

# Scheduling With Queue Length Guarantees For Shared Resource Systems

Gagan R. Gupta  
 School of Electrical and Computer Engineering  
 Purdue University  
 West Lafayette, IN, USA  
 grgupta@purdue.edu

Ness B. Shroff  
 Departments of ECE and CSE  
 The Ohio State University  
 Columbus, OH, USA  
 shroff@ece.osu.edu

## ABSTRACT

We develop a class of schemes called GMWM that guarantee optimal throughput for queuing systems with arbitrary constraints on the set of jobs that can be served simultaneously. We obtain an analytical upper bound on the expected queue length. To further tighten the upper bound, we formulate it as a convex optimization problem. We also show that whenever the arrival process is stabilizable, the scheme is guaranteed to achieve an expected queue length that is no larger than the expected queue length of any stationary randomized policy.

**Categories and Subject Descriptors:** C.4 [Performance Of Systems]: Modeling techniques; Performance attributes  
**General Terms:** Performance, Theory.

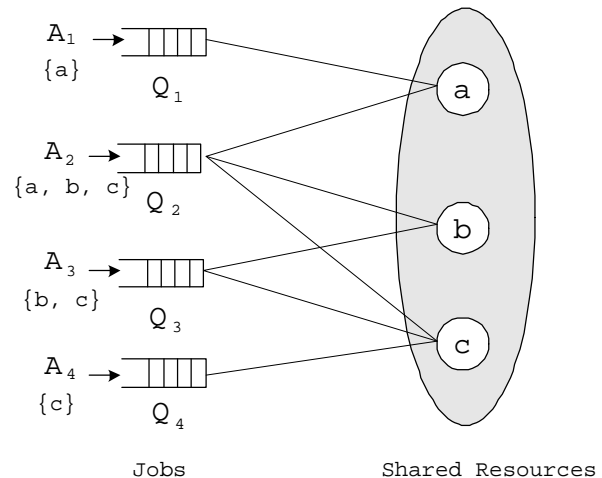
**Keywords:** Scheduling, Lyapunov Theory.

## 1. INTRODUCTION

In modern computer and communication systems, jobs compete with each other to access the limited resources in the system. An example is shown in Figure 1 where jobs are classified and queued according to the set of resources they need to acquire simultaneously for successful processing. Another familiar example is the wireless network where a successful transmission at a given link, necessitates that all interfering link stop transmitting. The system imposes constraints on the set of jobs that can be served simultaneously at any given time. In this work we focus our attention on a basic model, which we call the constrained queuing system (CQS) (proposed by Tassiulas and Ephremides [2]). It is well known that the Maximum Weighted Matching (MWM) scheduling algorithm [2, 4] stabilizes the system for any input load within the capacity region of the system. However, the queuing analysis of such a scheduler is mainly limited to providing stability guarantees and order results. It has been shown in [4] that MWM algorithm is asymptotically delay optimal in the heavy traffic regime, but it is not known if such a result holds for arbitrary load in the system. Further, no bounds on the delay performance have been provided.

## 2. SYSTEM MODEL

We consider a set of  $N$  parallel queues each having its own exogenous *i.i.d.* arrival stream  $\{A_l(t)\}_{t=1}^{\infty}$ . Different input streams may be correlated with each other. Let  $A(t) =$



**Figure 1: Example of a Constrained Queuing System with four classes of jobs. The resources requested by each job are indicated.**

$(A_1(t), \dots, A_N(t))$  represent the vector of exogenous arrivals, where  $A_l(t)$  is the number of jobs that arrive to queue  $l$  during time slot  $t$  (for  $l \in 1, \dots, N$ ). Let  $\lambda = (\lambda_1, \dots, \lambda_N)$  represent the corresponding arrival rate vector. Each job has deterministic service time equal to one unit. Assume that the second moments of the arrival processes  $\mathbf{E}[A_l^2]$  are finite.

The vector of the scheduled queues is denoted by  $\mathbf{I}(t) = (I_n(t)) : n = 1, \dots, N$ . There are constraints on the combination of queues that can be activated. These constraints can be arbitrary.  $\mathbf{I}(t)$  is a valid activation vector if it satisfies the constraints. Let  $S$  be the collection of all activation vectors,  $I^j$ . At each slot an activation vector  $\mathbf{I}(t)$  is scheduled. In this paper we analyze the scheduling policy called the Generalized Weighted Maximum Matching (*GMWM*( $\mathbf{w}$ )) policy described in Figure 2. This describes a class of policies parameterized by the weights  $w_i$ .

Let  $Q_n(t)$  denote the queue length of queue  $n$ . The queue length vector  $\mathbf{Q}(t) = \{Q_n(t) : n = 1, 2, \dots, N\}$ . The queue is activated in a slot  $t$  only if  $Q_n(t) > 0$ . The evolution of the queue is as follows,

$$Q_n(t+1) = (Q_n(t) + A_n(t) - I_n(t))^+, n = 1, \dots, N, \quad (2.1)$$

$$\mathbf{I}(t) = \operatorname{argmax}_{\mathbf{I} \in S} \sum_{i=1}^N w_i Q_i I_i^j \quad (2.2)$$

where  $I_i^j$  is the  $i^{\text{th}}$  component of the  $j^{\text{th}}$  activation vector in set  $S$  and  $w_i > 0$  are fixed constants.

Figure 2: GMWM Scheduling Policy

where

$$(x)^+ = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

The following is a well known result [1], about the existence of a stable stationary randomized scheduling policy for arbitrary load in the capacity region.

LEMMA 2.1. *For any feasible input rate vector  $\lambda = (\lambda_1, \dots, \lambda_N)$  which lies in the interior of the capacity region,  $C$  there exists a vector  $\mu = (\mu_1, \dots, \mu_N) \in C$  such that  $\lambda_l < \mu_l$  for all queues  $l \in 1, \dots, N$ . Also, there exists a stationary randomized scheduling policy,  $\Pi_R$  which chooses activation vectors  $M(t)$  such that  $\mathbf{E}[M_l(t)] = \mu_l$  and hence stabilizes the system.*

### 3. ANALYSIS OF GMWM

Using the Foster-Lyapunov drift criteria for countable Markov chains, it can be shown that the GMWM policy achieves 100% throughput for every choice of  $w$ , such that for all  $i$ ,  $w_i > 0$ .

THEOREM 3.1. *For any input load  $\lambda \in C$ , the GMWM scheduling algorithm ensures that the resulting DTMC is positive recurrent and ergodic.*

PROOF. Straight-forward and omitted for brevity.  $\square$

Using the Lemma 3 from [3] we prove the following result that bounds the sum of expected queue lengths in the system.

THEOREM 3.2. *Given any input load vector  $\lambda \in C$  and any vector  $\mu \in C : \mu > \lambda$ , the following bound on the expectation of the sum of lengths of queues holds true in a system operating under the GMWM policy where the weights  $w_i$  are chosen as  $w_i = \frac{1}{(\mu_i - \lambda_i)}$ :*

$$\sum_{i=1}^N \mathbf{E}[Q_i] \leq \sum_{i=1}^N \frac{(\lambda_i + \mathbf{E}[A_i^2]) - 2\lambda_i^2}{(\mu_i - \lambda_i)}$$

PROOF. Omitted for brevity.  $\square$

The above analysis naturally leads us to the question of which  $\mu > \lambda$  should be selected in the capacity region  $C$  such that the upper bound is minimized. Intuitively, this means that the distance between the load vector and the service process should be as large as possible. We formulate this as an optimization problem to compute the value of  $\mu$  that minimizes the upper bound (see Figure 3).

$$\begin{aligned} & \text{Minimize } \sum_{i=1}^N \frac{(\lambda_i + \mathbf{E}[A_i^2]) - 2\lambda_i^2}{2(\mu_i - \lambda_i)} \\ & \text{subject to } \mu \in C \end{aligned}$$

Figure 3: Optimization Problem for Minimizing the Upper Bound

### 3.1 Comparison with a randomized stationary policy

The stationary randomized scheduler  $\Pi_R$ , is unaware of the backlog and chooses to schedule queue  $l$  independent of whether the queue is empty or not. In every slot, if the queue is scheduled, exactly one job is served, otherwise the jobs in the queue wait for the next available slot. The system evolves as follows

$$q_l(t+1) = (q_l(t) + A_l(t) - M_l(t))^+ \quad (3.3)$$

We can obtain the following result, which we state without proof.

THEOREM 3.3. *Given any admissible arrival process  $\{A_l(t)\}_{t=1}^{\infty}$ , there exists a class of scheduling policies  $GMWM^{opt}$  for which the sum of expected queue lengths  $Q_l$  is no worse than the sum of expected queue lengths  $q_l$  of any other stabilizing stationary randomized policy. In other words,*

$$\sum_{l=1}^N \mathbf{E}[Q_l] \leq \sum_{l=1}^N \mathbf{E}[q_l]$$

## 4. CONCLUSION

There is extensive work on designing schedulers that stabilize the system for maximum throughput. However, there have been only a few results which establish guarantees on the performance metrics such as expected queue lengths, delay, etc. In this paper, we have proposed a class of generalized max weighted schemes, GMWM, and derive analytical bounds on the sum of expected queue lengths in the system. We have also shown that for any given  $\lambda \in C$ , the  $GMWM^{opt}$  is no worse than any stationary randomized scheduling policy.

## 5. REFERENCES

- [1] L. Georgiadis, M. J. Neely, and L. Tassiulas. *Resource Allocation and Cross-Layer Control in Wireless Networks, Foundations and Trends in Networking*, volume 1. Now Publishers, 2006.
- [2] T. Leandros and A. Ephremides. Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks. *IEEE Trans. Aut. Contr.* 37, 37(12):1936–1948, 1992.
- [3] M. J. Neely. Order optimal delay for opportunistic scheduling in multi-user wireless uplinks and downlinks. In *44th Annual Allerton Conference on Communication, Control, and Computing*, September 2006.
- [4] A. L. Stolyar. Maxweight scheduling in a generalized switch: State space collapse and workload minimization in heavy traffic. *Annals of Applied Probability*, Vol.14(No.1):pp.1–53., 2004.