FAST QUEUING POLICIES FOR MULTIMEDIA APPLICATIONS

Anonymous ICME submission

ABSTRACT

We present an analytical framework for providing Quality of Service (QoS) using queuing policies that achieves a given target distribution of packets in a network queue. Queuing policies are employed in an attempt to control the sending and receiving rates subject to the uncertainties in the environments. To a large extent, the stationary distribution of packets in the queue resulted from employing a certain queuing policy directly controls the typical OoS metrics for multimedia applications. Therefore, using the packet distribution in the queue as the metric, the proposed framework allows for a more general and precise control of QoS beyond the standard metrics such as bandwidth, jitter, loss, and delay. Moreover, the proposed framework aims to find a fast queuing policy that achieves a given target stationary distribution. This fast adaptation is especially useful for multimedia applications in fast-changing network conditions. As an example, we present a general procedure for obtaining a queuing policy that optimizes for a given arbitrary objective along with the standard QoS requirements. Both theory and simulation results are presented to verify our framework.

Index Terms— QoS, Queueing theory, Distribution Shaping, Convex Optimization

1. INTRODUCTION

Since the development of packet-switched networks in early 1960s, queuing theory has been a critical part in the performance analysis for most if not all the modern transmission protocols. The performance of TCP/IP protocols for example, can be analyzed in the language of queuing theory [1]. Many current wireless transmission protocols such as the IEEE 802.11 protocols owed their analysis to queuing theory. In fact, queues are so universal that they are virtually found in every communication devices from the core Internet routers and broadband modems to wireless LAN and cellular devices. It is therefore, not a surprise that a point-to-point data flow is typically modeled as a single queue or a network of queues. Understanding the dynamics of packets in queues over time as a result of employing certain queuing policy, enables the system engineers to characterize and to predict various properties of the data flow such as bandwidth, packet loss and delay. With the advent of multimedia communication applications that require certain levels of Quality of Service (QoS), e.g.,

requirements on minimum bandwidth, maximum jitter, delay, or loss, the role of queuing policy is becoming increasingly more important.

Queue and Queuing Policy: In a typical packet-switched network, the instantaneous arrival rates of packets at an intermediate router can vary significantly. Hence packet loss occurs when the arrival rate exceeds the sending rate at a router. Therefore, a queue or a buffer is used to temporarily store a burst of incoming packets in an attempt to prevent packet loss. These packets waits for their turns in the queue to be transmitted to the next hop, or to read by an application if the queue is located at a receiving end device. Queuing policy is a mechanism used to control various operations of a queue that govern the packet's entrance, departure, and drop. It is directly responsible for shaping the dynamics of packets in the queue which characterizes the delay, loss, and bandwidth of a flow. Depending on certain constraints, some queuing polices are more limited in their operations than others. For example, a simple queuing policy is the First In First Out (FIFO) scheme which is typically implemented at the Internet core routers. A router using FIFO policy sends out packets in the order of their arrivals as fast as possible. Packets arriving at the router are dropped when the queue is full. One important observation is that the FIFO has no ability to control the sending or dequeue rate, nor it has the ability to provide feedback to the upstream node for adjusting the incoming or enqueue rate. On the other hand, a more sophisticated queuing policy would be able to control, at least probabilistically, the dequeue rate and the enqueue rate possibly via feedback in order to achieve some given objectives such as queue stability or average queue length.

The well-known Transmission Control Protocol (TCP) is an example of end-to-end flow control in which the feedback (ACK message) to the sender is used to control the sending (enqueue) rate. The IEEE 802.11 protocol family also employs feedback in the form of collisions to adjust the sending rates appropriately. Beyond network protocols, queues are also extensively used in rate control for video coding [2]. The objective of rate control is to produce a coded video bit stream with a certain average bit rate and variance. In this setting, a "conceptual" queue is connected to a video encoder. The feedback from the queue to the video encoder is used by the video encoder to adjust the coded video bit rate using the coding parameters such as quantization level and coding mode appropriately.

In this paper, we consider a general class of queuing polices that allows for the ability to adjust the sending and receiving rates *probabilistically*. The probabilistic framework arises naturally from the unavoidable uncertainties in when and how fast packets arrive due to the fluctuations in network traffic. Furthermore, in some scenarios the ability to send packets out (de-queue) successfully at any time is probabilistic. For example, in a Wi-Fi network, a wireless node might not be able to successfully send out a packet (de-queue) at a certain time slot due to possible collision with other node's transmission. Also, its random back-off mechanism after a collision can in fact be viewed as a dequeuing operation with a certain probability. In this paper, we also limit our discussion to the analysis of queuing policy for a single queue. We believe the analysis for this simple case is still useful since it is applicable to providing QoS in the last mile scenario or single-hop networks such as a Wi-Fi or access networks.

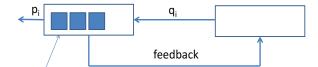
Our contributions include an analytic framework for providing Quality of Service (QoS) using fast a queuing policy that achieves a given target distribution of packets in a network queue. Using the packet distribution in the queue as the metric, the proposed framework allows for a more general and precise control of QoS beyond the standard metrics such as bandwidth, jitter, loss, and delay. The fast adaptive queuing policies are especially useful for multimedia applications in fast-changing network conditions. Finally, we show how even faster queuing policy can be achieved when the solution only need to produces the stationary distribution that is ϵ -close to the given target stationary distribution. Our framework is developed by the theory of fast mixing chain and convex optimization. As an example, we present a general procedure for obtaining a queuing policy that optimizes for a given arbitrary objective along with the standard QoS requirements.

Our paper is organized as follows. In Section 2, we provide some background on the theory of Markov Chain and Queuing as they are necessary for the development of our proposed framework. In Section 3, we describe a convex optimization framework with multiple formulations for finding fast queuing policies. As an example, in Section 4, we show an application of our framework to finding a queuing policy that optimizes for a given objective while ensuring the mean and variance of queuing delay are within given bounds. Finally, we provide a few concluding remarks.

2. PRELIMINARIES

2.1. Queues, Markov Chain, and Mixing Times

Informally, a finite and discrete Markov chain is a set of sequence of random variables $X_1, X_2, ... X_n$ such that given the present states, the past and the future states are independent. A finite state time-homogeneous Markov chain is formally characterized by a time-invariant transition probability matrix P. Let N be the number of the maximum physical



Queuing Policy: Enqueue and dequeue probabilities (rates)

$$\begin{pmatrix} r_1 & q_1 & & & & & & \\ p_2 & r_2 & q_2 & \cdots & & & & \\ & p_3 & r & & & & \vdots & & \\ & \vdots & & \ddots & & \vdots & & \\ & & & r_{n-2} & p_1 & & \\ & & & \cdots & p_{n-1} & r_{n-1} & q_{n-1} \\ & & & & p_{n-1} & r_{n-1} & q_{n-1} \end{pmatrix}$$

Fig. 1. Queuing policy can be viewed as a tridiagonal transition probability matrix

queue length, X_n be the number of packets in the queue at time step n, then dynamics of the number of packets in the queue can be mathematically represented by a Markov chain with a square tridiagonal probability matrix $P^{N\times N}$. Each entry P_{ij} denotes the conditional probability that the chain moves to state j in the next time step given that it is in state i in the current time step. Note that P_{ij} does not depend on time n. Fig. 1 shows a queuing policy that adjusts the sending and receiving rates probabilistically at different states and the corresponding triadiagonal transition probability matrix.

We also consider queuing policies such that the diagonal and off-diagonal entries in the corresponding tridiagonal matrix is different from 1. This ensures that the chain in aperiodic and irreducible. Informally, an aperiodic and irreducible chain has the properties that the chain can reach any state with non-zero probability at some point, and that the time the chain starts in any state i and returns to the same state i must not be a multiple of K>1. This assumption is important as it allows for the characterization of the stationary distribution with the following Proposition:

Proposition 1 For an irreducible, aperiodic, finite and discrete Markov chain with a transition probability matrix P, there exists a unique stationary distribution π such that

$$\lim_{n \to \infty} \nu^T P^n = \pi^T. \tag{1}$$

The stationary distribution approximately represents the probabilities of the chain being in different states after a sufficiently large number of time steps regardless of the initial state of the chain.

In order to quantify "fast" queuing policy, i.e., how fast a queuing policy drive an initial distribution to a given target stationary distribution, it is necessary to define a similarity measure between two distributions. One common similarity measure is the total variance distance defined below: **Definition 1 (Total variation distance)** For any two probability distributions ν and π on a finite state space Ω , we define the total variation distance as:

$$\|\nu - \pi\|_{TV} = \frac{1}{2} \sum_{i \in \Omega} |\nu(i) - \pi(i)|.$$

We now use the similarity measure to define an important notion called mixing time below:

Definition 2 (Mixing time) For a discrete, aperiodic and irreducible Markov chain with transition probability P and stationary distribution π , given an $\epsilon > 0$, the mixing time $t_{mix}(\epsilon)$ is defined as

$$t_{mix}(\epsilon) = \inf \{ n : \| \nu^T P^n - \pi^T \|_{TV} \le \epsilon, \text{for all } probability \text{ distributions } \nu \}.$$

Essentially, the mixing time of a discrete time Markov chain is the minimum number of time step n until the total variance distance between the n-step distribution ad the stationary distribution is less than ϵ . We will use the mixing time to characterize the convergence rate of a queuing policy. One of the successful techniques for bounding the mixing time of a stochastic matrix is via its spectral characterization, i.e., its eigenvalues.

Eigenvalues and Eigenvectors. A non-zero vector v_i is called a right (left) eigenvector of a square matrix P if there is a scalar λ_i such that: $Pv_i = \lambda_i v_i$ or $(v_i^T P = \lambda v_i^T)$. The scalar λ_i is said to be an eigenvalue of P. If P is a stochastic matrix, then $|\lambda_i| \leq 1, \forall i$. Denote the set of eigenvalues in non-increasing order:

$$1 = \lambda_1(P) \ge \lambda_2(P) \ge \dots \ge \lambda_{|\Omega|}(P) \ge -1$$

Definition 3 (Second largest eigenvalue modulus) The second largest eigenvalue modulus (SLEM) of a matrix P is defined as:

$$\mu(P) = \max_{i=2,\dots,|\Omega|} |\lambda_i(P)| = \max\{\lambda_2(P), -\lambda_{|\Omega|}(P)\} \quad (2)$$

In this paper, we also make use the reversibility property of Markov chain defined as follows:

Definition 4 (Reversible Markov Chain) *A* discrete Markov chain with a transition probability *P* is said to be reversible if

$$P_{ij}\pi(i) = P_{ji}\pi(j) \tag{3}$$

We now show an important bound that relates mixing time of the Markov chain to the SLEM of a reversible matrix P.

Theorem 1 (Bound on mixing time) [3]. Let P be the transition matrix of a reversible, irreducible and aperiodic Markov chain with state space Ω , and let $\pi_{min} := \min_{x \in \Omega} \pi(x)$. Then

$$t_{mix}(\epsilon) \le \frac{1}{1 - \mu(P)} \log\left(\frac{1}{\epsilon \pi_{min}}\right).$$
 (4)

It is not difficult to see that from Theorem 1, the error ϵ reduces over time at a rate of no greater than $\frac{e^{-(1-\mu(P))t}}{\pi_{min}}$. Thus, finding the matrix P with minimum $\mu(P)$ would result in the fastest convergence time. Next, we discuss previous results on how to find reversible matrices or queuing policies with fast convergence rates.

2.2. Finding Queuing Policy with Fast Convergence Rate

For a reversible, irreducible, aperiodic chain P with stationary distribution π , it was shown in [4] that

$$\mu(P) = ||D_{\pi}^{1/2} P D_{\pi}^{-1/2} - \sqrt{\pi} (\sqrt{\pi})^T||_2, \tag{5}$$

where D_{π} denotes the square diagonal matrix whose diagonal entries are taken from each elements of π , and ||.|| denote l_2 -induced matrix norm.

Then given a target distribution π^* , it is not difficult to see that $\mu(P)$ is a convex function in P. Thus, in [4], the problem of finding the reversible matrix P with the smallest SLEM, or Fastest Mixing Markov Chain (FMMC) is the following convex optimization:

FMMC framework.

Minimize
$$||D_{\pi^*}^{1/2}PD_{\pi^*}^{-1/2} - \sqrt{\pi^*}(\sqrt{\pi^*})^T||_2$$

Subject to :
$$\begin{cases} P\mathbf{1} = \mathbf{1} \\ D_{\pi^*}P = P^TD_{\pi^*} \end{cases}$$
 (6)

We note that the first constraint guarantee the matrix P to be a valid transition probability matrix, while reversibility is enforce in the second constraint. P is the only optimization variable.

An extension of the FMMC problem is also considered in [5], called the EFMMC problem. In the EFMMC problem the goal is to produce even a faster mixing Markov chain that the one the one obtained by the FMMC. However, the resultling stationary distribution is no longer exactly the given target distribution, but is an ϵ -approximation to the target stationary distribution. Specifically, it have shown that the solution of EFMMC can be obtained using the following convex optimization:

Extended FMMC framework.

Minimize
$$||D_{\pi^*}^{1/2}PD_{\pi^*}^{-1/2} - \sqrt{\pi^*}(\sqrt{\pi^*})^T||_2$$

Subject to :
$$\begin{cases} P\mathbf{1} = \mathbf{1} \\ ||\pi^{*T}P - \pi^{*T}||_2 \le \delta \\ \text{Other convex constraints on } P. \end{cases}$$
(7)

_

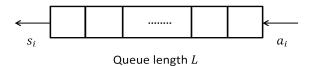


Fig. 2. Discrete queue model

By choosing appropriate value of $\delta(\epsilon)$ as discussed in [5], we guarantee the solution to our EFMMC problem will produce a stationary distribution π that is ϵ - approximation of π^* , specifically $|\pi - \pi^*| \leq \epsilon$.

Now we make the connection to the queuing policy and reversible matrix with the following Proposition:

Proposition 2 Any tridiagonal transition matrix corresponds to a reversible Markov Chain.

Since every queuing policy corresponds to a tridiagonal transition probability matrix, from the Proposition 2 all the queuing policies that we considered are reversible. Also, it is not difficult to add in additional convex constraints to ensure that the solutions of the convex optimization problems above to have the solution as a tridiagonal matrix.

However, it is important to note that for a given tridiagonal transition probability matrix, there might not be a valid queuing policy for specific settings. Therefore, even if one find the fastest reversible matrix for either two convex formulations above, there might not be readily a feasible queueing policy. In the next section, we show a method for find an approximate queuing policy based on the P found in the two formulations.

3. TRIDIAGONAL MATRICES AND QUEUE POLICY

Depending on specific settings, the tridiagonal will not produce a valid queuing policy, more precisely, produce a feasible way for conrolling the enqueue and dequeue rates. Let us consider the following scenario in which the arrival and depature rates at the queue can controlled to some extent by a queueing policy. Let us assume that as a result of the queuing policy, the probabilities of a packet arriving at the queue and departing from the queue when the queue length i are a_i and s_i , respectively. We assume that packets can only arrive and depart in each discrete time slot. The ability to control the arrival rate seems impossible for physical queues in the Internet routers, however, it is frequently implemented in high level network protocols such as TCP in which virtual queues are typically used to provide feedback to the sender for the purpose of rate control. Using this queuing model as shown in Fig. 2, let us denote:

- $|\Omega|$: Maximum queue length
- $s = (s_0, \dots, s_{|\Omega|})$ where $s_0 = 0$: Departing probability vector

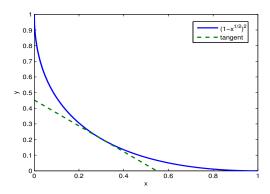


Fig. 3. Tangent at $x_0 = 0.3$ of the function $f(x) = (1 - \sqrt{x})^2$

• $a = (a_0, \dots, a_{|\Omega|})$ where $a_{|\Omega|} = 0$: Arrival probability vector.

Then it is not difficult to see that the dynamics of the number of packets in a queue over time is governed by a discrete Markov chain with the transition probability matrix below:

$$Q = \begin{pmatrix} 1 - a_0 & a_0 \\ s_1(1 - a_1) & 1 - s_1 - a_1 + 2s_1a_1 & (1 - s_1)a_1 \\ & \ddots & \ddots & \ddots \\ & & s_{|\Omega|} & 1 - s_{|\Omega|} \end{pmatrix}$$

Note that for each non-zero entry of each row, the left, middle, and right entries denote the probabilities that the number of packets in the queue decreases by 1, stays the same, or increases by 1, respectively.

Now, let us compare the above matrix Q to P the solution obtained from the problem FMMC (or EFMMC) above. In general, P is tridiagonal matrix with the entries: r_i, q_i, p_i .

$$P = \begin{pmatrix} r_0 & p_0 & & & \\ q_1 & r_1 & p_1 & & & \\ & \ddots & \ddots & \ddots & \\ & & q_{|\Omega|} & r_{|\Omega|} \end{pmatrix}$$
(9)

The main challenge is given r_i, q_i, p_i , how to find the corresponding s_i and a_i , i.e., enqueues and dequeue rates. It turns out that s_i and a_i might be negative or complex numbers which cannot be used in a feasible queuing policy. However, we dermine the conditions on q_i and p_i for which there exist real and non-negative solutions for s_i and a_i , leading to a feasible queuing policy. We proceed to derive the conditions as follows.

From (8) and (9) to compute vector s and vector a, we need to solve these following equations:

$$\begin{cases} s_i(1-a_i) = q_i \to a_i = 1 - q_i/s_i \\ (1-s_i)a_i = p_i \to a_i = p_i/(1-s_i) \end{cases}$$

$$\iff 1 - q_i/s_i = p_i/1 - s_i \text{ for } i = 1, \dots, |\Omega| - 1$$

$$\iff (1-s_i)s_i = (1-s_i)q_i + s_ip_i \text{ for } i = 1, \dots, |\Omega| - 1$$

$$\iff s_i^2 - s_i(1 + q_i - p_i) + q_i = 0 \text{ for } i = 1, \dots, |\Omega| - 1$$
 (10)

In order to guarantee the existence of feasible solution of (10), we need:

$$\Delta = (1 + q_i - p_i)^2 - 4q_i \ge 0 \text{ for } i = 1, \dots, |\Omega| - 1 \quad (11)$$

In theory, we can add these contraints to the two convex formulations above. However, these contraints are not convex constraint making it hard to solve in general. Therefore, our approach is to relax (11) to a convex constraint as follows.

(11)
$$\iff$$
 $1 + q_i - p_i > 2\sqrt{q_i} \text{ since } q_i > 0$ \iff $(1 - \sqrt{q_i})^2 > p_i$ (12)

Consider function $f=(1-\sqrt{x})^2$ for $x\in(0,1)$, we can find an approximate lower bound function of f in the form of tangent y=ax+b where $a=f'(x_0)$ and $f'(x)=(\sqrt{x}-1)/\sqrt{x}$ (See (Fig. 3)).

Hence, (12) is equivalent to the following convex constraints:

$$a(x_0)q_i + b(x_0) > p_i \text{ for } i = 1, \dots, |\Omega| - 1$$
 (13)

Now, we can incorporate these constraints in (13) to the FMMC and/or EFMMC problems, and still have convex formulations to find feasible queuing policies.

4. OPTIMIZING A GIVEN OBJECTIVE VIA QUEUING POLICY

4.1. Approach Illustration

In this section, we provide an example of applying our proposed framework to find fast queuing policy that optimizes a given objective while still satisfying other standard QoS requirements. Our approach consists of two steps. In the first step, we to find a stationary distribution π^* that optimizes a given objective subject to all the given constraints assuming that the given objective and the constraints are convex in π , and thus π^* can be determined efficiently. In the second step, we substitute π^* into either the FMMC or EFMMC with the convex constraints in (13) to find the fastest queuing policy. We give a specific example below.

Step 1. Let X be discrete random variable representing the number of packets in the queue $(X \in [0, \dots, L])$. Suppose a video application requires that the queuing delay average and second moment must be bounded by must be bounded within a range. For example,

$$\left\{ \begin{array}{l} X1 < E[X] < Y1 \\ X2 < E[X^2] < Y2 \end{array} \right.$$

Then E[X] and $E[X^2]$ can be computed from the stationary distribution π :

$$\begin{cases} E[X] = \sum_{x=0}^{L} \pi(x)x \\ E[X^{2}] = \sum_{x=0}^{L} \pi(x)x^{2} \end{cases}$$

Furthermore, suppose that there is a cost function c(x) where x denotes the number of packet in the queue. c(x) could be any arbitrary convex function that might represent energy, resources that depends on the queue occupancy. Now, suppose we want to minimize the total expected cost,

$$T = \sum_{x=0}^{x=L} c(x)\pi(x).$$

Then the optimization problem can be formulated as follows.

Minimize
$$\sum_{x=0}^{x=L} c(x)\pi(x)$$

$$\begin{cases}
X1 < \sum_{x=0}^{L} \pi(x)x < Y1 \\
X2 < \sum_{x=0}^{L} \pi(x)x^{2} < Y2
\end{cases}$$

$$\begin{cases}
\sum_{x=0}^{L} \pi(x) = 1 \\
\pi_{min} < \pi(x) \ \forall x = 0, 1, \dots, L
\end{cases}$$
(14)

Step 2. The solution of (14) gives us the target stationary distribution π^* satisfying the QoS requirements and the given objective. Now, we apply the FMMC and EFMMC formulations to find tridiagonal matrices with fast mixing rates. From these matrices, we can find the dequeuing and dequeuing rates as function of the number of packets in the queue that achieves the target distribution quickly.

4.2. Performance Evaluation

In this section, we present the performance evaluations of our approach using the example above with specific parameters. We assume the maximum physical queue length L=9. To demonstrate the flexibility of our approach, the cost function c(x) is chosen arbitrarily as shown in Fig. 5; X1=0; Y1=5; X2=0; Y2=19; $\pi_{min}=0.01$; Using the approximation method for obtaining a feasible queuing policy in Section 3, we choose the tangent at $x_0=0.2$; we set $\delta=0.001$ in the EFMMC framework.

Fig. 4 shows the shape of the target stationary distribution π^* and π as the results of steps 1 and 2 in Section 4.1, respectively. As seen, π^* and π are very close indicating very good approximation of our approach.

In addition, Fig. 6 shows that the EFMMC framework has a faster convergence rate as expected than the that of FMMC

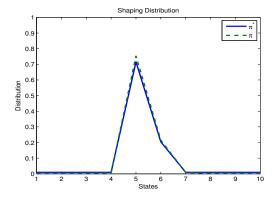


Fig. 4. Target and resulted distribution

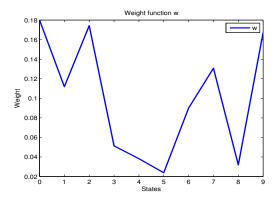


Fig. 5. Cost function c(x)

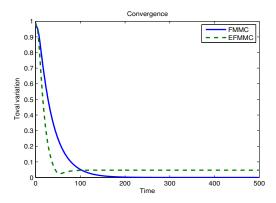


Fig. 6. Comparison of the convergence times in two cases

as expected. Importantly, faster convergence rate is especially useful in non-stationary settings.

To illustrate this point, Fig. 7 shows the total variance between the current distributions produced by the FMMC and EFMMC frameworks, and the target stationary distribution in a non-stationary environment. The non-stationary environment is simulated based on the bursty traffic Poisson patterns with $\lambda=30$. Specifically, in addition to the regular traffic,

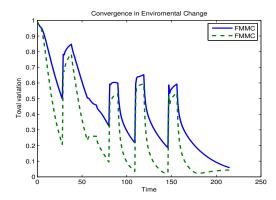


Fig. 7. Convergence of the system during environmental change

there are bursts of 5 packets arriving at the queue. On average, the time duration between these bursts are 30 time slots. As shown in Fig. 7, both curves have spikes when the bursts of packets arrive. This prevents the current distributions in both frameworks from approaching the target stationary distribution (i.e, the curves approaching zero). On the other hand, the queuing policy based on EFMMC framework is better than that of FMMC since it produces as close as possible to the target distribution quickly.

5. CONCLUSION

In this paper, we have proposed a framework for finding fast queuing policies that can provide both flexible QoS requirements as well as optimize for a given objective. The analysis and simulation results show that our framework is especially useful in fast-changing network conditions.

6. REFERENCES

- [1] L. Kleinrock, Queueing systems, volume 1: Theory, John Wiley & Sons, New York, 1975.
- [2] Wei Ding and Bede Liu, "Rate control of mpeg video coding and recording by rate-quantization modeling," *IEEE Trans. Cir. and Sys. for Video Technol.*, vol. 6, no. 1, pp. 12–20, Feb. 1996.
- [3] David A.Levin, Yuval Peres, and Elizabeth L.Wilmer, Markov Chains and Mixing Times, American Mathematical Society, 2008.
- [4] S. Boyd, P. Diaconis, and L. Xiao, "Fastest mixing markov chain on a graph," *SIAM Review*, vol. 46(4), pp. 667–689, December 2004.
- [5] D. Nguyen-Huu, T. Duong, and T. Nguyen, "Achieving quality of service via packet distribution shaping," in *GLOBECOM*, 2012.