Comparison of Multilevel Methods for Kronecker Based Markovian Representations

P. Buchholz, Dresden and T. Dayar, Ankara

Abstract

A popular way of dealing with large state spaces in Markovian modeling and analysis is to employ Kronecker based representations. This paper presents various multilevel (ML) methods for the steady state analysis of hierarchical Markovian models (HMMs) whose generator matrices have such representations. The particular ML methods are inspired by multigrid and aggregation-disaggregation techniques, and differ among each other by the type of multigrid cycle, the type of smoother, and the order of component aggregation they use. Numerical experiments demonstrate that so far ML methods with successive over-relaxation as smoother provide the most effective solvers for huge Markov chains modeled as HMMs with multiple macrostates.

AMS Subject Classifications: 65F10 (primary), 60J27 (secondary)

Keywords: Multilevel methods, multigrid, aggregation-disaggregation, Markov chains, Kronecker based numerical techniques.

1. Introduction

The attraction for Markov chains (MCs) lies in the exact results they provide up to computer precision for performance or reliability measures through numerical analysis [24]. The major problem in Markovian modeling and analysis is the exponential growth of the state space with the number of components (or subsystems) in the system. A popular way of dealing with this problem is to employ Kronecker [27] (or tensor) based representations.

In the Kronecker based approach, the system of interest is modeled so that it is formed of smaller interacting components, and its larger underlying MC is neither generated nor stored but rather represented using Kronecker products of the smaller component matrices. This introduces considerable storage savings at the expense of some overhead in the analysis phase. The concept of using Kronecker operations to define large MCs underlying structured representations appears in hierarchical Markovian models (HMMs) [4, 9, 11], or in compositional Markovian models such as stochastic automata networks (SANs) [21, 22, 14] and different classes of superposed Stochastic Petri Nets (SPNs) [13, 18]. In order to analyze large, structured Markovian models efficiently, various algorithms for vector-Kronecker product multiplication are devised [14, 10] and used as kernels in iterative solution techniques proposed for HMMs [4, 5, 7, 8], SANs [24, 25, 5, 10] and superposed Generalized SPNs [18].

In this paper, we consider the steady state analysis of HMMs, which consist of multiple low level models (LLMs) and a high level model (HLM) that defines the interaction among LLMs. SANs in the absence of functional transition rates are HMMs having one HLM state (that is, one macrostate). By introducing additional synchronized transitions, a SAN with

functional transitions can be transformed to a new SAN description without functional transitions [22]. Therefore, HMMs with multiple macrostates discussed in this paper have considerable expressive power.

Our aim is to solve

$$\pi Q = 0, \quad \sum_{i=0}^{n-1} \pi_i = 1,$$
 (1)

where Q is the infinitesimal generator (i.e., continuous-time Markov chain, CTMC) of order n underlying an HMM and π is its (row) stationary probability vector. We number the states of Q starting from 0 and assume it is irreducible implying π is also its steady state vector [24]. The matrix Q has nonnegative off-diagonal elements known as (exponential) transition rates and diagonal elements that are negated row sums of its off-diagonal elements; hence, Q has row sums of zero. Equation (1) can be viewed as a homogeneous linear system with a singular coefficient matrix of rank (n-1) subject to a normalization condition so that its solution vector can be uniquely determined.

The CTMC underlying an HMM can be expressed using sums of Kronecker products thereby facilitating the representation of considerably large Markovian models compactly. Typically the number of LLMs in an HMM is at least three, implying a problem dimension of at least four when there are multiple macrostates. The most effective solvers known for HMM problems with multiple macrostates and of dimension four or larger are block successive over-relaxation (BSOR) [7] preconditioned BiConjugate Gradient STABilized (BiCGStab) [26] and Transpose Free Quasi-Minimal Residual (TFQMR) [15] methods (see [8]). Unfortunately, these solvers are sensitive to the ordering of LLMs, the block partitionings chosen, and the amount of fill-in in the factorized diagonal blocks. This paper aims at improving the state-of-the-art in steady state solvers for such HMMs using ideas from multigrid [16, 28] and aggregation-disaggregation [24] techniques.

Multigrid (MG) techniques are iterative algorithms defined on multiple grids of increasing coarseness for the problem at hand through which the solution process proceeds in cycles until some predetermined stopping criteria are met. One cycle of MG consists of the traversal of these grids from the finest to the coarsest and back to the finest in some order. The finest grid is where the solution vector is required. The coarser grids are where smaller, approximate versions of the original problem are solved.

The general MG algorithm may be formulated recursively (see the multigrid algorithm in [23]). As boundary case, a linear system at the coarsest grid is solved. At intermediate, finer grids there are a number of consecutive recursive calls to the next coarser grid which are preceded by smoothing, residual computation, and restriction operations and are followed by interpolation (or prolongation), correction, and smoothing operations. We will refer to smoothing operations before the recursive call(s) as pre-smoothing iterations, those after as post-smoothing iterations, and the method used in the process as the smoother. At an intermediate grid, for a V-cycle (which is the standard) the number of recursive calls to the next coarser grid is one, whereas that for a W-cycle is two. An F-cycle at an intermediate grid is slightly more complicated, but can be viewed as a recursive call to a W-cycle followed by a recursive call to a V-cycle on the next coarser grid. The MG idea has

been shown to provide effective solvers for (partial) differential equations when the grids are chosen appropriately.

A multilevel (ML) algorithm inspired by MG has been presented for the steady state analysis of large, sparse MCs in [17]. Therein the restriction and interpolation operations of MG are replaced respectively with aggregation and disaggregation [24, 19], and accordingly residual computation and correction operations are omitted. The number of grids employed in the corresponding ML solver is about log_2n owing it to fact that the number of unknowns in each grid is halved at the next coarser grid. Unfortunately, this solver is hindered by the sparse generation and storage of the intermediate aggregated matrices for general MC problems.

The Kronecker structure of an HMM suggests a natural definition for the grids. Since HMM components form a hierarchy, one grid can be associated with each level of the hierarchy implying as many grids as the number of LLMs plus one for the HLM when it has multiple states (that is, the case of multiple macrostates). With this choice of grids, if one replaces the restriction and interpolation operations respectively with aggregation and disaggregation, then residual computation and correction operations disappear from the MG algorithm as in [17] and one can still utilize the Kronecker structure on intermediate aggregated matrices. Such a view formed the basis of the ML algorithm in [6] for HMMs with one macrostate whose generators can be represented as sums of Kronecker products. The particular ML solver therein employed a V-cycle and could use either the power or the Jacobi over-relaxation (JOR) method as smoother.

This paper extends the ML solver in [6] to HMMs with multiple macrostates and with the capability of using (V, W, F) cycles, (power, JOR, SOR) methods as smoothers, and (fixed, cyclic, dynamic) orders in which LLMs can be aggregated in a cycle. Then it provides the results of numerical experiments on a set of HMMs showing that the ML method with SOR smoother provides the most effective solver for HMMs with multiple macrostates so far. The storage requirements of the proposed solver can be forecasted from the HMM description and are nearly insensitive to the ordering of LLMs to the contrary of BSOR preconditioned projection methods.

The next section introduces the Kronecker based description of CTMCs underlying HMMs on a model from the literature. The third section presents the proposed class of ML methods for HMMs with multiple macrostates and discusses how they work. The fourth section provides results of numerical experiments. The fifth section concludes the paper.

2. Hierarchical Markovian Models

We introduce HMMs on an example from the literature so that the inherent structure of the Kronecker representation and ideas related to multigrid and iterative aggregation-disaggregation can be closely observed. Yet, the interested reader can find a formal definition of HMMs in [5, pp. 387-390]. We recall that the Kronecker product of two matrices $A \in \mathbb{R}^{r_A \times c_A}$ and $B \in \mathbb{R}^{r_B \times c_B}$ results in the matrix $U \in \mathbb{R}^{r_A r_B \times c_A c_B}$, which is written

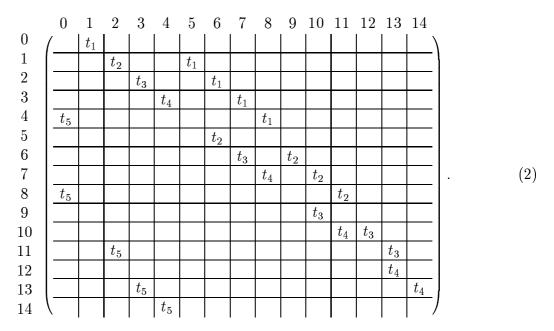
$_{\rm HLM}$	LLM 1	LLM 2	LLM 3	LLM 4	$_{ m LLM}$ 5	# of microstates						
0	17:31	0.5	0.5	0.5	0.5	15 .	6.	6.	6.	6	=	19,440
1	6:16	6:16	0.5	0.5	0.5	11 .	11 .	6.	6.	6	=	$26,\!136$
2	6:16	0.5	6:16	0.5	0.5	11 .	6.	11.	6.	6	=	$26,\!136$
3	6:16	0.5	0.5	6:16	0.5	11 .	6.	6.	11 .	6	=	$26,\!136$
4	6:16	0.5	0.5	0.5	6:16	11 .	6.	6.	6.	11	=	$26,\!136$
5	0.5	17:31	0.5	0.5	0.5	6.	15 .	6.	6.	6	=	19,440
6	0.5	6:16	6:16	0.5	0.5	6.	11.	11.	6.	6	=	$26,\!136$
7	0.5	6:16	0.5	6:16	0.5	6.	11 .	6.	11.	6	=	$26,\!136$
8	0.5	6:16	0.5	0.5	6:16	6.	11.	6.	6.	11	=	$26,\!136$
9	0.5	0.5	17:31	0.5	0.5	6.	6.	15 .	6.	6	=	$19,\!440$
10	0.5	0.5	6:16	6:16	0.5	6.	6.	11.	11.	6	=	$26,\!136$
11	0.5	0.5	6:16	0.5	6:16	6.	6.	11.	6.	11	=	$26,\!136$
12	0.5	0.5	0.5	17:31	0.5	6.	6.	6.	15 .	6	=	19,440
13	0.5	0.5	0.5	6:16	6:16	6.	6.	6.	11 .	11	=	$26,\!136$
14	0.5	0.5	0.5	0.5	17:31	6.	6.	6.	6.	15	=	19,440

Table 1: Mapping between LLM states and HLM states in $msmq_medium$.

as $U = A \otimes B$ and whose elements satisfy $u_{i_A r_B + i_B, j_A c_B + j_B} = a_{i_A, j_A} b_{i_B, j_B}$ [27]. The Kronecker sum of two square matrices $E \in \mathbb{R}^{r_E \times r_E}$ and $F \in \mathbb{R}^{r_F \times r_F}$ results in the matrix $V \in \mathbb{R}^{r_E r_F \times r_E r_F}$, which is written as $V = E \oplus F$ and defined in terms of two Kronecker products as $V = E \otimes I_{r_F} + I_{r_E} \otimes F$. Here I_{r_E} and I_{r_F} respectively denote identity matrices of orders r_E and r_F .

Example 1 We consider a model of the multiserver multiqueue discussed in [1] and name it as msmq_medium. Its HLM of 15 states describes the interaction among five LLMs, each corresponding to a finite queue of capacity 5. Customers arrive to each queue according to a Poisson process and those that arrive to a full queue get lost. There are 2 servers serving the 5 queues in a round-robin manner. When a server arrives at a queue with customers, it serves the customer at the head of the queue and travels to the next queue in line. A server that arrives to an empty queue, immediately moves to the next queue in line. Service times at queues and traveling times from one queue to the next are exponentially distributed. Note that there can be two servers simultaneously serving two different customers at the same queue. Each LLM describes a queue together with the arrival and service process, and has 32 states. The state space of each LLM is partitioned into three subsets depending on the number of servers momentarily serving the corresponding queue or traveling from the particular queue to the next in line. All states are numbered starting from 0. We name the states of the HLM as macrostates and those of Q as microstates. The mapping between LLM states and HLM states is given in Table 1. Macrostates in an HLM may have different numbers of microstates when LLMs have partitioned state spaces, as in this example. For more information about the HMM components in this example and the corresponding matrices see [5, pp. 390–392].

Six transitions denoted by t_0 , t_1 , t_2 , t_3 , t_4 , and t_5 take place in the HLM and affect the LLMs. The last five of these transitions are captured by the following (15 × 15) HLM matrix which will define the coarsest grid in the ML solver:



To each transition in the HLM matrix corresponds a Kronecker product of five (i.e., number of LLMs) LLM matrices. The matrices associated with those LLMs that do not participate in a transition are all identity. LLM 1 participates in t_1 and t_5 respectively with the matrices $Q_{t_1}^{(1)}$ and $Q_{t_5}^{(1)}$; LLM 2 participates in t_1 and t_2 respectively with the matrices $Q_{t_1}^{(2)}$ and $Q_{t_2}^{(2)}$; LLM 3 participates in t_2 and t_3 respectively with the matrices $Q_{t_3}^{(3)}$ and $Q_{t_4}^{(3)}$; LLM 4 participates in t_3 and t_4 respectively with the matrices $Q_{t_4}^{(4)}$ and $Q_{t_4}^{(4)}$; and LLM 5 participates in t_4 and t_5 respectively with the matrices $Q_{t_4}^{(5)}$ and $Q_{t_5}^{(5)}$. In general, these matrices are very sparse and therefore held in row sparse format [24]. In this example, each of the transitions t_1 , t_2 , t_3 , t_4 , t_5 affects exactly two LLMs. For instance, the Kronecker product associated with t_5 in element (4,0) of the HLM matrix in equation (2) is

$$Q_{t_5}^{(1)}(6:16,17:31)\otimes I_6\otimes I_6\otimes I_6\otimes Q_{t_5}^{(5)}(6:16,0:5),$$

where $Q_{t_5}^{(1)}(6:16,17:31)$ denotes the submatrix of $Q_{t_5}^{(1)}$ that lies between states 6 through 16 rowwise and states 17 through 31 columnwise, I_6 denotes the identity matrix of order 6, $Q_{t_5}^{(5)}(6:16,0:5)$ denotes the submatrix of $Q_{t_5}^{(5)}$ that lies between states 6 through 16 rowwise and states 0 through 5 columnwise, and \otimes is the Kronecker product operator [27]. Hence, this particular Kronecker product yields a $(26,136\times19,440)$ matrix. The rates associated with the 25 transitions in (2) are all 1. The transition rates are scalars that multiply the corresponding Kronecker products.

Other than Kronecker products due to the transitions in (2), there is a Kronecker sum implicitly associated with each diagonal element of the HLM matrix. Each Kronecker sum is formed of five LLM matrices corresponding to *local transition* t_0 . For instance, the Kronecker sum associated with element (5,5) of the HLM matrix is

$$Q_{t_0}^{(1)}(0:5,0:5) \oplus Q_{t_0}^{(2)}(17:31,17:31) \oplus Q_{t_0}^{(3)}(0:5,0:5) \oplus Q_{t_0}^{(4)}(0:5,0:5) \oplus Q_{t_0}^{(5)}(0:5,0:5),$$

where \oplus is the Kronecker sum operator. Each Kronecker sum is a sum of five Kronecker products in which all but one of the matrices are identity. The non-identity matrix in

each Kronecker product appears in the same position as in the Kronecker sum. That state changes do not take place in any but one of the LLM matrices with t_0 in each such Kronecker product is the reason behind naming t_0 a local transition. The particular Kronecker sum associated with element (5,5) of the HLM matrix is $(19,440 \times 19,440)$.

In the HLM matrix of $msmq_medium$, there do not exist any non-local transitions along the diagonal. In general, this need not be so. Therefore, we introduce the following definition.

Definition 1 In a given HMM, let K be the number of LLMs, $\mathcal{S}_{j}^{(k)}$ be the subset of states of LLM k mapped to macrostate j, $\mathcal{T}_{i,j}$ be the set of LLM non-local transitions in element (i,j) of the HLM matrix, $rate_{t_e}(i,j)$ be the rate associated with transition $t_e \in \mathcal{T}_{i,j}$, and D_j be the diagonal (correction) matrix that sums the rows of Q corresponding to macrostate j to zero. Then the diagonal block (j,j) of Q corresponding to element (j,j) of the HLM matrix is given by

$$Q_{j,j} = \bigoplus_{k=1}^{K} Q_{t_0}^{(k)}(\mathcal{S}_j^{(k)}, \mathcal{S}_j^{(k)}) + \sum_{t_e \in \mathcal{T}_{j,j}} rate_{t_e}(j,j) \bigotimes_{k=1}^{K} Q_{t_e}^{(k)}(\mathcal{S}_j^{(k)}, \mathcal{S}_j^{(k)}) + D_j,$$
(3)

and, when there are multiple macrostates, the off-diagonal block (i, j) of Q corresponding to element (i, j) of the HLM matrix is given by

$$Q_{i,j} = \sum_{t_e \in \mathcal{T}_{i,j}} rate_{t_e}(i,j) \bigotimes_{k=1}^K Q_{t_e}^{(k)}(\mathcal{S}_i^{(k)}, \mathcal{S}_j^{(k)}).$$
(4)

When there are multiple macrostates, Q is a block matrix having as many blocks in each dimension as the number of macrostates (i.e., order of the HLM matrix). The diagonal of Q is formed of its negated off-diagonal row sums, and may be stored explicitly or can be generated as needed.

In Example 1, the second term in equation (3) is missing. Although Q in $msmq_medium$ is of order 358,560 and has 2,135,160 nonzeros, its Kronecker representation needs to store 1 HLM matrix having 25 nonzeros and 15 LLM matrices (since identity matrices are not stored) having a total of 370 nonzeros. This is a substantial saving in storage.

If we neglect the diagonal of Q which is handled separately, from Definition 1 it follows that each nonzero element of the HLM matrix is essentially a sum of Kronecker products. This has a very nice implication on the choice of grids in the proposed ML solver when component aggregation is used in forming the coarser grids. The HLM and LLMs 1 through K define the least coarsest (in other words, the finest) grid. This grid is Q. Since the HLM holds the LLMs together, it will define the coarsest grid, which is not aggregated. Regarding the intermediate grids, let us assume that LLMs are aggregated starting from 1 up to K. Then the HLM and LLMs 2 through K define the first coarser grid when LLM 1 is aggregated. The HLM and LLMs 3 through K define the second coarser grid when LLMs 1 and 2 are aggregated; and so on. We consider a form of aggregation in which all grids are irreducible and have row sums of zero.

Now, let us concentrate on the sizes of the grids defined by the HLM and LLMs for the assumed order in which LLMs are aggregated. In Example 1, the grids defined in this way by HLM, HLM and LLMs 5, HLM and LLMs 4-5, HLM and LLMs 3-5, HLM and LLMs 2-5, HLM and LLMs 1-5 have respectively the sizes (15×15) , (119×119) , (913×913) , $(6,822 \times 6,822)$, $(49,896 \times 49,896)$, $(358,560 \times 358,560)$ (see Table 1 and equations (3)-(4)). Obviously, one is not restricted to aggregating LLMs in the order 1 through K, and can consider other orders. In the next section, we introduce the ML method with the grid choices suggested by the Kronecker structure of HMMs and remark that the grids are never explicitly generated.

3. ML Methods for HMMs with Multiple Macrostates

The class of ML methods introduced in this section for HMMs with multiple macrostates have the capability of using (V, W, F) cycles, (power, JOR, SOR) methods as smoothers, and (fixed, cyclic, dynamic) orders in which LLMs can be aggregated in a cycle. These parameters are respectively denoted by C, S, and O. We remark that $C \in \{V, W, F\}$, $S \in \{POWER, JOR, SOR\}$, and $O \in \{FIXED, CYCLIC, DYNAMIC\}$. In a particular ML solver, C, S, and O are fixed at the beginning.

In Algorithm 1, we give the recursive ML function that is invoked for LLMs. It is the driver in Algorithm 2 where the particular ML solver starts executing at the finest grid involving the HLM and all the LLMs, and then invokes the recursive ML function with the order of aggregation in the list \mathcal{C} . Each pass through the body of the repeat-until loop in Algorithm 2 corresponds to one cycle of the ML method. One will notice that steps 3-9 in Algorithm 1 are almost identical to the statements between step 3 and 4 in Algorithm 2. Nevertheless, Algorithm 2 is coded separately since the finest grid is treated somewhat differently as we next explain.

The order of aggregating LLMs in each ML cycle is determined by the list \mathcal{J} defined in Algorithm 2. The elements of \mathcal{J} from its head to its tail are denoted respectively by $\mathcal{J}_1, \mathcal{J}_2, \ldots, \mathcal{J}_{K+1}$. The subscripts of these elements indicate their orders in \mathcal{J} . In each ML cycle, HLM is always the last model to be handled due to its special position in the hierarchy. Hence, \mathcal{J}_{K+1} is given the value K+1 and is associated with the HLM; this never changes. Initially, LLM k is associated with element \mathcal{J}_k which has the value k for $k=1,2,\ldots,K$ (see step 1 of Algorithm 2). In each ML cycle, LLMs are aggregated according to these values starting from the element at the head of the list (see the second statement in the repeat-until loop of Algorithm 2). Hence, LLM \mathcal{J}_1 is the first LLM to be aggregated.

In the FIXED order of aggregating LLMs, the initial assignment of values to the elements of \mathcal{J} does not change after the ML method starts executing; this is the default order. In the CYCLIC order, at the end of each ML cycle a circular shift of elements \mathcal{J}_1 through \mathcal{J}_K in the list are performed; this ensures some kind of fairness in aggregating LLMs in the next ML cycle. On the other hand, the DYNAMIC order sorts the elements \mathcal{J}_1 through

 \mathcal{J}_K according to the residual norms projected (or restricted) to the corresponding LLM at the end of the ML cycle, and aggregates the LLMs in this sorted order in the next ML cycle (see step 8 of Algorithm 2). This ensures that LLMs which have smaller residual norms are aggregated earlier at finer grids. We expect small residual norms to be indicative of good approximations in those components. Note that at each intermediate grid, the recursive ML function is invoked for the next coarser grid with the list of LLMs in $\mathcal C$ which is formed by removing the LLM at the head of the incoming list \mathcal{D} by aggregation (see step 4 in Algorithm 1). Once the list of LLMs is exhausted, that is K+1 is the only value remaining in the list \mathcal{D} , backtracking from the recursion starts by solving a linear system as large as the HLM matrix (see the first if statement in Algorithm 1).

Algorithm 1 Recursive ML function on LLMs in \mathcal{D}

function $\mathrm{ML}(Q_{\mathcal{D}}, x_{\mathcal{D}}, \mathcal{D}, \gamma)$

if
$$(|\mathcal{D}| == 1)$$
 then
$$y_{\mathcal{D}} = \text{solve}(Q_{\mathcal{D}}, x_{\mathcal{D}}); \qquad \text{(step 1)}$$
if $(C == F)$ then
$$\gamma = 1;$$

else

else
$$x_{\mathcal{D}} = S(Q_{\mathcal{D}}, x_{\mathcal{D}}, w, MIN_IN_PRE, MAX_IN_PRE, \rho, RES_COUNT, \eta_1); \quad (\text{step 3})$$

$$\mathcal{C} = \mathcal{D} - [head(\mathcal{D})]; \quad (\text{step 4})$$

$$\text{compute } Q_{\mathcal{C}} \text{ from } Q_{\mathcal{D}} \text{ and } x_{\mathcal{D}}; \quad (\text{step 5})$$

$$x_{\mathcal{C}}(s_{\mathcal{C}}) = \sum_{s_{\mathcal{D}} \in \mathcal{S}_{\mathcal{D}}, proj(s_{\mathcal{D}}, \mathcal{D}, \mathcal{C}) = = s_{\mathcal{C}}} x_{\mathcal{D}}(s_{\mathcal{D}}) \text{ for all } s_{\mathcal{C}} \in \mathcal{S}_{\mathcal{C}}; \quad (\text{step 6})$$

$$\text{if } (\gamma == 1) \text{ then} \quad (\text{step 7})$$

$$y_{\mathcal{C}} = \text{ML}(Q_{\mathcal{C}}, x_{\mathcal{C}}, \mathcal{C}, \gamma);$$

$$\text{else}$$

$$y_{\mathcal{C}} = \text{ML}(Q_{\mathcal{C}}, x_{\mathcal{C}}, \mathcal{C}, \gamma);$$

$$y_{\mathcal{C}} = \text{ML}(Q_{\mathcal{C}}, x_{\mathcal{C}}, \mathcal{C}, \gamma);$$

$$y_{\mathcal{D}}(s_{\mathcal{D}}) = x_{\mathcal{D}}(s_{\mathcal{D}}) \frac{y_{\mathcal{C}}(proj(s_{\mathcal{D}}, \mathcal{D}, \mathcal{C}))}{x_{\mathcal{C}}(proj(s_{\mathcal{D}}, \mathcal{D}, \mathcal{C}))} \text{ for all } s_{\mathcal{D}} \in \mathcal{S}_{\mathcal{D}};$$

$$y_{\mathcal{D}} = S(Q_{\mathcal{D}}, y_{\mathcal{D}}, w, MIN_IN_POST, MAX_IN_POST, \rho, RES_COUNT, \eta_2); \text{ (step 9)}$$

$$\text{return}(y_{\mathcal{D}});$$

Now we discuss the operation that computes the next coarser grid $Q_{\mathcal{C}}$ from the grid $Q_{\mathcal{D}}$ using the vector $x_{\mathcal{D}}$ (see step 5 in Algorithm 1) by aggregating the component at the head of list \mathcal{D} (i.e., $head(\mathcal{D})$) and thus forming the list of components in \mathcal{C} . Let $\mathcal{S}_{\mathcal{D}}$ and $\mathcal{S}_{\mathcal{C}}$ respectively denote the state spaces of components in \mathcal{D} and \mathcal{C} . A positive number of states $s_{\mathcal{D}} \in \mathcal{S}_{\mathcal{D}}$ are projected to each state $s_{\mathcal{C}} \in \mathcal{S}_{\mathcal{C}}$. We represent this in the algorithms as $\exists s_{\mathcal{D}} \in \mathcal{S}_{\mathcal{D}}, proj(s_{\mathcal{D}}, \mathcal{D}, \mathcal{C}) == s_{\mathcal{C}}$ for all $s_{\mathcal{C}} \in \mathcal{S}_{\mathcal{C}}$. We also remark that there are no unreachable states in HMMs and their underlying CTMCs are always irreducible. Hence, this projection is surjective (or onto).

Algorithm 2 ML driver

```
main()
\mathcal{J} = [1, 2, \dots, K+1]; x = \text{initial approximation}; it = 0; cyc = 0; stop = FALSE; (step 1)
if (C == W \text{ or } C == F) then
                                                                                                                                     (step 2)
     \gamma = 2;
else
     \gamma = 1;
                                                                                                                                     (step 3)
repeat
     x = S(Q, x, w, MIN\_OUT\_PRE, MAX\_OUT\_PRE, \rho, RES\_COUNT, \nu_1);
     \mathcal{C} = \mathcal{J} - [head(\mathcal{J})];
     compute Q_{\mathcal{C}} from Q and x;
     x_{\mathcal{C}}(s_{\mathcal{C}}) = \sum_{s \in \mathcal{S}, proj(s, \mathcal{J}, \mathcal{C}) = = s_{\mathcal{C}}} x(s) \text{ for all } s_{\mathcal{C}} \in \mathcal{S}_{\mathcal{C}};
     if (\gamma == 1) then
          y_{\mathcal{C}} = \mathrm{ML}(Q_{\mathcal{C}}, x_{\mathcal{C}}, \mathcal{C}, \gamma);
     else
          y_{\mathcal{C}} = \mathrm{ML}(Q_{\mathcal{C}}, x_{\mathcal{C}}, \mathcal{C}, \gamma);
         y_{\mathcal{C}} = \mathrm{ML}(Q_{\mathcal{C}}, x_{\mathcal{C}}, \mathcal{C}, \gamma);
     y(s) = x(s) \frac{y_{\mathcal{C}}(proj(s,\mathcal{J},\mathcal{C}))}{x_{\mathcal{C}}(proj(s,\mathcal{J},\mathcal{C}))} for all s \in \mathcal{S};
     y = S(Q, y, w, MIN\_OUT\_POST, MAX\_OUT\_POST, \rho, RES\_COUNT, \nu_2);
     if (C == F) then
                                                                                                                                     (step 4)
          \gamma = 2;
     x = y; it = it + \nu_1 + \nu_2; cyc = cyc + 1;
                                                                                                                                     (step 5)
     normalize(x); r = xQ;
                                                                                                                                     (step 6)
     if (it \ge MAX\_IT \text{ or } time \ge MAX\_TIME \text{ or } ||r|| \le STOP\_TOL)
                                                                                                                                     (step 7)
          stop = TRUE;
     else
         if (O == DYNAMIC) then
                                                                                                                                     (step 8)
              sort LLM indices \mathcal{J}_1, \mathcal{J}_2, \ldots, \mathcal{J}_K into increasing order of ||r_k||,
               where r_k is the residual associated with LLM k and is computed from r;
          else if (O == CYCLIC) then
               circular_shift(\mathcal{J}_1, \mathcal{J}_2, \ldots, \mathcal{J}_K);
until (stop);
```

For each state $s_{\mathcal{C}} \in \mathcal{S}_{\mathcal{C}}$, the columns of the grid $Q_{\mathcal{D}}$ corresponding to the states in $\mathcal{S}_{\mathcal{D}}$ that get projected to the same state $s_{\mathcal{C}}$ are summed. This yields a column aggregated grid whose row sums are zero given that $Q_{\mathcal{D}}$ has row sums of zero. For each state $s_{\mathcal{C}} \in \mathcal{S}_{\mathcal{C}}$, the rows of this column aggregated grid corresponding to the states in $\mathcal{S}_{\mathcal{D}}$ that are projected

take x as the steady state vector π of the HMM;

to the same state $s_{\mathcal{C}}$ are multiplied with the corresponding elements of the row vector $x_{\mathcal{D}}$ and summed. This yields the square grid $Q_{\mathcal{C}}$ which has row sums of zero regardless of the norm of $x_{\mathcal{D}}$. We remark that the grid $Q_{\mathcal{C}}$ is irreducible as long as $x_{\mathcal{D}} > 0$ and $Q_{\mathcal{D}}$ is irreducible (see also [6, pp. 346-348]). In practice, $Q_{\mathcal{C}}$ is not explicitly generated, but represented by allocating a vector per each non-local transition in the HLM matrix that is as long as the number of states in the state space of the LLMs in $\mathcal{C} - [tail(\mathcal{C})]$ that get mapped to the particular row of the HLM matrix. Note that the tail of \mathcal{C} has the value K+1 corresponding to the HLM. These vectors are used so as to facilitate the operations on the coarser grid $Q_{\mathcal{C}}$ in the recursive ML function. Hence, 25 vectors need to be kept for each intermediate grid except the coarsest in Example 1. Since these vectors are generally much shorter than n, they do not bring considerable storage overhead to the ML method.

The aggregation on the columns of $Q_{\mathcal{D}}$ is also performed on the columns of the row vector $x_{\mathcal{D}}$ yielding the vector $x_{\mathcal{C}} > 0$ when $x_{\mathcal{D}} > 0$ (see step 6 in Algorithm 1). In [6, p. 348] it has been shown that $x_{\mathcal{C}}$ is the stationary vector of $Q_{\mathcal{C}}$ if $x_{\mathcal{D}}$ is the stationary vector of $Q_{\mathcal{D}}$. Step 8 in Algorithm 1 corresponds to the opposite of what is done in step 6; that is, it performs disaggregation using $x_{\mathcal{D}}$, $x_{\mathcal{C}}$, and the newly computed vector $y_{\mathcal{C}}$ to obtain the vector $y_{\mathcal{D}}$. Similar aggregation and disaggregation operations are performed in Algorithm 2 at the finest grid Q.

The variable γ in the two algorithms determines the number of recursive calls to the ML function. In step 2 of the initialization statements before the repeat-until loop in Algorithm 2, γ is set to 2 for a W- or an F- cycle and set to 1 for a V-cycle. After this point, there are two places where the value of γ changes, and these happen only for an F-cycle. Hence, for a V-cycle γ remains 1 and for a W-cycle it remains 2, meaning for V- and W-cycles respectively 1 and 2 recursive calls are made to the ML function on the next coarser grid. On the other hand, for an F-cycle γ is set to 1 at the boundary case of the recursion (see step 2 in Algorithm 1). Recall that an F-cycle can be seen as a recursive call to a W-cycle followed by a recursive call to a V-cycle. After the F-cycle is over, γ is reset to 2 in step 4 of Algorithm 2 so as to be ready for a new cycle [28, pp. 174-175].

Each ML cycle starts and ends with some number of iterations using the smoother S. See respectively the two statements right after step 3 and right before step 4 in Algorithm 2. The same is true for each execution of the recursive ML function at intermediate grids as can be seen in steps 3 and 9 of Algorithm 1. The first two arguments of the call to the smoothers in both algorithms represent respectively the grid to be used in the smoothing process and the vector to be smoothed. The user is given the flexibility to specify different numbers of pre- and post-smoothings in the two algorithms. Hence, we have the nonnegative integer pairs of parameters $(MIN_OUT_PRE, MAX_OUT_PRE)$, $(MIN_OUT_POST, MAX_OUT_POST)$ for the finest grid handled by Algorithm 2, and (MIN_IN_PRE, MAX_IN_PRE) , $(MIN_IN_POST, MAX_IN_POST)$ for the coarser intermediate grids handled by Algorithm 1.

For each pair of parameters (MIN_*, MAX_*) , S performs MAX_* smoothings when $MIN_* \ge MAX_*$. When $MIN_* < MAX_*$, S performs an adaptive number of smoothings using the two parameters ρ and RES_COUNT as follows. Upon entry to the smoother,

the residual norm of the current solution vector is computed and recorded. Then MIN_* smoothings are performed and the residual norm of the solution vector is recomputed. If the ratio of the two residual norms is less than ρ , then S stops executing; otherwise, smoothings continue till MAX_* iterations or the ratio of residual norms of two solution vectors RES_COUNT iterations apart are less than ρ . Note that the computation of the residual vector requires an extra implicit vector-grid multiply when S is SOR. However, this is performed only every RES_COUNT smoothings once the smoother is beyond MIN_* smoothings. The parameter w in the call to the smoother is the relaxation parameter for JOR and SOR. The parameters (ν_1, ν_2) and (η_1, η_2) can be used to keep track respectively of the number of (pre-, post-) smoothings at the finest and coarser grids.

We start the ML iteration with x set to the uniform distribution. At the end of each ML cycle the solution vector x is normalized and the residual vector r = xQ is computed. Global convergence of iterative aggregation-disaggregation on CTMCs with the possibility of using iterative methods to solve the aggregated matrices appears in [20]. This result suggests that as long as a sufficient number of smoothings are performed at each grid, the ML method should converge. See also the related comments in [6, p. 350]. The ML iteration continues until the total number of smoothings at the finest grid (i.e., it) exceeds MAX_IT, the CPU time exceeds MAX_TIME, or the residual norm of the solution vector (i.e., ||r||) meets the prespecified stopping tolerance, $STOP_TOL$. At that point, x is taken as the approximation of the steady state vector π . The variable cyc counts the number of ML cycles performed until stopping in case this information cannot be obtained from it. Note that this is possible for an adaptive number of pre- or post-smoothings at the finest grid. Finally, we remark that the smoothers of choice require two vectors of length n and two vectors (three in SOR) as long as the maximum number of microstates per macrostate in the HMM. One of the vectors of length n in SOR is required for the computation of residuals in the implementation of DYNAMIC ordering of LLMs for aggregation. The next section presents numerical results with the proposed class of ML solvers.

4. Numerical Experiments

We implemented the ML method as discussed in the previous section in C as part of the APNN toolbox [3, 2]. Now the toolbox has three ML solvers named ML_POWER, ML_JOR, and ML_SOR depending on the smoother used. We report the results of three sets of numerical experiments although other experiments that provide results along the same direction have also been performed.

First are the results of the msmq_medium problem introduced in Example 1 and the courier_medium problem [7], which are discussed in the next subsection. In these medium sized problems, we compare the results of ML solvers among each other and with those of STR_POWER, STR_JOR, STR_RSOR, which respectively implement power, JOR, SOR methods, and also with those of STR_BSOR, BSOR_BICGSTAB, and BSOR_TFQMR available in the APNN toolbox. Here, STR_BSOR is the two-level version of the BSOR solver in [7] that takes advantage of various techniques to reduce the amount of fill-in when factorizing diagonal blocks of the chosen partitioning. The solvers BSOR_BICGSTAB and

Attribute	$msmq_large$	$courier_large$	$qh_real control$
HLM states	$\{0:34\}$	$\{0:12\}$	$\{0:8\}$
nz_{HLM}	75	65	18
	$(rates \in \{10\})$	$(rates \in \{1, 1449.3, 4821.6, 8771.9\})$	$(rates \in \{10,000\})$
LLM 1 states	$\{0:6,7:19,$	$\{0:29\}$	$\{0:76,77:123,$
	20:37,38:59		$124:170,171:202$ }
LLM 2 states	$\{0:6,7:19,$	$\{0,1:140,141,142:201,202,$	$\{0:61,62:99,$
	20:37,38:59	$203:222,223,224:227,228\}$	$100:137,138:163$ }
LLM 3 states	$\{0:6,7:19,$	$\{0:321,322:326,327:470,$	$\{0:56,57:91,$
	20:37,38:59	471:474,475:526,527:529,	$92:126,127:150\}$
		$530:542,543:544,545\}$	
LLM 4 states	$\{0:6,7:19,$	$\{0:14\}$	None
	20:37,38:59		
LLM 5 states	$\{0:6,7:19,$	None	None
	20:37,38:59		
LLM matrices	15	14	15
nz_{LLMs}	790	2,333	$1,\!486$
$\mathcal{T}(i,j),\ i \neq j$	$\{t_{14},t_{15},t_{16},t_{17},t_{18}\}$	$\{t_0,t_{28},t_{29},t_{30}\}$	$\{t_{17},t_{18},t_{19},t_{21},t_{24},t_{27}\}$
$\mathcal{T}(j,j)$	None	$\{t_{17},t_{23}\}$	None
n	2,945,880	1,632,600	$399,\!476$
nz	19,894,875	9,732,330	1,871,004

Table 2: Benchmark problems.

BSOR_TFQMR are respectively BSOR preconditioned versions of the projection methods BiCGStab and TFQMR in [8]. To the best of our knowledge, BSOR_BICGSTAB and BSOR_TFQMR are the most competitive solvers for HMMs with multiple macrostates when they utilize favorable block partitionings.

Second are the results of experiments that shed light to the scalability of the ML method. These appear in subsection 4.2. Having observed in a multitude of experiments that SOR performs best among the three smoothers, we run ML_SOR on different dimensioned versions of the msmq problem and compare the results with those of STR_RSOR.

Third are the results of experiments in subsection 4.3 on three problems named $msmq_large$, $courier_large$, and $qh_realcontrol$, again with multiple macrostates, using ML_SOR, STR_RSOR, STR_BSOR, BSOR_BICGSTAB, and BSOR_TFQMR. The first two of these problems are larger versions of the two problems used in the first set of experiments in subsection 4.1. The medium version of the msmq problem has been discussed in detail in section 2. We introduce the medium version of the courier problem in subsection 4.1.

The characteristics of the three problems used in subsection 4.3 are given in Table 2. For each of them, we provide the macrostates (HLM states), the number of nonzeros in the HLM matrix (nz_{HLM}) and their values (rates), the state space partition of each LLM (LLM states), the number of LLM matrices (LLM matrices), the total number of nonzeros in LLM matrices (nz_{LLMs}) , the transitions in the off-diagonal part $(\mathcal{T}(i,j), i \neq j)$ and the diagonal part $(\mathcal{T}(j,j))$ of the HLM matrix, the number of states (n) and the number of nonzeros (nz) of the underlying CTMC.

The CTMCs underlying all problems in this paper are irreducible. The techniques pro-

posed in [7] to reduce the amount of fill-in require modest storage for the factors of the diagonal blocks in BSOR and BSOR preconditioned projection methods in these problems. Albeit smaller, $qh_realcontrol$ is known to be rather difficult to solve. BSOR preconditioned projection methods perform particularly well on this problem. Hence, the results of the last set of experiments will especially indicate the effectiveness of the class of ML solvers.

In each set of experiments we use w=1.0 wherever required, implying Jacobi, Gauss-Seidel (GS), and block GS (BGS) methods rather than JOR, SOR, and BSOR, respectively. Furthermore, as smoother parameters in the ML solvers we let

```
\{(1,1,1,1,1,1,1,1),
(MIN\_IN\_PRE,
                             (3, 3, 3, 3, 1, 1, 1, 1),
MAX\_IN\_PRE,
                             (1, 1, 1, 1, 3, 3, 3, 3),
MIN\_IN\_POST.
                             (3, 3, 3, 3, 3, 3, 3, 3),
MAX\_IN\_POST,
                            (5, 5, 5, 5, 1, 1, 1, 1, 1),
MIN\_OUT\_PRE,
                             (1, 1, 1, 1, 5, 5, 5, 5),
MAX\_OUT\_PRE.
                             (5, 5, 5, 5, 5, 5, 5, 5),
MIN\_OUT\_POST,
                             (10, 10, 10, 10, 1, 1, 1, 1),
MAX\_OUT\_POST)
                             (1, 1, 1, 1, 10, 10, 10, 10),
                             \{10, 10, 10, 10, 10, 10, 10, 10\}
```

in the cases where a fixed number of smoothings are performed. Besides, we set $RES_COUNT = 2$ and $\rho = 0.9$, and let $(MIN_IN_PRE, MAX_IN_PRE, MIN_IN_POST, MAX_IN_POST, MIN_OUT_PRE, MAX_OUT_PRE, MIN_OUT_POST, MAX_OUT_POST) \in \{(1,3,1,3,1,1,1,1), (1,5,1,5,1,1,1,1), (1,10,1,10,1,1,1,1)\}$ to experiment with an adaptive number of inner pre- and post-smoothing iterations. For each setting of the smoother parameters, we experiment with (V, W, F) cycles and (FIXED, CYCLIC, DYNAMIC) ordering of LLMs for aggregation. Hence, we perform $(10+3) \times 3 \times 3 = 117$ experiments per smoother.

All experiments are performed on a 550 MHz Pentium III processor and a 1 GBytes main memory under Linux. The large main memory is necessary due to the large number of vectors of length n used in BSOR preconditioned projection methods. All times are reported as seconds of CPU time. In Tables 4 through 7, we report the times spent in setup and iterative parts of the solvers respectively under columns Setup and Solve, and indicate the fastest solvers in bold. For each smoother, we list the best three ML solvers out of the 117 considered. The it column indicates the number of (outer) iterations it takes the solvers to stop and the res column indicates the infinity norm of the residual upon stopping. We set $MAX_IT = 1,000$, $MAX_TIME = 1,000$ seconds, and $STOP_TOL = 10^{-8}$ in all solvers. In BSOR_BICGSTAB and BSOR_TFQMR, each pass through the body of the code counts as two iterations rather than one since two vector-matrix products are computed. We choose to normalize the solution vector and compute the residual every 10 iterations in the solvers STR_POWER, STR_JOR, STR_RSOR and STR_BSOR.

For the problems in which convergence is observed due to the stopping tolerance of 10^{-8} but the norm of the residual is found to be larger than 10^{-8} , we continued the iterative process by decreasing the stopping tolerance one order of magnitude at a time until we

			~	
Solver	it	res	Setup	Solve
STR_POWER	1,000	10^{-5}	0	303
$ML_POWER(3,3,3,3,1,1,1,1), CYCLIC, W$	52	10^{-9}	0	45
$ML_POWER(5,5,5,5,1,1,1,1), CYCLIC, W$	46	10^{-9}	0	48
$ML_POWER(5,5,5,5,1,1,1,1), DYNAMIC, W$	44	10^{-9}	0	47
STR_JOR	720	10^{-9}	0	237
$ML_{JOR}(3,3,3,3,1,1,1,1), CYCLIC, W$	26	10^{-9}	0	23
$ML_{JOR}(3,3,3,3,1,1,1,1), CYCLIC, F$	32	10^{-9}	0	28
$ML_{JOR}(5,5,5,5,1,1,1,1), CYCLIC, F$	32	10^{-9}	0	32
STR_RSOR	240	10^{-9}	0	104
$ML_SOR(1,1,1,1,1,1,1,1,1), CYCLIC, W$	24	10^{-9}	0	21
$ML_SOR(1,1,1,1,1,1,1,1,1), CYCLIC, F$	22	10^{-9}	0	19
$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, W$	18	10^{-9}	0	21
STR_BSOR	120	10^{-9}	2	62
BSOR_BICGSTAB	47	10^{-9}	2	43
$BSOR_TFQMR$	46	10^{-10}	2	42

Table 3: ML method on msmq_medium

Table 4: ML method on courier_medium

Solver	it	res	Setup	Solve
STR_POWER	1,000	10^{-2}	1	841
$\operatorname{ML_POWER}(5,5,5,5,5,5,5,5), CYCLIC, W$	590	10^{-9}	1	693
$\text{ML_POWER}(5,5,5,5,5,5,5,5), CYCLIC, F$	610	10^{-9}	1	709
$ML_{POWER}(10,10,10,10,10,10,10,10), CYCLIC, W$	640	10^{-9}	1	701
STR_JOR	1,000	10^{-5}	1	924
$ML_JOR(1,1,1,1,1,1,1,1), CYCLIC, F$	116	10^{-9}	1	213
$ML_JOR(3,3,3,3,1,1,1,1), CYCLIC, V$	106	10^{-9}	1	199
$ML_{JOR}(5,5,5,5,1,1,1,1), DYNAMIC, V$	108	10^{-9}	1	213
STR_RSOR	360	10^{-9}	1	503
$ML_SOR(5,5,5,5,1,1,1,1), CYCLIC, V$	28	10^{-9}	1	82
$ML_SOR(5,5,5,5,5,5,5,5), FIXED, V$	50	10^{-9}	1	78
$ML_SOR(5,5,5,5,5,5,5,5), FIXED, W$	50	10^{-9}	1	84
STR_BSOR	60	10^{-10}	4	154
BSOR_BICGSTAB	37	10^{-9}	4	124
BSOR_TFQMR	40	10^{-9}	4	133

encountered a residual norm less then 10^{-8} . Such a situation is witnessed among BSOR preconditioned projection methods since we work with unnormalized solution vectors and the underlying CTMCs are not scaled. Recall that the system we solve is singular and a non-scaled coefficient matrix with considerably large entries may result in the residual norm being larger than what the (unnormalized) solution vector actually implies (see [12, p. 1697]) especially when convergence takes place rapidly. In only one of the problems we are not able to reduce the residual norm below 10^{-8} by iterating in this manner, and that happens to be with the BSOR_TFQMR solver in $qh_realcontrol$ in subsection 4.3 where the residual norm is computed to be 1.2×10^{-8} .

4.1 Performance of ML Solvers on msmq_medium and courier_medium

Other than $msmq_medium$, in this subsection we consider a a typical benchmark from the literature which is introduced in [29]. We name this model as $courier_medium$. Its HLM,

which has 10 states, describes the interaction among four LLMs. LLM 1 has 15 states, LLM 2 has 217 states, LLM 3 has 88 states, and LLM 4 has 30 states. The state spaces of LLM 2 and 3 are respectively partitioned as 0:1, 2, 3:5, 6:18, 19:22, 23:74, 75:216 and 0, 1, 2:5, 6, 7:26, 27, 28:87. In the particular HLM under consideration, six transitions denoted by t_0 through t_5 take place and affect the LLMs. LLM 2 and 3 each participates in four transitions while each of the other two LLMs participates in one transition. Contrary to $msmq_medium$, the HLM matrix of $courier_medium$ has non-local transitions along its diagonal. The rates associated with all HLM transitions are 1. Although the underlying CTMC has 419,400 states and 2,281,620 nonzeros, the Kronecker representation associated with the HMM needs to store 1 HLM matrix having 47 nonzeros and 14 LLM matrices having a total of 845 nonzeros. Note that medium sized problems have in the order of 100,000 states.

In Tables 3 and 4 we provide the results of numerical experiments with the $msmq_medium$ and $courier_medium$ problems, respectively. The setup times of the ML solvers are all negligible. The $courier_medium$ problem seems to be more difficult to solve than $msmq_medium$ due to the longer time it takes to be solved by a particular solver, and benefits relatively more from a larger number of pre- and post-smoothings at intermediate grids in ML solvers. The ML_POWER and ML_SOR solvers in Table 4 also benefit relatively more than those in Table 3 from a larger number of smoothings at the finest grid.

The winner in both problems happens to be an ML_SOR solver, which is significantly better than BSOR_BICGSTAB and BSOR_TFQMR. The winners in msmq_medium and courier_medium require respectively 11 and 10 ML cycles. The quality of the three smoothers are observed to increase in the order POWER, JOR, and SOR. We also see that adaptive number of smoothings do not yield the best ML solvers. This seems to be due to the extra effort spent in computing the residual norms (at least twice in each call to the smoother) to facilitate adaptiveness although a small number of iterations are required of the smoothers for rapid convergence. At least six ML solvers in each of the two tables use CYCLIC order of aggregating LLMs. In msmq_medium W and F cycles are used in the best ML solvers, whereas V cycle is associated with the two fastest ML_SOR solvers in courier_medium. Interestingly, in the more difficult of the two problems, FIXED ordering of LLMs for aggregation works better with the two of the three fastest ML_SOR solvers. Finally, we also observe that increasing the number of pre- and post-smoothings at the finest grid tends to reduce the number of ML cycles to convergence (see the ML_POWER and ML_SOR results in Table 4). A similar statement can also be made for smoothings at intermediate grids. However, a smaller number of ML cycles does not necessarily imply a shorter solution time.

4.2 Scalability of ML_SOR

In this subsection we investigate the scalability of the ML_SOR solver. We consider five HMMs of different dimensions related to the model introduced in section 2. These HMMs are named $msmq_c3$, $msmq_c4$, $msmq_c5$, $msmq_c6$, $msmq_c7$ and have respectively 3, 4, 5, 6, 7 LMMs. Hence, the largest model among these HMMs with multiple macrostates

Problem	Solver	it	res	Setup	Solve
$msmq_c3$	STR_RSOR	180	10^{-9}	0	0
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, W$	44	10^{-9}	0	0
	$ML_SOR(1,1,1,1,1,1,1,1,1), DYNAMIC, W$	44	10^{-9}	0	0
	$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, F$	28	10^{-9}	0	0
$msmq_c4$	STR_RSOR	260	10^{-9}	0	1
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, W$	40	10^{-9}	0	1
	$ML_SOR(1,1,1,1,1,1,1,1), DYNAMIC, W$	38	10^{-9}	0	1
	$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, F$	26	10^{-9}	0	1
$msmq_c5$	STR_RSOR	290	10^{-9}	0	13
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, W$	34	10^{-9}	0	4
	$ML_SOR(1,1,1,1,1,1,1,1), DYNAMIC, W$	32	10^{-9}	0	4
	$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, F$	24	10^{-9}	0	4
$msmq_c6$	STR_RSOR	360	10^{-9}	0	109
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, W$	30	10^{-9}	0	22
	$ML_SOR(1,1,1,1,1,1,1,1,1), DYNAMIC, W$	30	10^{-9}	0	23
	$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, F$	22	10^{-9}	0	23
$msmq_c7$	STR_RSOR	420	10^{-9}	1	873
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, W$	30	10^{-9}	1	148
	$ML_SOR(1,1,1,1,1,1,1,1), DYNAMIC, W$	30	10^{-9}	1	153
	$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, F$	22	10^{-9}	1	154

Table 5: Scalability of ML_SOR

has a problem dimension of eight. The number of LLMs indicate the number of queues in the corresponding HMM. There are 2 servers serving the queues in each HMM in a round-robin manner as in $msmq_medium$. In all HMMs, each LLM matrix is of order 20 (due to a finite queueing capacity of 3) and the 20 states are partitioned into the three subsets 0-3, 4-10, and 11-19. Each LLM participates in two non-local transitions as in $msmq_medium$. The HLM matrices corresponding to the five HMMs have respectively 6, 10, 15, 21, 28 macrostates and 9, 16, 25, 36, 49 nonzeros. The values of all nonzeros in the HMM matrices are 10. The number of LLM matrices and their total number of nonzeros are respectively 9, 12, 15, 18, 21 and 132, 176, 220, 264, 308. The number of microstates and number of nonzeros in the underlying MCs of the HMMs are respectively 1,020, 7,008, 42,880, 243,456, 1,311,744 and 3,969, 32,976, 235,680, 1,527,552, 9,241,344.

Table 5 presents the results of STR_RSOR and the best three ML_SOR solvers for the five msmq problems. Note that, the ML_SOR solvers across all problems in Table 5 have the same parameters. Each of the solvers performs two smoothings at the finest grid, use either CYLIC or DYNAMIC order of aggregation, and utilize W or F cycle. The number of ML cycles to convergence range between 11 and 22, but do not vary significantly for a specific ML_SOR solver. On the other hand, the number of iterations performed by STR_RSOR increases as the problem size increases. The setup times for the ML solvers are all negligible. The solution times of $msmq_c3$ and $msmq_c4$ are too small to say something. However, the ratio of solution times of consecutive problems among $msmq_c5$, $msmq_c6$, and $msmq_c7$ seem to resemble the ratio of the number of states in the corresponding problems. Hence, ML_SOR is clearly scalable in msmq.

Problem	Solver	it	res	Setup	Solve
$msmq_large$	STR_POWER	280	10^{-3}	3	1,004
	STR_JOR	260	10^{-6}	3	1,030
	STR_RSOR	190	10^{-9}	3	902
	STR_BSOR	120	10^{-9}	24	751
	BSOR_BICGSTAB	46	10^{-9}	24	487
	${ m BSOR_TFQMR}$	46	10^{-10}	24	470
	$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, W$	14	10^{-9}	3	152
	$ML_SOR(3,3,3,3,1,1,1,1), CYCLIC, F$	14	10^{-9}	3	150
	$ML_SOR(10,10,10,10,1,1,1,1), CYCLIC, V$	14	10^{-9}	3	171
$courier_large$	STR_POWER	240	10^{0}	3	1,005
	STR_JOR	220	10^{-5}	3	1,017
	STR_RSOR	130	10^{-9}	3	819
	STR_BSOR	48	10^{-7}	21	1,028
	BSOR_BICGSTAB	42	10^{-9}	21	1,020
	$BSOR_TFQMR$	42	10^{-9}	21	992
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, V$	48	10^{-9}	3	483
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, W$	44	10^{-9}	3	480
	$ML_SOR(1,1,1,1,1,1,1,1), CYCLIC, F$	44	10^{-9}	3	474
$qh_real control$	STR_POWER	1,000	10^{0}	0	432
	STR_JOR	1,000	10^{-3}	0	475
	STR_RSOR	1,000	10^{-4}	0	553
	STR_BSOR	1,000	10^{-4}	6	550
	BSOR_BICGSTAB	276	10^{-9}	6	273
	${ m BSOR_TFQMR}$	246	10^{-8}	6	239
	$ML_SOR(5,5,5,5,1,1,1,1), CYCLIC, W$	292	10^{-9}	0	262
	$ ext{ML_SOR}(5,5,5,5,1,1,1,1), \ CYCLIC, \ F$	302	10^{-9}	0	271
	$ML_SOR(10,10,10,10,1,1,1,1), CYCLIC, W$	272	10^{-9}	0	274

Table 6: ML_SOR on benchmark problems

4.3 Performance of ML_SOR on Benchmark Problems

In this subsection we consider three problems. The first two problems named $msmq_large$ and $courier_large$ are respectively larger versions of $msmq_medium$ and $courier_medium$, which are analyzed in subsection 4.1, and have in the order of 1,000,000 states. The third problem named $qh_realcontrol$ is the model of token based scheduling in a queueing network and has in the order of 100,000 states. See Table 2 for characteristics of these problems. The $qh_realcontrol$ and $msmq_large$ problems do not have any non-local transitions along the diagonal of their HLM matrices, whereas $courier_large$ does. Regarding non-local transitions, each LLM in $qh_realcontrol$ and $msmq_large$ respectively participates in four and two transitions. In $courier_large$, LLM 2 and 3 each participate in four transitions while each of the other two LLMs participate in one transition. The $qh_realcontrol$ problem is especially difficult to solve owing it to the existence of nonzeros in its HLM and LLM matrices that have considerably different orders of magnitude.

In msmq_large and courier_large, ML_SOR with CYCLIC order of aggregating LLMs provides clear winners (see Table 6). The best ML solvers in qh_realcontrol also employ the CYCLIC order. The numbers of pre- and post-smoothings at the finest grid performed by the best ML_SOR solvers in all three problems are each one. The numbers of pre- and post-smoothings at the intermediate grids in the winning ML_SOR solvers are three to five.

Regarding the type of cycles, the F cycle provides two winners among the ML_SOR solvers. It is also the cycle of choice in the second best ML_SOR solver for $qh_realcontrol$. Nevertheless, the W cycle appears the most among the nine ML_SOR solvers in Table 6. Note that it takes about 150 cycles to solve $qh_realcontrol$ using ML_SOR. For this more difficult problem, BSOR_TFQMR and BSOR_BICGSTAB provide very strong solvers, which are not easy to beat, but the ML solvers demonstrate similar performance.

5. Conclusion

This paper has proposed a class of ML methods for HMMs with multiple macrostates. The ML solvers are capable of using (V, W, F) cycles, (power, JOR, SOR) methods as smoothers, and (fixed, cyclic, dynamic) orders in which LLMs can be aggregated in each cycle. Extensive numerical experiments on three benchmark HMMs have been performed. Results demonstrate that ML with SOR as the smoother provides the most competitive solver for HMMs with multiple macrostates so far. In almost all cases, three to five preand post-smoothings at intermediate grids and one pre- and post-smoothing at the finest grid of the ML_SOR solver are sufficient to obtain the solution in a small number of ML cycles. Among the three different orders of aggregating LLMs, the cyclic order seems to be favored the most. Regarding cyle type, W or F can be recommended. The storage requirements of the proposed ML solvers are modest and nearly insensitive to the ordering of LLMs in the given HMM description. This is to the contrary of the situation in BSOR preconditioned projection methods. However, as it is observed in one of the test problems, there may be HMMs that are difficult to solve for which BSOR preconditioned projection methods also provide effective solvers.

Acknowledgments

This work has been carried out at Dresden University of Technology, where the second author was a research fellow of the Alexander von Humboldt Foundation.

References

- [1] Ajmone-Marsan, M., Donatelli, S., Neri, F.: GSPN models of Markovian multiserver multiqueue systems, Performance Evaluation 11, 227–240 (1990).
- [2] APNN-Toolbox case studies: http://www4.cs.uni-dortmund.de/APNN-TOOLBOX/case_studies/
- [3] Bause, F., Buchholz, P., Kemper, P.: A toolbox for functional and quantitative analysis of DEDS, In: Puigjaner, R. Savino, N. N., Serra, B., eds., Quantitative Evaluation of Computing and Communication Systems, Lecture Notes in Computer Science 1469, Springer Verlag, 356–359 (1998).
- [4] Buchholz, P.: A class of hierarchical queueing networks and their analysis, Queueing Systems 15, 59–80 (1994).

- [5] Buchholz, P.: Structured analysis approaches for large Markov chains, Applied Numerical Mathematics 31, 375–404 (1999).
- [6] Buchholz, P.: Multilevel solutions for structured Markov chains, SIAM Journal on Matrix Analysis and Applications 22, 342–357 (2000).
- [7] Buchholz, P., Dayar, T.: Block SOR for Kronecker structured representations, Technical Report TUD-FI03-02, Department of Computer Science, Dresden University of Technology, Dresden, Germany (2003). http://www.inf.tu-dresden.de/Fak/berichte.html
- [8] Buchholz, P., Dayar, T.: Block SOR preconditioned projection methods for Kronecker Structured Markovian Representations, Technical Report TUD-FI03-05, Department of Computer Science, Dresden University of Technology, Dresden, Germany (2003). http://www.inf.tu-dresden.de/Fak/berichte.html
- [9] Buchholz, P., Kemper, P.: On generating a hierarchy for GSPN analysis, Performance Evaluation Review 26, 5–14 (1998).
- [10] Buchholz, P., Ciardo, G., Donatelli, S., Kemper, P.: Complexity of memory-efficient Kronecker operations with applications to the solution of Markov models, INFORMS Journal on Computing 12, 203–222 (2000).
- [11] Campos, J., Donatelli, S., Silva, M.: Structured solution of asynchronously communicating stochastic models, IEEE Transactions on Software Engineering 25, 147–165 (1999).
- [12] Dayar, T., Stewart, W. J.: Comparison of partitioning techniques for two-level iterative solvers on large, sparse Markov chains, SIAM Journal on Scientific Computing 21, 1691–1705 (2000).
- [13] Donatelli, S.: Superposed stochastic automata: a class of stochastic Petri nets with parallel solution and distributed state space, Performance Evaluation 18, 21–26 (1993).
- [14] Fernandes, P., Plateau, B., Stewart, W. J.: Efficient descriptor-vector multiplications in stochastic automata networks, Journal of the ACM 45, 381–414 (1998).
- [15] Freund, R. W.: A transpose-free quasi-minimal residual method for non-Hermitian linear systems, SIAM Journal on Scientific Computing 14, 470–482 (1993).
- [16] Hackbusch, W.: Multi-grid Methods and Applications, Berlin: Springer Verlag 1985.
- [17] Horton, G., Leutenegger, S.: A multi-level solution algorithm for steady state Markov chains, Performance Evaluation Review 22, 191–200 (1994).
- [18] Kemper, P.: Numerical analysis of superposed GSPNs, IEEE Transactions on Software Engineering 22, 615–628 (1996).
- [19] Krieger, U. R.: Numerical solution of large finite Markov chains by algebraic multigrid techniques, In: Stewart, W. J., ed., Computations with Markov Chains, Kluwer Academic Publishers, 403–424 (1995).

- [20] Marek, I., Mayer, P.: Convergence analysis of an iterative aggregation/disaggregation method for computing stationary probability vectors of stochastic matrices, Numerical Linear Algebra with Applications 5, 253–274 (1998).
- [21] Plateau, B.: On the stochastic structure of parallelism and synchronization models for distributed algorithms, In: Proceedings of the ACM SIGMETRICS Conference on Measurement and Modelling of Computer Systems, Austin, Texas, 147–154 (1985).
- [22] Plateau, B., Fourneau, J.-M.: A methodology for solving Markov models of parallel systems, Journal of Parallel and Distributed Computing 12, 370–387 (1991).
- [23] Rüde, U.: The Multigrid Workbench. http://www.mgnet.org/mgnet/tutorials/xwb.html
- [24] Stewart, W. J.: Introduction to the Numerical Solution of Markov Chains, Princeton: Princeton University Press 1994.
- [25] Uysal, E., Dayar, T.: Iterative methods based on splittings for stochastic automata networks, European Journal of Operational Research 110, 166–186 (1998).
- [26] van der Vorst, H. A.: BI-CGSTAB: A fast and smoothly converging variant of BI-CG for the solution of nonsymmetric linear systems, SIAM Journal on Scientific and Statistical Computing 13 631–644 (1992) .
- [27] Van Loan, C. F.: The ubiquitious Kronecker product, Journal of Computational and Applied Mathematics 123, 85-100 (2000).
- [28] Wesseling, P.: An Introduction to Multigrid Methods, Chichester: John Wiley & Sons 1992. http://www.mgnet.org/mgnet-books-wesseling.html
- [29] Woodside, C. M., Li, Y., Performance Petri net analysis of communications protocol software by delay equivalent aggregation, In: Proceedings of the 4th International Workshop on Petri Nets and Performance Models, IEEE CS-Press, 64–73 (1991).

Peter Buchholz
Department of Computer Science
Dresden University of Technology
D-01062 Dresden
Germany

e-mail: p.buchholz@inf.tu-dresden.de

Tuğrul Dayar Department of Computer Engineering Bilkent University TR-06800 Bilkent, Ankara Turkey

e-mail: tugrul@cs.bilkent.edu.tr