable configurations for the PaToH and MeTiS partitioner. However, it is likely that there are better algorithms and further research might produce partitioning strategies that extend the range of semi-Markov models for which atomic aggregation can be used.

It is also possible that there are ways of channelling the outgoing transitions of predecessor states through an extra vanishing state that keep the error term lower than our algorithm does. One way of doing this might be to introduce more than one extra vanishing state. This would allow us to refine the connectivity of the graph with the extra states to reflect the original structure of the network more accurately than a graph with only one extra vanishing state does. Finally, it would be interesting to investigate the extent to which atomic aggregation can be applied to transient and steady-state probability analysis.

## References

- [1] J. Abate, G.L. Choudhury, and W. Whitt. On the Laguerre method for numerically inverting Laplace transforms. *INFORMS Journal on Computing*, 8(4):413–427, 1996.
- [2] J. Abate and W. Whitt. The Fourier-series method for inverting transforms of probability distributions. *Queueing Systems*, 10(1):5–88, 1992.
- [3] J. Abate and W. Whitt. Numerical inversion of Laplace transforms of probability distributions. *ORSA Journal on Computing*, 7(1):36–43, 1995.
- [4] C. Berge. *Hypergraphs: Combinatorics of Finite Sets*. North-Holland, Amsterdam, 1989.
- [5] J.T. Bradley. A passage-time preserving equivalence for semi-Markov processes. In *Lecture Notes in Computer Science 2324: Proceedings of the 12th International Conference on Modelling, Techniques and Tools (TOOLS'02)*, pages 178–187, London, April 14th–17th 2002. Springer-Verlag.
- [6] J.T. Bradley, N.J. Dingle, P.G. Harrison, and W.J. Knottenbelt. Distributed computation of passage time quantiles and transient state distributions in large Semi-Markov models. In *Proceedings of the International Workshop on Performance Modeling, Evaluation and Optimization of Parallel and Distributed Systems (PMEO-PDS'03)*, Nice, April 26th 2003.
- [7] J.T. Bradley, N.J. Dingle, P.G. Harrison, and W.J. Knottenbelt. Performance queries on semi-Markov stochastic Petri nets with an extended Continuous Stochastic Logic. In *Proceedings of 10th International Workshop on Petri Nets and Performance Models (PNPM'03)*, pages 62–71, Urbana-Champaign IL, USA, September 2nd–5th 2003.

- [8] J.T. Bradley, N.J. Dingle, and W.J. Knottenbelt. Strategies for exact iterative aggregation of semi-Markov performance models. In *Proceedings of International Symposium on Performance Evaluation of Computer and Telecommunication Systems (SPECTS'03)*, pages 755–762, Montreal, Canada, July 20th–24th 2003.
- [9] J.T. Bradley, N.J. Dingle, W.J. Knottenbelt, and H.J. Wilson. Hypergraph-based parallel computation of passage time densities in large semi-Markov models. In *Proceedings of the 4th International Conference on the Numerical Solution of Markov Chains (NSMC'03)*, pages 99–120, Urbana-Champaign IL, USA, September 2nd–5th 2003.
- [10] J.T. Bradley, N.J. Dingle, W.J. Knottenbelt, and H.J. Wilson. Hypergraph-based parallel computation of passage time densities in large semi-Markov models. *Linear Algebra and its Applications*, 386:311–334, 2004.
- [11] P. Buchholz. Hierarchical Markovian models: Symmetries and aggregation. *Performance Evaluation*, 22:93–110, 1995.
- [12] W-L. Cao and W.J. Stewart. Iterative aggregation/disaggregation techniques for nearly uncoupled Markov chains. *Journal of the ACM*, 32(3):702–719, July 1985.
- [13] U.V. Catalyürek and C. Aykanat. Hypergraph-partitioning-based decomposition for parallel sparse-matrix vector multiplication. *IEEE Transactions on Parallel and Distributed Systems*, 10(7):673–693, July 1999.
- [14] U.V. Catalyürek and C. Aykanat. PaToH: A multilevel hypergraph partitioning tool. Technical Report BU-CE-9915, Version 3.0, Department of Computer Engineering, Bilkent University, Ankara, 06800, Turkey, 1999.
- [15] N.J. Dingle. *Parallel Computation of Response Time Densities and Quantiles in Large Markov and Semi-Markov Models*. PhD thesis, Imperial College, London, United Kingdom, 2004. To appear.
- [16] G. Karypis and V. Kumar. *METIS: A Software Package for Partitioning Unstructured Graphs, Partitioning Meshes, and Computing Fill-Reducing Orderings of Sparse Matrices, Version 4.0*. University of Minnesota, September 1998.
- [17] J.G. Kemeny and J.L. Snell. *Finite Markov Chains*. Van Nostrand, 1960.
- [18] R. Neapolitan. *Probabilistic Reasoning in Expert Systems*. John Wiley, 1990.
- [19] R. Pyke. Markov renewal processes: Definitions and preliminary properties. *Annals of Mathematical Statistics*, 32(4):1231–1242, December 1961.
- [20] R. Pyke. Markov renewal processes with finitely many states. *Annals of Mathematical Statistics*, 32(4):1243–1259, December 1961.
- [21] C.M. Woodside and Y. Li. Performance Petri net analysis of communication protocol software by delay-equivalent aggregation. In *Proceedings of the 4th International Workshop on Petri nets and Performance Models (PNPM'91)*, pages 64–73, Melbourne, Australia, 2–5 December 1991. IEEE Computer Society Press.