

Max Weight Learning Algorithms with Application to Scheduling in Unknown Environments

Michael J. Neely

University of Southern California

<http://www-rcf.usc.edu/~mjneely>

Abstract

We consider a discrete time stochastic queueing system where a controller makes a 2-stage decision every slot. The decision at the first stage reveals a hidden source of randomness with a control-dependent (but unknown) probability distribution. The decision at the second stage incurs a penalty vector that depends on this revealed randomness. The goal is to stabilize all queues and minimize a convex function of the time average penalty vector subject to an additional set of time average penalty constraints. This setting fits a wide class of stochastic optimization problems. This includes problems of opportunistic scheduling in wireless networks, where a 2-stage decision about channel measurement and packet transmission must be made every slot without knowledge of the underlying transmission success probabilities. We develop a simple max-weight algorithm that learns efficient behavior by averaging functionals of previous outcomes. The algorithm yields performance that can be pushed arbitrarily close to optimal, with a tradeoff in convergence time and delay.

Index Terms

Opportunistic scheduling, stochastic optimization, dynamic control, queueing analysis

I. INTRODUCTION

We consider a stochastic queueing system that operates in discrete time with unit timeslots $t \in \{0, 1, 2, \dots\}$. Every slot t , a controller makes a 2-stage control decision that affects queue dynamics and incurs a random penalty vector. Specifically, the controller first chooses an action $k(t)$ from a finite set of K “stage-1” control actions, given by an action set $\mathcal{K} = \{1, \dots, K\}$. After the action $k(t) \in \mathcal{K}$ is chosen, a random vector $\omega(t)$ is revealed, which represents a collection of system parameters for slot t (such as channel states for a wireless system). The

Michael J. Neely is with the Electrical Engineering department at the University of Southern California, Los Angeles, CA (web: <http://www-rcf.usc.edu/~mjneely>).

This material is supported in part by one or more of the following: the DARPA IT-MANET program grant W911NF-07-0028, the NSF grant OCE 0520324, the NSF Career grant CCF-0747525.

random vectors $\omega(t)$ are conditionally i.i.d. with distribution function $F_k(\omega)$ over all slots for which $k(t) = k$, where $F_k(\omega)$ is defined:

$$F_k(\omega) \triangleq Pr[\omega(t) \leq \omega \mid k(t) = k] \text{ for } k \in \mathcal{K} \quad (1)$$

where vector inequality is taken entrywise. However, the distribution functions $F_k(\omega)$ are unknown. Based on knowledge of the revealed $\omega(t)$ vector, the controller makes an additional decision $I(t)$, where $I(t)$ is chosen from some abstract (possibly infinite) set \mathcal{I} . This decision affects the service rates and arrival processes of the queues on slot t , and additionally incurs an M -dimensional *penalty vector* $\mathbf{x}(t) = (x_1(t), \dots, x_M(t))$, where each entry $m \in \{1, \dots, M\}$ is a function of $I(t)$, $k(t)$, and $\omega(t)$ according to known functions $\hat{x}_m(k(t), \omega(t), I(t))$:

$$x_m(t) = \hat{x}_m(k(t), \omega(t), I(t)) \text{ for } m \in \{1, \dots, M\} \quad (2)$$

The penalties can be either positive, zero, or negative (negative penalties can be used to represent *rewards*). Let $\bar{\mathbf{x}}$ be the *time average penalty vector* that results from the control actions made over time (assuming temporarily that this time average is well defined). The goal is to develop a control policy that minimizes a convex function $f(\bar{\mathbf{x}})$ of the time average penalty vector, subject to queue stability and to an additional set of N linear constraints of the type $h_n(\bar{\mathbf{x}}) \leq b_n$ for $n \in \{1, \dots, N\}$, where the constants b_n are given and the functions $h_n(\mathbf{x})$ are linear over $\mathbf{x} \in \mathbb{R}^M$.¹ This objective is similar to the objectives treated in [1] [2] [3] for stochastic network optimization problems, and the problem can be addressed using the techniques given there in the following special cases:

- (Special Case 1) There is no “stage-1” control action $k(t)$, so that the revealed randomness $\omega(t)$ does not depend on any control decision.
- (Special Case 2) The distribution functions $F_k(\omega)$ are known.

An example of Special Case 1 is the problem of minimizing time average power expenditure in a multi-user wireless downlink (or uplink) with random time-varying channel states that are known at the beginning of every slot. Simple max-weight transmission policies are known to solve such problems, even without knowledge of the probability distributions for the channels or packet arrivals [4]. An example of Special Case 2 is the same system with the additional assumption that there is a cost to measuring channels at the beginning of each slot. In this example, we have the option of either measuring the channels (and thus having the hidden random channel states revealed to us) or transmitting blindly. Such a problem is treated in [5], and a related problem with partial channel measurement is treated in [6]. Both [5] and [6] solve the problem via max-weight algorithms that include an expectation with respect to the known joint channel state distribution. While it is reasonable to estimate the joint channel state distribution when channels are independent and/or when the number of channels M is small (and the number of possible states per channel is also small), such estimation becomes intractable in cases when channels are correlated and there are, say, 1024 possible states per channel (and hence there are 1024^M probabilities to be estimated in the joint channel state distribution).

¹For simplicity we treat the case of linear $h_n(\mathbf{x})$ functions here, although the analysis can be extended to treat convex (possibly non-linear) $h_n(\mathbf{x})$ functions, as considered in [1] for the case without “stage 1” control decisions. See also Remark 1 in Section II-D for a further discussion.

Another important example is that of dynamic packet routing and transmission scheduling in a multi-commodity, multi-hop network with probabilistic channel errors and multi-receiver diversity. The Diversity Backpressure Routing (DIVBAR) algorithm of [7] reduces this problem to a 2-stage max-weight problem where each node decides which of the K commodities to transmit at the first stage. After transmission, the random vector of neighbor successes is revealed, and the “stage-2” packet forwarding decision is made. If there is a single commodity ($K = 1$), the problem of maximizing throughput reduces to a problem without “stage-1” decisions, while if there is more than one commodity the solution given in [7] requires knowledge of the joint transmission success probabilities for all neighboring nodes. It is of considerable interest to design a modified algorithm that does not require such probability information.

In this paper, we provide a framework for solving such problems without having a-priori knowledge of the underlying probability distributions. For simplicity, we focus primarily on 1-hop networks, although the techniques extend to multi-hop networks using the techniques of [1] [8]. Our approach uses the observation that, rather than requiring an estimate of the full probability distributions, all that is needed is an estimate of a set of expected *max-weight functionals* that depend on these distributions. These can be efficiently estimated using penalties incurred on previous transmissions to learn optimal behavior.

Related stochastic network optimization problems (without the 2-stage decision and learning component) appear in [9] [1] [3] [2]. Work in [9] considers optimization of a utility function of time average throughput in an opportunistic scheduling scenario but without queues or stability constraints. Work in [1] [3] treats joint queue stability and performance optimization using Lyapunov optimization, and work in [2] treats similar problems in a fluid limit sense using primal-dual methods. Sequential channel probing techniques via dynamic programming are treated in [10] [11] [12]. General methods for Q-learning, based on approximate dynamic programming, are presented in [13]. Our approach is different and is based on simpler Lyapunov optimization techniques, which, due to the special structure of the problem, provide strong (polynomial) bounds on convergence even for high dimensional state spaces. Simple methods of pursuit learning and reinforcement learning, which try to converge to the repeated selection of an optimal single index that provides a maximum mean reward (without a-priori knowledge of the average rewards for each index), are considered in [14] and applied to wireless rate selection in [15]. Our stage-1 decision options can be viewed as a finite set of indices, and hence our problem is related to [14] [15]. However, our 2-stage problem structure and the underlying stochastic queues, convex cost optimization, and multi-dimensional inequality constraints, make our problem much more complex. Further, the optimal policy may (and typically does) result in a probabilistic mixture of many different action modes, rather than a single fixed action.

II. THE MAX WEIGHT LEARNING PROBLEM

Consider a collection of L discrete time queues $\mathbf{Q}(t) = (Q_1(t), \dots, Q_L(t))$ with dynamic equation:

$$Q_l(t+1) = \max[Q_l(t) - \mu_l(t), 0] + A_l(t) \quad (3)$$

where $A_l(t)$ is the amount of new arrivals to queue l on slot t , and $\mu_l(t)$ is the queue l server rate on slot t . These quantities are possibly affected by the two-stage control decision at slot t . Specifically, let $\mathcal{K} \triangleq \{1, \dots, K\}$ represent

the set of stage-1 decision options, and let $k(t)$ represent the stage-1 decision made by the controller at time t , for $t \in \{0, 1, 2, \dots\}$. Recall that the corresponding random vector $\boldsymbol{\omega}(t)$ that is revealed is conditionally i.i.d. over all slots for which $k(t) = k$, with distribution function $F_k(\boldsymbol{\omega})$ given by (1). The $F_k(\boldsymbol{\omega})$ distributions are unknown to the controller. Let \mathcal{I} be the (possibly infinite) set of stage-2 control actions, and let $I(t) \in \mathcal{I}$ denote the stage-2 control action at time t .

The arrival and service vectors $\mathbf{A}(t) = (A_1(t), \dots, A_L(t))$ and $\boldsymbol{\mu}(t) = (\mu_1(t), \dots, \mu_L(t))$ are determined by $k(t)$, $\boldsymbol{\omega}(t)$, $I(t)$ according to (known) functions $\hat{a}_l(k(t), \boldsymbol{\omega}(t), I(t))$ and $\hat{\mu}_l(k(t), \boldsymbol{\omega}(t), I(t))$:²

$$A_l(t) = \hat{a}_l(k(t), \boldsymbol{\omega}(t), I(t))$$

$$\mu_l(t) = \hat{\mu}_l(k(t), \boldsymbol{\omega}(t), I(t))$$

Likewise, the penalty vector $\mathbf{x}(t) = (x_1(t), \dots, x_M(t))$ is determined by the (known) penalty functions $x_m(t) = \hat{x}_m(k(t), \boldsymbol{\omega}(t), I(t))$ for each $m \in \{1, \dots, M\}$. The penalties are (possibly negative) real numbers, and we assume that the penalty functions are bounded by finite constants x_m^{\min} and x_m^{\max} for all $m \in \{1, \dots, M\}$, so that:

$$x_m^{\min} \leq x_m(t) \leq x_m^{\max} \text{ for all } t$$

Likewise, the queue arrivals and service rates are bounded as follows:

$$0 \leq A_l(t) \leq A_l^{\max} \text{ for all } t$$

$$0 \leq \mu_l(t) \leq \mu_l^{\max} \text{ for all } t$$

Aside from this boundedness, the functions $\hat{a}_l(\cdot)$, $\hat{\mu}_l(\cdot)$, and $\hat{x}_m(\cdot)$ are otherwise arbitrary (possibly nonlinear, non-convex, and discontinuous). Define the time average penalty $\bar{\mathbf{x}}(t)$, averaged over the first t slots, as follows:

$$\bar{\mathbf{x}}(t) \triangleq \frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E} \{ \mathbf{x}(\tau) \}$$

Let $f(\mathbf{x})$ be a convex and continuous function over $\mathbf{x} \in \mathbb{R}^M$ (possibly negative, non-monotonic, and non-differentiable). Let $h_n(\mathbf{x})$ for $n \in \{1, \dots, N\}$ be a collection of linear functions over $\mathbf{x} \in \mathbb{R}^M$. Note that since the $\mathbf{x}(t)$ penalties are bounded, the values of $f(\mathbf{x}(t))$ and $h_n(\mathbf{x}(t))$ are also bounded. The goal is to design a control

²The analysis is the same if $\hat{a}_l(\cdot)$, $\hat{\mu}_l(\cdot)$, $\hat{x}_m(\cdot)$ outcomes are random but i.i.d. given $k(t)$, $\boldsymbol{\omega}(t)$, $I(t)$, with known means $\bar{a}_l(\cdot)$, $\bar{\mu}_l(\cdot)$, $\bar{x}_m(\cdot)$ that are used in the decision making part of the algorithm.

policy that makes 2-stage decisions over time so as to solve the following problem:³

$$\text{Minimize:} \quad \limsup_{t \rightarrow \infty} f(\bar{\mathbf{x}}(t)) \quad (4)$$

$$\text{Subject to:} \quad \limsup_{t \rightarrow \infty} h_n(\bar{\mathbf{x}}(t)) \leq b_n \text{ for } n \in \{1, \dots, N\} \quad (5)$$

$$\text{Stability of all queues } Q_1(t), \dots, Q_L(t) \quad (6)$$

In cases when the time average penalty vector converges to some value $\bar{\mathbf{x}}$, the \limsup is equal to the regular limit and the above problem can be more simply stated as minimizing $f(\bar{\mathbf{x}})$ subject to $h_n(\bar{\mathbf{x}}) \leq b_n$ for all $n \in \{1, \dots, N\}$ and to stability of all queues. The following notion of queue stability is used:

Definition 1: A discrete time queue is *strongly stable* if:

$$\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E} \{|Q(\tau)|\} < \infty$$

We shall use the term *stability* throughout to refer to strong stability. The definition above uses the absolute value of queue size because we shall soon introduce additional virtual queues that can take negative values.

A. Auxiliary Variables for Nonlinear Cost Functions

It is useful to write the cost function $f(\mathbf{x})$ as a sum of linear (or affine) and non-linear components. Specifically, define $\tilde{\mathcal{M}}$ as the set of all indices $m \in \{1, \dots, M\}$ for which there are penalty variables $x_m(t)$ that participate in a *non-linear component* of $f(\mathbf{x})$. Then we can write $f(\mathbf{x})$ as follows:

$$f(\mathbf{x}) = l(\mathbf{x}) + \tilde{f}(\tilde{\mathbf{x}})$$

where $l(\mathbf{x})$ is a linear (or affine) function, $\tilde{\mathbf{x}} = (x_m)_{m \in \tilde{\mathcal{M}}}$ is a “sub-vector” of \mathbf{x} that contains only entries x_m for $m \in \tilde{\mathcal{M}}$, and $\tilde{f}(\tilde{\mathbf{x}})$ are convex functions (and typically non-linear). Such a decomposition is always possible, and in principle we can choose the trivial decomposition $\tilde{\mathcal{M}} = \{1, \dots, M\}$, $l(\mathbf{x}) = 0$, $\tilde{\mathbf{x}} = \mathbf{x}$, which does not attempt to exploit linearity even if it exists in the cost function. However, it is useful to separate out the linear components, because we shall require one *auxiliary variable* $\gamma_m(t)$ for each penalty $x_m(t)$ that participates in a non-linear component of a cost function, while no such auxiliary variable is required for penalties that do not participate in any non-linear components.⁴

For each $m \in \tilde{\mathcal{M}}$, let $\gamma_m(t)$ be a new variable that can be chosen as desired on each timeslot t , subject only to the constraint that:

$$x_m^{min} - \sigma \leq \gamma_m(t) \leq x_m^{max} + \sigma \text{ for all } m \in \tilde{\mathcal{M}} \quad (7)$$

³While we assume the objective function $f(\mathbf{x})$ is a general convex (possibly non-linear) function, for simplicity we assume the cost functions $h_n(\mathbf{x})$ are linear (see Remark 1 in Section II-D for extensions to non-linear $h_n(\mathbf{x})$ functions). Example linear constraints for a wireless system are *average power constraints* at each node, where $h_n(\mathbf{x})$ is a linear function that sums the relevant components of the penalty vector $\mathbf{x}(t)$ that correspond to instantaneous power expenditure at node n , and b_n represents the average power constraint of node n . A typical non-linear objective for networks is the maximization of a concave utility function $g(\mathbf{x})$ of the time average throughput, where $g(\mathbf{x})$ selects only those entries x_m that correspond to throughput, and $f(\mathbf{x}) = -g(\mathbf{x})$.

⁴While it is possible to always define one auxiliary variable per penalty, exploiting linearity and reducing the number of auxiliary variables can be more direct and may lead to faster convergence times.

for some positive value $\sigma > 0$ (to be chosen later). Let $\gamma(t) = (\gamma_m(t))|_{m \in \tilde{\mathcal{M}}}$ be a vector of $\gamma_m(t)$ components for $m \in \tilde{\mathcal{M}}$. Define the time average $\bar{\gamma}(t)$ as follows:

$$\bar{\gamma}(t) \triangleq \frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E} \{ \gamma(\tau) \}$$

Then it is not difficult to show that the problem (4)-(6) is equivalent to the following:

$$\text{Minimize:} \quad \limsup_{t \rightarrow \infty} [l(\bar{\mathbf{x}}(t)) + \tilde{f}(\bar{\gamma}(t))] \quad (8)$$

$$\text{Subject to:} \quad \limsup_{t \rightarrow \infty} h_n(\bar{\mathbf{x}}(t)) \leq b_n \quad \text{for } n \in \{1, \dots, N\} \quad (9)$$

$$\lim_{t \rightarrow \infty} [\bar{x}_m(t) - \bar{\gamma}_m(t)] = 0 \quad \text{for } m \in \tilde{\mathcal{M}} \quad (10)$$

$$\text{Stability of all queues } Q_1(t), \dots, Q_L(t) \quad (11)$$

Indeed, the equality constraint (10) indicates that the auxiliary variable $\gamma_m(t)$ can be used as a proxy for $x_m(t)$ for all $m \in \tilde{\mathcal{M}}$, so that the above problem is equivalent to (4)-(6). This is useful for stochastic optimization because $\gamma_m(t)$ can be chosen deterministically as any real number that satisfies (7), whereas the penalty $x_m(t)$ has random outcomes. These auxiliary variables are similar to those introduced in [3] [1] for optimizing a convex and non-linear function of a time average penalty in a stochastic network, which is a more general (and more complex) problem than that of optimizing a time average of a non-linear penalty function. In the special case when the objective function $f(\mathbf{x})$ is itself linear (so that $\tilde{f}(\mathbf{x}) = 0$ and $f(\mathbf{x}) = l(\mathbf{x})$), then no auxiliary variables are needed, the set $\tilde{\mathcal{M}}$ is empty, and the constraints (10) are irrelevant.

B. Virtual Queues for Time Average Inequalities and Equalities

To satisfy the time average inequality constraints in (9), we define one *virtual queue* $U_n(t)$ for each $n \in \{1, \dots, N\}$, with dynamic queueing equation:

$$U_n(t+1) = \max[U_n(t) + h_n(\mathbf{x}(t)) - b_n, 0] \quad (12)$$

This can be viewed as a discrete time queueing system with a constant “service rate” b_n and with arrivals $h_n(\mathbf{x}(t))$, although we note in this case that the “arrivals” and/or the “service rate” can potentially be negative on a given slot t . The intuition is that stabilizing this virtual queue ensures that the time average “arrival rate” is less than or equal to b_n . This is similar to the virtual queues used for average power constraints in [4] and average penalty constraints in [1].

To satisfy the time average equality constraints in (10), we introduce a *generalized virtual queue* $Z_m(t)$ for each $m \in \tilde{\mathcal{M}}$, with dynamic equation:

$$Z_m(t+1) = Z_m(t) - \gamma_m(t) + x_m(t) \quad (13)$$

This has a different structure because it enforces an equality constraint, and it can be either positive or negative. The following lemma shows that stabilizing the queues $U_n(t)$ and $Z_m(t)$ ensures that the corresponding inequality and equality constraints are satisfied.

Lemma 1: (Queue Stability Lemma) If the queues $U_n(t)$ and $Z_m(t)$ satisfy the following (for all $n \in \{1, \dots, N\}$ and $m \in \tilde{\mathcal{M}}$):

$$\lim_{t \rightarrow \infty} \frac{\mathbb{E}\{U_n(t)\}}{t} = 0, \quad \lim_{t \rightarrow \infty} \frac{\mathbb{E}\{|Z_m(t)|\}}{t} = 0 \quad (14)$$

Then all inequality constraints (9) and (10) are satisfied. Further, the condition (14) holds whenever the queues are strongly stable.

Proof: Omitted for brevity (see [4] for a related proof). \square

C. Lyapunov Functions

Define $\Theta(t) \triangleq [Q(t); U(t); Z(t)]$ as the vector of all actual and virtual queue backlogs. To stabilize the queues, we define the following Lyapunov function:

$$L(\Theta(t)) \triangleq \frac{1}{2} \sum_{l=1}^L Q_l(t)^2 + \frac{1}{2} \sum_{n=1}^N U_n(t)^2 + \frac{1}{2} \sum_{m \in \tilde{\mathcal{M}}} Z_m(t)^2$$

Note that this Lyapunov function grows large when the absolute value of queue size is large, and hence keeping this function small also maintains stable queues. Define the *one-step conditional Lyapunov drift* as follows:⁵

$$\Delta(\Theta(t)) \triangleq \mathbb{E}\{L(\Theta(t+1)) - L(\Theta(t)) \mid \Theta(t)\} \quad (15)$$

Let V be a non-negative parameter used to control the proximity of our algorithm to the optimal solution of (8)-(11). Using the framework of [1], we consider a control policy that observes the queue backlogs $\Theta(t)$ and takes control actions on each slot t that minimize a bound on the following “drift plus penalty” expression:

$$\Delta(\Theta(t)) + \mathbb{E}\{Vl(\mathbf{x}(t)) + V\tilde{f}(\gamma(t)) \mid \Theta(t)\}$$

Computing the Lyapunov drift $\Delta(\Theta(t))$ by squaring the queueing update equations (12), (13), (3) and taking conditional expectations leads to the following lemma.

Lemma 2: (The RHS(\cdot) Bound) For a general control policy we have:

$$\begin{aligned} \Delta(\Theta(t)) + \mathbb{E}\{Vl(\mathbf{x}(t)) + V\tilde{f}(\gamma(t)) \mid \Theta(t)\} &\leq B \\ &+ \mathbb{E}\{Vl(\mathbf{x}(t)) + V\tilde{f}(\gamma(t)) \mid \Theta(t)\} \\ &- \sum_{n=1}^N U_n(t) \mathbb{E}\{b_n - h_n(\mathbf{x}(t)) \mid \Theta(t)\} \\ &- \sum_{m \in \tilde{\mathcal{M}}} Z_m(t) \mathbb{E}\{\gamma_m(t) - x_m(t) \mid \Theta(t)\} \\ &- \sum_{l=1}^L Q_l(t) \mathbb{E}\{\mu_l(t) - A_l(t) \mid \Theta(t)\} \end{aligned} \quad (16)$$

⁵Strictly speaking, notation should be $\Delta(\Theta(t), t)$, as the drift may be non-stationary. However, we use the simpler notation $\Delta(\Theta(t))$ as a formal representation of the right hand side of (15).

where B is a finite constant that satisfies the following for all t and all possible control actions that can be taken on slot t :

$$\begin{aligned} B \geq & \sum_{n=1}^N \mathbb{E} \{ (b_n - h_n(\mathbf{x}(t)))^2 \mid \Theta(t) \} \\ & + \sum_{m \in \bar{M}} \mathbb{E} \{ (\gamma_m(t) - x_m(t))^2 \mid \Theta(t) \} \\ & + \sum_{l=1}^L \mathbb{E} \{ (\mu_l(t) - A_l(t))^2 \mid \Theta(t) \} \end{aligned}$$

Such a constant B exists because of the boundedness assumptions of the penalty and cost functions, and an explicit bound can be determined by considering the maximum squared values attained by the penalties and costs.

Proof: The proof is a straightforward drift computation (see, for example, [1]), and is omitted for brevity. \square

The next section analyzes the performance of policies that choose control actions every slot to (approximately) minimize the right hand side of the drift expression (16).

D. The Performance Theorem

Define f^* as the optimal solution for the problem (4)-(6) (i.e., it is the infimum cost over all policies that satisfy the constraints). Define a value θ such that $0 \leq \theta < 1$, and consider the class of restricted policies that have random *exploration events* independently with probability θ every slot. If a given slot t is an exploration event, the stage-1 decision $k(t)$ is chosen independently and uniformly over $\{1, \dots, K\}$ (regardless of the state of the system at this time). We say that the slot is an *exploration event of type k* if the exploration event leads to the random choice of option k . Hence, exploration events of type k occur independently with probability θ/K every slot. We note that the stage-2 decision $I(t)$ and the auxiliary variables $\gamma(t)$ can be chosen arbitrarily on every slot, regardless of whether or not the slot is an exploration event.

If $\theta > 0$, the exploration events ensure that each stage-1 control option is tested infinitely often. Define f_θ^* as the optimal solution of (4)-(6) subject to the additional constraint that such random exploration events are imposed. It shall be convenient to define optimality in terms of f_θ^* . It is clear that $f_0^* = f^*$, and intuitively one expects that $f_\theta^* \rightarrow f^*$ as $\theta \rightarrow 0$.⁶ Further, in systems where the optimal f^* can be achieved by a policy that chooses each stage-1 control option a positive fraction of time, it can be shown that there exists a positive value θ^* such that $f^* = f_\theta^*$ whenever $0 \leq \theta \leq \theta^*$. We now assume the following properties hold concerning stationary and randomized control policies with random exploration events of probability θ .

Assumption 1 (Feasibility): There is a stationary and randomized policy that chooses a stage-1 control action $k^*(t) \in \mathcal{K}$ according to a fixed probability distribution such that each option is chosen with probability at least θ/K (revealing a corresponding random vector $\omega^*(t)$), and chooses a stage-2 control action $I^*(t) \in \mathcal{I}$ as a potentially

⁶Specifically, it can be shown that $f_\theta^* \rightarrow f^*$ whenever $\epsilon_{max} > 0$, where ϵ_{max} is defined in Assumption 2.

randomized function of $\omega^*(t)$, such that:

$$l(\mathbb{E}\{\mathbf{x}^*(t)\}) + \tilde{f}(\gamma^*) = f_\theta^* \quad (17)$$

$$b_n - h_n(\mathbb{E}\{\mathbf{x}^*(t)\}) \geq 0 \text{ for all } n \in \{1, \dots, N\} \quad (18)$$

$$\mathbb{E}\{\mu_l^*(t)\} - \mathbb{E}\{A_l^*(t)\} \geq 0 \text{ for all } l \in \{1, \dots, L\} \quad (19)$$

where $\mathbf{x}^*(t)$, $\boldsymbol{\mu}^*(t)$, $\mathbf{A}^*(t)$ are the penalty, service rate, and arrival vectors corresponding to the stationary and randomized policy, defined by:

$$\mathbf{x}^*(t) = \hat{x}(k^*(t), \omega^*(t), I^*(t))$$

$$\boldsymbol{\mu}^*(t) = \hat{\mu}(k^*(t), \omega^*(t), I^*(t))$$

$$\mathbf{A}^*(t) = \hat{a}(k^*(t), \omega^*(t), I^*(t))$$

and where γ^* is a vector with components $(\gamma_m^*)|_{m \in \tilde{\mathcal{M}}}$ such that $\gamma_m^* \triangleq \mathbb{E}\{x_m^*(t)\}$ for all $m \in \tilde{\mathcal{M}}$. Note that $x_m^{\min} \leq x_m(t) \leq x_m^{\max}$ always, and so $x_m^{\min} \leq \gamma_m^* \leq x_m^{\max}$ for all $m \in \tilde{\mathcal{M}}$. Thus, each component γ_m^* satisfies the required auxiliary variable constraint (7).

This assumption states that the problem is feasible, and that the optimal f_θ^* value can be achieved by a particular stationary and randomized policy that meets the time average penalty constraints and ensures the time average service rate is greater than or equal to the time average arrival rate in all queues.⁷ The next assumption states that the constraints are not only feasible, but have a useful slackness property.

Assumption 2 (Slackness of Constraints): There is a value $\epsilon_{max} > 0$ together with a stationary and randomized policy that makes stage-1 and stage-2 control decisions $k'(t) \in \mathcal{K}$ and $I'(t) \in \mathcal{I}$ such that each stage-1 option is chosen with probability at least θ/K , and:

$$b_n - h_n(\mathbb{E}\{\mathbf{x}'(t)\}) \geq \epsilon_{max} \text{ for all } n \in \{1, \dots, N\} \quad (20)$$

$$\mathbb{E}\{\mu_l'(t)\} - \mathbb{E}\{A_l'(t)\} \geq \epsilon_{max} \text{ for all } l \in \{1, \dots, L\} \quad (21)$$

where $\mathbf{x}'(t)$, $\boldsymbol{\mu}'(t)$, $\mathbf{A}'(t)$ are the penalty, service rate, and arrival vectors corresponding to the decisions $k'(t)$ and $I'(t)$.

Now define $RHS(t, \Theta(t), k(t), I(t), \gamma(t))$ as the right hand side of the drift bound (16) with a given queue state $\Theta(t)$ and control actions $k(t)$, $I(t)$, $\gamma(t)$ at time t . Given a particular queue state $\Theta(t)$, define the *max-weight* control decisions $k^{mw}(t)$, $I^{mw}(t)$, $\gamma^{mw}(t)$ as the ones that minimize the following conditional expectation over all alternative feasible control actions that can be made on slot t (subject to the θ exploration probability):⁸

$$\mathbb{E}\{RHS(t, \Theta(t), k(t), I(t), \gamma(t)) \mid \Theta(t)\} \quad (22)$$

⁷See [4] for a proof that optimality can be defined over the class of stationary, randomized policies for minimum power problems.

⁸For simplicity, we implicitly assume that the infimum of (22) over all feasible control actions is achieved by a particular set of decisions, called the max-weight decisions. Else, the results can be recovered by defining the max-weight decisions according to a sequence of policies that converge to the infimum.

Note that the $k^{mw}(t)$ decisions are still determined randomly in the case of exploration events of probability θ , but are chosen to maximize the above expression whenever the current slot does not have an exploration event.

The auxiliary vector $\gamma(t)$ appears in separable terms on the right hand side of (16), and so the policy $\gamma^{mw}(t)$ can be determined separately from the $k^{mw}(t)$ and $I^{mw}(t)$ decisions. It is computed by first observing the queue backlogs $Z_m(t)$ on each slot t , and choosing $\gamma^{mw}(t)$ as the solution to the following deterministic convex optimization:

$$\text{Minimize:} \quad V\tilde{f}(\gamma(t)) - \sum_{m \in \tilde{\mathcal{M}}} Z_m(t)\gamma_m(t) \quad (23)$$

$$\text{Subject to:} \quad x_m^{min} - \sigma \leq \gamma_m(t) \leq x_m^{max} + \sigma \text{ for all } m \in \tilde{\mathcal{M}} \quad (24)$$

If the non-linear function $\tilde{f}(\gamma)$ is separable in the γ vector (as is the case in many network optimization problems), the above optimization amounts to separately finding $\gamma_m^{mw}(t)$ (for each $m \in \tilde{\mathcal{M}}$) as the minimum of a convex single-variable function over the closed interval defined by (24).

While the $\gamma^{mw}(t)$ can thus be computed, it is more challenging to determine the stage-1 and stage-2 decisions that minimize the right hand side of (16), as this would require knowledge of the probability distributions $F_k(\omega)$. We thus seek an *approximation* to the $k^{mw}(t)$ and $I^{mw}(t)$ policies. Suppose the following additional assumption holds concerning such an approximation.

Assumption 3 (Approximate Scheduling): Every slot t the queue backlogs $\Theta(t)$ are observed and control decisions $k(t) \in \mathcal{K}$ (subject to exploration events with probability θ), $I(t) \in \mathcal{I}$, and $\gamma(t)$ satisfying (7) are made to ensure the following:

$$\begin{aligned} \mathbb{E} \{RHS(t, \Theta(t), k(t), I(t), \gamma(t))\} &\leq \mathbb{E} \{RHS(t, \Theta(t), k^{mw}(t), I^{mw}(t), \gamma^{mw}(t))\} \\ &+ C + V\epsilon_V \\ &+ \sum_{n=1}^N \mathbb{E} \{U_n(t)\} \epsilon_U + \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E} \{|Z_m(t)|\} \epsilon_Z + \sum_{l=1}^L \mathbb{E} \{Q_l(t)\} \epsilon_Q \end{aligned} \quad (25)$$

where C , ϵ_V , ϵ_U , ϵ_Z , ϵ_Q are non-negative constants (independent of t). The expectation on the left hand side is with respect to the current queue state $\Theta(t)$ and the actual decisions $k(t)$, $I(t)$, $\gamma(t)$ implemented, while the expectation on the right is with respect to the current queue state $\Theta(t)$ and the (possibly not implemented) max-weight decisions $k^{mw}(t)$, $I^{mw}(t)$, $\gamma^{mw}(t)$ that minimize the right hand side of (16).

We note that the structure of the approximation bound in (25) is typical for algorithms that attempt to select a control action based on imperfect knowledge of the probability distributions of the resulting $\mathbf{x}(t)$, $\boldsymbol{\mu}(t)$, $\mathbf{A}(t)$ vectors, as the resulting approximations are typically proportional to the V constant and the $U_n(t)$, $|Z_m(t)|$, and $Q_l(t)$ queue sizes on the right hand side of (16). In the case of perfect implementation of the max-weight policy $k^{mw}(t)$, $I^{mw}(t)$, $\gamma^{mw}(t)$, we have $\epsilon_V = \epsilon_U = \epsilon_Z = \epsilon_Q = 0$ and $C = 0$.

Theorem 1: (Performance Theorem) Suppose Assumptions 1 and 2 hold, and that a control algorithm is implemented that satisfies Assumption 3 with fixed control parameters $V \geq 0$ and $\sigma > 0$. Suppose ϵ_Q , ϵ_Z , ϵ_U are small enough and σ is chosen large enough to satisfy the following:

$$\epsilon_U < \epsilon_{max}, \epsilon_Z < \sigma, \epsilon_Q < \epsilon_{max} \quad (26)$$

Then all time average constraints (9)-(11) hold. In particular, all queues are strongly stable and satisfy for all t :

$$\frac{1}{t} \sum_{\tau=0}^{t-1} \left[\sum_{n=1}^N \mathbb{E} \{U_n(\tau)\} + \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E} \{|Z_m(\tau)|\} + \sum_{l=1}^L \mathbb{E} \{Q_l(\tau)\} \right] \leq \frac{B + C + V(l_{diff} + \tilde{f}_{diff} + \epsilon_V)}{\epsilon_{approx}} + \frac{\mathbb{E} \{L(\Theta(0))\}}{\epsilon_{approx} t} \quad (27)$$

where ϵ_{approx} is defined:

$$\epsilon_{approx} \triangleq \min[\epsilon_{max} - \epsilon_U, \epsilon_{max} - \epsilon_Q, \sigma - \epsilon_Z]$$

and where l_{diff} and \tilde{f}_{diff} are finite bounds that satisfy:

$$\begin{aligned} l(\mathbf{x}_1) - l(\mathbf{x}_2) &\leq l_{diff} && \text{for any } \mathbf{x}_1, \mathbf{x}_2 \text{ in the set } \mathbf{x}^{min} \leq \mathbf{x} \leq \mathbf{x}^{max} \\ \tilde{f}(\gamma_1) - \tilde{f}(\gamma_2) &\leq \tilde{f}_{diff} && \text{for any } \gamma_1, \gamma_2 \text{ in the set } \mathbf{x}^{min} - \sigma \leq \gamma \leq \mathbf{x}^{max} + \sigma \end{aligned}$$

Further, the time average cost satisfies:⁹

$$\limsup_{t \rightarrow \infty} f(\bar{\mathbf{x}}(t)) \leq f_{\theta}^* + \epsilon_V + \delta + (B + C)/V \quad (28)$$

where we recall that $f(\mathbf{x}) = l(\mathbf{x}) + \tilde{f}(\tilde{\mathbf{x}})$ and f_{θ}^* is the optimal solution of (8)-(11) subject to exploration events with probability θ , and where δ is defined:

$$\delta \triangleq (l_{diff} + \tilde{f}_{diff}) \max \left[\frac{\epsilon_U}{\epsilon_{max}}, \frac{\epsilon_Z}{\sigma}, \frac{\epsilon_Q}{\epsilon_{max}} \right]$$

Theorem 1 states that, under the given approximation assumptions, the algorithm stabilizes all queues and yields a time average cost that is within $\epsilon_V + \delta + O(1/V)$ of the optimal value f_{θ}^* . Hence, this bound can be made arbitrarily close to $f_{\theta}^* + \epsilon_V + \delta$ by choosing V suitably large, at the cost of a linear increase in average queue congestion with V . Further, we note that the terms ϵ_V and δ tend to zero as the error values ϵ_V , ϵ_U , ϵ_Z , ϵ_Q tend to zero. In the special case when the exact max-weight policy is implemented every slot (so that every slot t the controller makes decisions $k^{mw}(t)$, $I^{mw}(t)$, $\boldsymbol{\mu}^{mw}(t)$ that minimize the right hand side of (16)), then we have $C = 0$ and $\epsilon_V = \epsilon_U = \epsilon_Z = \epsilon_Q = \delta = 0$. In this case, we can also choose $\theta = 0$ so that performance is within $O(1/V)$ of the optimal value f^* . This special case is similar to the stochastic network optimization result of [1], with the exception that [1] assumes the convex cost function $f(\mathbf{x})$ is non-decreasing in each entry of \mathbf{x} (using auxiliary variables with “one-sided” virtual queues that are always non-negative), whereas here we treat a possibly non-monotonic cost function via (possibly negative) virtual queues $Z_m(t)$.

Proof: (Theorem 1) See Appendix A. □

The following related theorem uses a variable $V(t)$ parameter and allows for the uncertainty to tend to zero while achieving the exact penalty f_{θ}^* . Its proof follows as a simple consequence of the proof of Theorem 1.

⁹The expression (28) holds for all t (without the lim sup) in the special case when $\Theta(0) = \mathbf{0}$ and $f(\mathbf{x})$ is linear so that $f(\mathbf{x}) = l(\mathbf{x})$. The rate at which the limit converges in the general (non-linear) case is proportional to the rate at which the time average expectations of $\gamma_m(t)$ converge to the time average expectations of $x_m(t)$ for each $m \in \tilde{\mathcal{M}}$, which is roughly the average of $|Z_m(t)|/t$. This is highlighted in the proof of the theorem, see inequality (38).

Theorem 2: (Variable $V(t)$ parameter) Suppose Assumptions 1 and 2 hold. Let β_1 and β_2 be values such that $0 < \beta_1 < \beta_2 < 1$. Assume that after some finite time t_0 , we use a $V(t)$ parameter that increases with time, so that $V(t) = (t - t_0 + 1)^{\beta_2} V_0$ for all $t \geq t_0$ and for some constant $V_0 > 0$. Assume the queue states at time t_0 are arbitrary but finite, and assume we make control decisions $k(t)$, $I(t)$, $\gamma(t)$ such that the following holds for all $t \geq t_0$ (which is a modification of Assumption 3):

$$\begin{aligned} \mathbb{E}\{RHS(t, \Theta(t), k(t), I(t), \gamma(t))\} &\leq \mathbb{E}\{RHS(t, \Theta(t), k^{mw}(t), I^{mw}(t), \gamma^{mw}(t))\} \\ &+ C(t) + V(t)\epsilon_V(t) \\ &+ \sum_{n=1}^N \mathbb{E}\{U_n(t)\} \epsilon_U(t) + \sum_{m \in \mathcal{M}} \mathbb{E}\{|Z_m(t)|\} \epsilon_Z(t) + \sum_{l=1}^L \mathbb{E}\{Q_l(t)\} \epsilon_Q(t) \end{aligned}$$

where $C(t)$, $\epsilon_V(t)$, $\epsilon_U(t)$, $\epsilon_Z(t)$, $\epsilon_Q(t)$ are deterministic functions of time such that:

$$\lim_{t \rightarrow \infty} \epsilon_x(t) = 0$$

where $x \in \{V, U, Z, Q\}$, and where:

$$C(t) \leq O((t - t_0 + 1)^{\beta_1}) \text{ for } t \geq t_0$$

Then the time average constraints (9)-(10) hold, and all queues $Q_l(t)$ are *mean rate stable*, in the sense that:

$$\lim_{t \rightarrow \infty} \frac{\mathbb{E}\{Q_l(t)\}}{t} = 0 \text{ for all } l \in \{1, \dots, L\}$$

Further, the time average cost converges to the optimal value f_θ^* :

$$\lim_{t \rightarrow \infty} f(\bar{\mathbf{x}}(t)) = f_\theta^*$$

Proof: See Appendix B. □

This method of using an increasing $V(t)$ parameter can be viewed as a stochastic analogue of classic diminishing step-size methods for static optimization problems [16]. We note that $C(t)$ is assumed to increase at a rate slower than that of $V(t)$, while the $\epsilon_x(t)$ functions can converge to zero with any rate. Note that mean rate stability is a weak form of stability, and does *not* imply that average queue sizes and delays are finite. In fact, typically average congestion and delay are *necessarily* infinite when exact cost optimization is achieved [17] [18].

Remark 1: The results of Theorems 1 and 2 can be generalized to allow the $h_n(\mathbf{x})$ functions to be convex (possibly non-linear) by using one auxiliary variable $\gamma_m(t)$ for each penalty $x_m(t)$, in which case the constraints (10) can be enforced by modifying the virtual queues $U_n(t)$ in (12) to $\hat{U}_n(t)$ with dynamics:

$$\hat{U}_n(t+1) = \max[\hat{U}_n(t) + h_n(\gamma(t)) - b_n, 0]$$

This has the disadvantage of creating more virtual queues (one for each penalty $m \in \mathcal{M}$ rather than one for each penalty $m \in \tilde{\mathcal{M}}$), but has the advantage of allowing for non-linear $h_n(\mathbf{x})$ functions. It has the additional advantage of removing the uncertain $\mathbf{x}(t)$ penalties from the drift terms corresponding to the queues $\hat{U}_n(t)$. This ensures $\epsilon_U = 0$ whenever the auxiliary variables are chosen according to the max-weight rule $\gamma^{mw}(t)$ (which, due to separability, does not require knowledge of the $F_k(\omega)$ distributions). Similarly, one can also use auxiliary variables

in the cost function $f(\gamma(t))$ (as a proxy for the $f(x(t))$ values), so that $\epsilon_V = 0$. With these modifications, all uncertainty is isolated to ϵ_Z and ϵ_Q .

Remark 2: Theorems 1 and 2 can be used for any form of approximate scheduling, including cases when the optimal $I(t)$ decision involves a complex combinatorial choice that can only be approximated (or when the optimization for the auxiliary variable $\gamma(t)$ is approximate). This is related to similar approximate scheduling results developed for systems without stage-1 decisions in [1] [7] [19] [20]. However, our main interest is when the approximation is due to the uncertainty in the probability distributions $F_k(\omega)$, and max-weight learning algorithms for this context are developed in the next section.

III. ESTIMATING THE MAX-WEIGHT FUNCTIONAL

Theorem 1 suggests that our control policy should make decisions for $k(t)$, $I(t)$, $\gamma(t)$ every slot in an effort to minimize the right hand side of (16). The optimal auxiliary variable decisions $\gamma^{mw}(t)$ for this goal have already been established and are given by the solution of (23)-(24). Note that these decisions do not require knowledge of the $F_k(\omega)$ distribution. Likewise, the optimal $I^{mw}(t)$ decision does not require knowledge of the $F_k(\omega)$ distribution. Specifically, given a collection of observed queue backlogs $\Theta(t)$ and an observed outcome $\omega(t)$ (which is the result of the stage-1 decision $k(t)$ that is chosen), $I^{mw}(t)$ is defined as the optimal solution to the following (breaking ties arbitrarily):

$$\begin{aligned} \text{Minimize:} \quad & VI(\hat{x}(k(t), \omega(t), I(t))) + \sum_{n=1}^N U_n(t) h_n(\hat{x}(k(t), \omega(t), I(t))) + \\ & \sum_{m \in \tilde{\mathcal{M}}} Z_m(t) \hat{x}_m(k(t), \omega(t), I(t)) - \sum_{l=1}^L Q_l(t) [\hat{\mu}_l(k(t), \omega(t), I(t)) - \hat{a}_l(k(t), \omega(t), I(t))] \quad (29) \\ \text{Subject to:} \quad & I(t) \in \mathcal{I} \end{aligned}$$

The complexity of making these $I^{mw}(t)$ decisions depends on the physical structure of the network. The decisions are often trivial when the set \mathcal{I} contains only a finite (and small) number of control options (such as when the decisions are to remain idle or serve a single queue), in which case the function (29) is simply compared on each of the different choices in \mathcal{I} . For multi-hop networks with combinatorial resource allocation constraints, the choice of $I^{mw}(t)$ might be difficult, although constant-factor approximations are often possible (see [1] [7] [19] [20]).

The optimal $k^{mw}(t)$ decisions can be defined in terms of the $I^{mw}(t)$ decisions as follows: On each slot t , $k^{mw}(t)$ is chosen as k , according an independent type- k exploration event, with probability θ/K . If no exploration event occurs on slot t (which happens with probability $1 - \theta$), the queue backlogs $\Theta(t)$ are observed and $k^{mw}(t)$ is chosen as the value $k \in \{1, \dots, K\}$ with the lowest value of $e_k(t)$ (breaking ties arbitrarily), where $e_k(t)$ is defined:

$$e_k(t) \triangleq \mathbb{E} \left\{ \min_{I \in \mathcal{I}} [Y_k(I, \omega(t), \Theta(t))] \mid k(t) = k, \Theta(t) \right\} \quad (30)$$

where $\omega(t)$ is the random outcome that results from the stage-1 choice $k(t) = k$, and the function $Y_k(I, \omega, \Theta)$ is

defined for a particular stage-2 decision I , outcome ω , and queue state $\Theta = [\mathbf{Q}; \mathbf{U}; \mathbf{Z}]$, as follows:

$$\begin{aligned}
Y_k(I, \omega, \Theta) &\triangleq Vl(\hat{x}(k, \omega, I)) + \sum_{n=1}^N U_n h_n(\hat{x}(k, \omega, I)) \\
&+ \sum_{m \in \tilde{\mathcal{M}}} Z_m \hat{x}_m(k, \omega, I) \\
&- \sum_{l=1}^L Q_l [\hat{\mu}_l(k, \omega, I) - \hat{a}_l(k, \omega, I)]
\end{aligned} \tag{31}$$

Thus, $e_k(t)$ is the expected value of the expression (29) over the distribution $F_k(\omega)$ for the $\omega(t)$ random variable that arises from choosing $k(t) = k$, assuming that the optimal stage-2 decision $I^{mw}(t)$ is then made. However, computation of the exact $e_k(t)$ values would typically require full knowledge of the probability distributions $F_k(\omega)$ (and the computation may be difficult even if these distributions are fully known). Rather than using the exact conditional expectations, we consider two forms of estimates.

A. Estimating the $e_k(t)$ value — Approach 1

Define an integer W that represents a *moving average window size*. For each stage-1 option $k \in \{1, \dots, K\}$ and each time t , define $\omega_1^{(k)}(t), \dots, \omega_W^{(k)}(t)$ as the actual $\omega(\tau)$ outcomes observed over the last W type- k exploration events that took place before time t . Define the estimate $\hat{e}_k(t)$ as follows:

$$\hat{e}_k(t) \triangleq \frac{1}{W} \sum_{w=1}^W \min_{I \in \mathcal{I}} \left[Y_k(I, \omega_w^{(k)}(t), \Theta(t)) \right]$$

In the case when there have not yet been W previous type- k exploration events by time t , the estimate $\hat{e}_k(t)$ is taken with respect to the (fewer than W) events, and is set to zero if no such events have occurred. The estimates $\hat{e}_k(t)$ can be viewed as empirical averages of the function (31), using the current queue backlogs $\Theta(t) = [\mathbf{Q}(t); \mathbf{U}(t); \mathbf{Z}(t)]$ but using the outcomes $\omega_w^{(k)}(t)$ observed on previous type- k exploration events and the corresponding optimal stage-2 decisions.

Note that one might define $\hat{e}_k(t)$ according to an average over the past W slots on which stage-1 decision k has been made, rather than over the past W type- k exploration events. The reason we have used exploration events is to overcome the subtle “inspection paradox” issues involved in sampling the previous $\omega(\tau)$ outcomes. Indeed, even though $\omega(\tau)$ is generated in an i.i.d. way every slot in which $k(\tau) = k$ is chosen, the distribution of the last-seen outcome ω that corresponds to a particular decision k may be *skewed* in favor of creating larger penalties. This is because our algorithm may choose to avoid decision k for a longer period of time if this last outcome was non-favorable. Sampling at random type- k exploration events ensures that our samples indeed form an i.i.d. sequence. An additional difficulty remains: Even though these samples $\{\omega_w^{(k)}(t)\}$ form an i.i.d. sequence, they are *not* independent of the queue values $\Theta(t)$, as these prior outcomes have influenced the current queue states. We overcome this difficulty in Section III-D via a delayed-queue analysis.

This form of estimation does not require knowledge of the $F_k(\omega)$ distributions. However, evaluation of $\hat{e}_k(t)$ requires W computations of the type (29) on each slot t , according to the value of each particular $\omega_w^{(k)}(t)$ vector.

This can be difficult in the case when W is large, and hence the next subsection describes a second estimation approach that uses only one such computation per slot.

B. Estimating the $e_k(t)$ value — Approach 2

Again let W be an integer moving average window size. For each stage-1 decision $k \in \{1, \dots, K\}$, define $\omega_1^{(k)}(t), \dots, \omega_W^{(k)}(t)$ the same as in Approach 1. Further define $\Theta_1^{(k)}(t), \dots, \Theta_W^{(k)}(t)$ as the corresponding *queue backlogs* at the latest W type- k exploration events before time t . Define an estimate $\tilde{e}_k(t)$ as follows:

$$\tilde{e}_k(t) \triangleq \frac{1}{W} \sum_{w=1}^W \min_{I \in \mathcal{I}} [Y_k(I, \omega_w^{(k)}(t), \Theta_w^{(k)}(t))]$$

The $\tilde{e}_k(t)$ estimate is adjusted appropriately if fewer than W type- k exploration events have occurred (being set to zero initially). This approach is different from Approach 1 in that the current queue backlogs are not used. Hence, this is simply an empirical average over the past W samples of the actual cost achieved in the $I^{mw}(\tau)$ computation (29) at those particular sample times τ . Because $I^{mw}(\tau)$ (and its corresponding cost) was already computed on slot τ in order to make the stage-2 control decision, we can simply reuse the same value, without requiring any additional computation of problems of type (29).

C. The Max-Weight Learning Algorithm

Let θ be a given exploration probability (so that $0 \leq \theta < 1$ and exploration events of type K occur with probability θ/K). Let $\sigma > 0$ be a given parameter, and let $V(t)$ be a given (non-negative) control function of slot t (possibly a constant function). Let $\hat{W}(t)$ be a (possibly constant) function such that $\hat{W}(t) \geq 1$ for all t , and define $W_0 \triangleq \hat{W}(0)$. Define the actual window size used at slot t (for either Approach 1 or Approach 2) as follows:

$$W(t) \triangleq \min[\hat{W}(t), W_{rand}(t)]$$

where $W_{rand}(t)$ is the minimum number exploration events that have occurred for any type (minimized over the types $k \in \{1, \dots, K\}$), including the W_0 events that take place at initialization as described below. Thus, there are always at least $W(t)$ type- k exploration events by time t . The *Max-Weight Learning Algorithm* is as follows.

- (Initialization) For a given integer $W_0 > 0$, let $\Theta(-KW_0) = \mathbf{0}$, and run the system over slots $t = \{-W_0K, -W_0K + 1, \dots, -1\}$, choosing each stage-1 decision option $k \in \{1, \dots, K\}$ in a fixed round-robin order (and choosing $I^{mw}(t)$ according to (29) and $\gamma^{mw}(t)$ according to (23)-(24)). This ensures that we have W_0 independent samples by time 0, and creates a possibly non-zero initial queue state $\Theta(0)$. Next perform the following sequence of actions for each slot $t \geq 0$.
- (Stage-1 Decisions) Independently with probability θ , decide to have an exploration event. If there is an exploration event, choose $k(t)$ uniformly over all options $\{1, \dots, K\}$. If there is no exploration event, then under Approach 1 we observe current queue backlogs $\Theta(t)$ and compute $\hat{e}_k(t)$ for each $k \in \{1, \dots, K\}$ (using window size $W(t)$). We then choose $k(t)$ as the index $k \in \{1, \dots, K\}$ that minimizes $\hat{e}_k(t)$ (breaking

ties arbitrarily). Under Approach 2, if there is no exploration event we choose $k(t)$ to minimize $\tilde{e}_k(t)$ (using window size $W(t)$).

- (Stage-2 Decisions) Observe the queue backlogs $\Theta(t)$ and the outcome $\omega(t)$ that resulted from the stage-1 decision. Then choose $I^{mw}(t) \in \mathcal{I}$ according to (29). Choose auxiliary variables $\gamma^{mw}(t)$ according to (23)-(24).
- (Past Value Storage) For Approach 1, store the resulting $\omega(t)$ vector in memory as appropriate. For Approach 2, store the resulting cost from (29) in memory as appropriate.
- (Queue Updates) Update virtual queues $U_n(t)$ according to (12) and $Z_m(t)$ according to (13). Also allow the actual system queues $Q_l(t)$ to proceed according to (3).

Remark 3: For some systems, we may not require an exploration event for each of the K stage-1 decision options. For example, in an L -queue downlink where the decisions are to either measure all channels, blindly transmit over one of the L channels, or remain idle (as in [5]), there are $K = L + 2$ stage-1 options. However, the “idle” choice does not require any exploration events, as it clearly incurs a cost of 0. Further, the information gained by randomly choosing to blindly transmit over a given channel can also be gained by measuring all channels, as the outcome of the channel measurement can be used to determine if a blind transmission would have been successful. It is therefore more efficient to modify the algorithm by considering only *one type of exploration event*: the one that randomly chooses to measure all channels. Similarly, in DIVBAR-like situations where the K decisions involve sending a packet of one of the various commodities (as in [7]), the success/failure event observed after sending any particular packet does not depend on the packet commodity and hence can be used to update the max-weight estimates for each commodity.

D. Analysis of the Max-Weight Learning Algorithm

For brevity, we analyze only Approach 2.¹⁰ Let $k^{mw}(t)$ denote the (ideal) max-weight stage-1 decision on slot t , and let $\tilde{k}(t)$ denote the Approach 2 decision. Recall that Approach 2 also uses the (ideal) $I^{mw}(t)$ and $\gamma^{mw}(t)$ decisions. Our goal is to compute parameters C , ϵ_V , ϵ_U , ϵ_Z , ϵ_Q for (25) that can be plugged into Theorem 1.

Theorem 3: (Performance Under Approach 2 — Fixed Window) Suppose the Max-Weight Learning Algorithm with Approach 2 is implemented using an exploration probability $\theta > 0$. Suppose we use a fixed integer window size $W = W_0 > 0$ (so that $W(t) = W$ for all t , and our initialization takes W samples from each exploration type before time 0). Suppose that $V(t)$ is held constant, so that $V(t) = V$ for some $V > 0$. Then condition (25) of Assumption 3 holds with:

$$C = \frac{cWK^2(1+\theta)}{\theta}, \quad \epsilon_V = \epsilon_U = \epsilon_Z = \epsilon_Q = \frac{Ky_{diff}^{max}}{2\sqrt{W}}$$

where c and y_{diff}^{max} are constants that are independent of queue backlog and of V , W , θ (and depend on the maximum and minimum penalties and maximum queue changes that can occur on one slot).

¹⁰Bounds on the performance of Approach 1 can be obtained similarly. In practice, Approach 1 would typically have superior performance because it uses current queue backlogs.

Proof: See Appendix C. □

It follows that if the fixed window size W is chosen to be suitably large, then the $\epsilon_V, \epsilon_U, \epsilon_Z, \epsilon_Q$ constants will be small enough to satisfy the conditions $\epsilon_U < \epsilon_{max}, \epsilon_Z < \sigma, \epsilon_Q < \epsilon_{max}$ required for Theorem 1, and hence the result of Theorem 1 holds for this max-weight learning algorithm.

Theorem 4: (Performance Under Approach 2 — Variable $W(t)$ and $V(t)$) Suppose that we use the Max-Weight Learning algorithm (with Approach 2) using an exploration probability $\theta > 0$ and a variable $V(t)$ and $W(t)$ with initialization parameter $W_0 = 1$, and with:

$$V(t) = (t + 1)^{\beta_2} V_0, \quad W(t) = \min[(t + 1)^{\beta_1}, W_{rand}(t)]$$

where β_1 and β_2 are constants such that $0 < \beta_1 < \beta_2 < 1$, V_0 is a positive constant, and where we recall that $W_{rand}(t)$ is the minimum number exploration events of type k that have occurred, minimized over all $k \in \{1, \dots, K\}$. Then the time average constraints (9)-(10) hold, all queues $Q_l(t)$ are *mean rate stable*, and the time average cost converges to the optimal value f_θ^* :

$$\lim_{t \rightarrow \infty} f(\bar{\mathbf{x}}(t)) = f_\theta^*$$

Proof: The proof combines results from the proofs of Theorems 3 and 2, and is given in Appendix E. □

IV. CONCLUSION

This work extends the important max-weight framework for stochastic network optimization to a context with 2-stage decisions and unknown distributions that govern the stochastics at the first stage. This is useful in a variety of contexts, including transmission scheduling in wireless networks in unknown environments and with unknown channels. The learning algorithms developed here are based on estimates of expected max-weight functionals, and are much more efficient than algorithms that would attempt to learn the complete probability distributions associated with the system. Our analysis provides explicit bounds on the deviation from optimality in terms of the sample size W and the control parameter V . The W and V parameters also affect an explicit tradeoff in average congestion and delay. A modified algorithm with time-varying $W(t)$ and $V(t)$ parameters was shown to converge to exact optimal performance while keeping all queues mean-rate stable, at the cost of incurring a possibly infinite average congestion and delay.

APPENDIX A — PROOF OF THEOREM 1

Proof: (Theorem 1 — The Queue Stability Inequality (27)) Writing the drift inequality (16) using the $RHS(\cdot)$ function yields:

$$\mathbb{E} \left\{ Vl(\mathbf{x}(t)) + V\tilde{f}(\gamma(t)) \mid \Theta(t) \right\} + \Delta(\Theta(t)) \leq \mathbb{E} \{ RHS(t, \Theta(t), k(t), I(t), \gamma(t)) \mid \Theta(t) \}$$

Taking expectations of both sides with respect to the queue state distribution for $\Theta(t)$ and using the law of iterated expectations yields:

$$\begin{aligned}
V\mathbb{E}\left\{l(\mathbf{x}(t)) + \tilde{f}(\gamma(t))\right\} + \mathbb{E}\{L(\Theta(t+1))\} - \mathbb{E}\{L(\Theta(t))\} &\leq \mathbb{E}\{RHS(t, \Theta(t), k(t), I(t), \gamma(t))\} \\
&\leq \mathbb{E}\{RHS(t, \Theta(t), k^{mw}(t), I^{mw}(t), \gamma^{mw}(t))\} \\
&\quad + C + V\epsilon_V + \epsilon_Z \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{|Z_m(t)|\} \\
&\quad + \epsilon_U \sum_{n=1}^N \mathbb{E}\{U_n(t)\} + \epsilon_Q \sum_{l=1}^L \mathbb{E}\{Q_l(t)\} \quad (32) \\
&\leq \mathbb{E}\{RHS(t, \Theta(t), k'(t), I'(t), \gamma'(t))\} \\
&\quad + C + V\epsilon_V + \epsilon_Z \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{|Z_m(t)|\} \\
&\quad + \epsilon_U \sum_{n=1}^N \mathbb{E}\{U_n(t)\} + \epsilon_Q \sum_{l=1}^L \mathbb{E}\{Q_l(t)\} \quad (33)
\end{aligned}$$

where (32) holds by Assumption 3, and (33) holds because the max-weight policy minimizes the expectation of $RHS(\cdot)$ over all alternative decisions for slot t . The decisions $k'(t)$, $I'(t)$, $\gamma'(t)$ can be chosen as any feasible control decisions for slot t (where a feasible control decision for $k'(t)$ must also respect the random exploration events of probability θ). Suppose that $k'(t)$ and $I'(t)$ are the decisions given in Assumption 2, so that properties (20) and (21) hold. Choose auxiliary decision variables $\gamma'(t) = (\gamma'_m(t))_{m \in \tilde{\mathcal{M}}}$ as follows:

$$\gamma'_m(t) = \begin{cases} \mathbb{E}\{x'_m(t)\} + \sigma & \text{if } Z_m(t) \geq 0 \\ \mathbb{E}\{x'_m(t)\} - \sigma & \text{if } Z_m(t) < 0 \end{cases} \quad (34)$$

Note that these $\gamma'_m(t)$ decisions satisfy the required constraints (7). That is because for each $m \in \tilde{\mathcal{M}}$ we have $x_m^{min} \leq \mathbb{E}\{x'_m(t)\} \leq x_m^{max}$ and therefore:

$$x_m^{min} - \sigma \leq \mathbb{E}\{x'_m(t)\} - \sigma \leq \mathbb{E}\{x'_m(t)\} + \sigma \leq x_m^{max} + \sigma$$

Using these $\gamma'_m(t)$ decisions and the definition of $RHS(\cdot)$ in the inequality (33) yields:¹¹

$$\begin{aligned}
V\mathbb{E}\left\{l(\mathbf{x}(t)) + \tilde{f}(\gamma(t))\right\} + \mathbb{E}\{L(\Theta(t+1))\} - \mathbb{E}\{L(\Theta(t))\} &\leq B + C + V\epsilon_V \\
&+ \mathbb{E}\left\{Vl(\mathbf{x}'(t)) + V\tilde{f}(\gamma'(t))\right\} \\
- \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{Z_m(t)[\mathbb{E}\{x'_m(t)\} - x'_m(t)]\} & \\
- \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{|Z_m(t)|[\sigma - \epsilon_Z]\} & \\
- \sum_{n=1}^N \mathbb{E}\{U_n(t)[b_n - h_n(\mathbf{x}'(t)) - \epsilon_U]\} & \\
- \sum_{l=1}^L \mathbb{E}\{Q_l(t)[\mu'_l(t) - A'_l(t) - \epsilon_Q]\} & \tag{35}
\end{aligned}$$

Note that because the policies $k'(t)$ and $I'(t)$ are stationary, randomized, and independent of the queue backlog vector $\Theta(t)$, and because the functions $h_n(\mathbf{x})$ are linear or affine, we have:

$$\begin{aligned}
\mathbb{E}\{U_n(t)h_n(\mathbf{x}'(t))\} &= \mathbb{E}\{U_n(t)\}h_n(\mathbb{E}\{\mathbf{x}'(t)\}) \\
\mathbb{E}\{Z_m(t)x'_m(t)\} &= \mathbb{E}\{Z_m(t)\}\mathbb{E}\{x'_m(t)\} \\
\mathbb{E}\{Q_l(t)[\mu'_l(t) - A'_l(t)]\} &= \mathbb{E}\{Q_l(t)\}\mathbb{E}\{\mu'_l(t) - A'_l(t)\}
\end{aligned}$$

Using these identities together with properties (20)-(21) directly in the right hand side of (35) and rearranging terms yields:

$$\begin{aligned}
\mathbb{E}\{L(\Theta(t+1))\} - \mathbb{E}\{L(\Theta(t))\} &\leq B + C + V[l_{diff} + \tilde{f}_{diff} + \epsilon_V] \\
&- (\epsilon_{max} - \epsilon_U) \sum_{n=1}^N \mathbb{E}\{U_n(t)\} \\
&- (\sigma - \epsilon_Z) \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{|Z_m(t)|\} \\
&- (\epsilon_{max} - \epsilon_Q) \sum_{l=1}^L \mathbb{E}\{Q_l(t)\} \tag{36}
\end{aligned}$$

where we have used the following fact:

$$\mathbb{E}\{l(\mathbf{x}'(t)) - l(\mathbf{x}(t))\} \leq l_{diff}, \quad \mathbb{E}\left\{\tilde{f}(\gamma') - \tilde{f}(\gamma(t))\right\} \leq \tilde{f}_{diff}$$

The inequality (36) holds for all slots $t \in \{0, 1, 2, \dots\}$. Summing the telescoping series over $\tau \in \{0, 1, \dots, t-1\}$ (as in [1]) and dividing by t yields:

$$\begin{aligned}
\frac{\mathbb{E}\{L(\Theta(t))\} - \mathbb{E}\{L(\Theta(0))\}}{t} &\leq B + C + V[l_{diff} + \tilde{f}_{diff} + \epsilon_V] \\
- \frac{1}{t} \sum_{\tau=0}^{t-1} \left[(\epsilon_{max} - \epsilon_U) \sum_{n=1}^N \mathbb{E}\{U_n(\tau)\} + (\sigma - \epsilon_Z) \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{|Z_m(\tau)|\} + (\epsilon_{max} - \epsilon_Q) \sum_{l=1}^L \mathbb{E}\{Q_l(\tau)\} \right] &
\end{aligned}$$

¹¹Recall that $RHS(\cdot)$ is defined as the right hand side of (16).

Using non-negativity of the Lyapunov function $L(\cdot)$ in the above inequality proves (27). Taking the lim sup of (27) as $t \rightarrow \infty$ proves that the queues $Q_l(t)$, $Z_m(t)$, $U_n(t)$ are strongly stable (for all $l \in \{1, \dots, L\}$, $m \in \tilde{\mathcal{M}}$, $n \in \{1, \dots, N\}$). Hence (by Lemma 1), the inequality constraints (9)-(11) are satisfied. \square

Proof: (Theorem 1 — The Utility Inequality (28)) Recall that the inequality (33) holds for any alternative set of feasible control decisions $k''(t)$, $I''(t)$, $\gamma''(t)$. Re-writing (33) using this notation and using the definition of $RHS(\cdot)$ yields:

$$\begin{aligned} V\mathbb{E}\{l(\mathbf{x}(t)) + \tilde{f}(\boldsymbol{\gamma}(t))\} + \mathbb{E}\{L(\boldsymbol{\Theta}(t+1))\} - \mathbb{E}\{L(\boldsymbol{\Theta}(t))\} &\leq B + C + V\epsilon_V + \epsilon_Z \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{|Z_m(t)|\} \\ &+ \mathbb{E}\{Vl(\mathbf{x}''(t)) + V\tilde{f}(\boldsymbol{\gamma}''(t))\} \\ &- \sum_{n=1}^N \mathbb{E}\{U_n(t)(b_n - h_n(\mathbf{x}''(t)) - \epsilon_U)\} \\ &- \sum_{m \in \tilde{\mathcal{M}}} \mathbb{E}\{Z_m(t)(\gamma''_m(t) - x''_m(t))\} \\ &- \sum_{l=1}^L \mathbb{E}\{Q_l(t)(\mu''_l(t) - A''_l(t) - \epsilon_Q)\} \end{aligned}$$

Let α be a probability (to be chosen later), and define joint control actions $(k''(t); I''(t); \gamma''(t))$ as follows:

$$(k''(t); I''(t); \gamma''(t)) = \begin{cases} (k'(t); I'(t); \gamma'(t)) & \text{with prob. } \alpha \\ (k^*(t); I^*(t); \gamma^*) & \text{with prob. } 1 - \alpha \end{cases}$$

where $k'(t)$, $I'(t)$ are as defined in Assumption 2 (and satisfy (20)-(21)), variables $\gamma'_m(t)$ are as defined in (34), and $I^*(t)$, $k^*(t)$, γ^* are as defined in Assumption 1 (and satisfy properties (17)-(19)). Note that the $k''(t)$ decision defined here still has random exploration events with probability θ , as both $k'(t)$ and $k^*(t)$ have such events. Also note that $\gamma''_m(t)$ satisfies (7) because both $\gamma'_m(t)$ and γ^*_m satisfy (7). Further, we have:

$$\begin{aligned} \mathbb{E}\{\mathbf{x}''(t)\} &= \alpha \mathbb{E}\{\mathbf{x}'(t)\} + (1 - \alpha) \mathbb{E}\{\mathbf{x}^*(t)\} \\ \mathbb{E}\{\boldsymbol{\gamma}''(t)\} &= \alpha \mathbb{E}\{\boldsymbol{\gamma}'(t)\} + (1 - \alpha) \boldsymbol{\gamma}^* \end{aligned}$$

It follows from properties (20)-(21) and (17)-(19) (together with linearity of $l(\mathbf{x})$ and $h_n(\mathbf{x})$ and the fact that the randomized $k''(t)$ and $I''(t)$ choices are independent of queue backlog) that:

$$\begin{aligned} V\mathbb{E}\left\{l(\mathbf{x}(t)) + \tilde{f}(\gamma(t))\right\} + \mathbb{E}\{L(\Theta(t+1))\} - \mathbb{E}\{L(\Theta(t))\} &\leq B + C + V\epsilon_V \\ &+ (1-\alpha)Vl(\mathbb{E}\{\mathbf{x}^*(t)\}) + (1-\alpha)V\tilde{f}(\gamma^*) \\ &+ \alpha Vl(\mathbb{E}\{\mathbf{x}'(t)\}) + \alpha V\mathbb{E}\left\{\tilde{f}(\gamma'(t))\right\} \\ &- \sum_{n=1}^N \mathbb{E}\{U_n(t)\}(\alpha\epsilon_{max} - \epsilon_U) \\ &- \sum_{m \in \mathcal{M}} \mathbb{E}\{|Z_m(t)|\}(\alpha\sigma - \epsilon_Z) \\ &- \sum_{l=1}^L \mathbb{E}\{Q_l(t)\}(\alpha\epsilon_{max} - \epsilon_Q) \end{aligned}$$

Now choose α as follows:

$$\alpha = \max\left[\frac{\epsilon_U}{\epsilon_{max}}, \frac{\epsilon_Z}{\sigma}, \frac{\epsilon_Q}{\epsilon_{max}}\right]$$

This is a valid probability because we have assumed that $\epsilon_U < \epsilon_{max}$, $\epsilon_Z < \sigma$, $\epsilon_Q < \epsilon_{max}$. The above inequality reduces to:

$$\begin{aligned} V\mathbb{E}\left\{l(\mathbf{x}(t)) + \tilde{f}(\gamma(t))\right\} + \mathbb{E}\{L(\Theta(t+1))\} - \mathbb{E}\{L(\Theta(t))\} &\leq \\ B + C + V\epsilon_V + Vf_\theta^* + \alpha V(l_{diff} + \tilde{f}_{diff}) &\quad (37) \end{aligned}$$

The above inequality holds for all t . Taking a telescoping series over $\tau \in \{0, 1, \dots, t-1\}$ yields:

$$\frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E}\left\{l(\mathbf{x}(\tau)) + \tilde{f}(\gamma(\tau))\right\} + \frac{\mathbb{E}\{L(\Theta(t))\} - \mathbb{E}\{L(\Theta(0))\}}{Vt} \leq f_\theta^* + \epsilon_V + \alpha(l_{diff} + \tilde{f}_{diff}) + \frac{B+C}{V}$$

Therefore, using $\delta \triangleq \alpha(l_{diff} + \tilde{f}_{diff})$, non-negativity of $L(\cdot)$, and Jensen's inequality with convexity of $l(\mathbf{x})$ and $\tilde{f}(\mathbf{x})$, we have:

$$l(\bar{\mathbf{x}}(t)) + \tilde{f}(\bar{\gamma}(t)) \leq f_\theta^* + \epsilon_V + \delta + \frac{B+C}{V} + \frac{\mathbb{E}\{L(\Theta(0))\}}{Vt}$$

However, we have:

$$\tilde{f}(\bar{\gamma}(t)) \geq \tilde{f}(\bar{\mathbf{x}}(t)) - \tilde{M}\nu\|\bar{\mathbf{x}}(t) - \bar{\gamma}(t)\|$$

where ν is the magnitude of the largest left or right partial derivative of the $\tilde{f}(\cdot)$ function and \tilde{M} is the cardinality of $\tilde{\mathcal{M}}$.¹² Combining the above two inequalities and using the fact that $f(\mathbf{x}) \triangleq l(\mathbf{x}) + \tilde{f}(\mathbf{x})$ yields:

$$f(\bar{\mathbf{x}}(t)) - \tilde{M}\nu\|\bar{\mathbf{x}}(t) - \bar{\gamma}(t)\| \leq f_\theta^* + \epsilon_V + \delta + \frac{B+C}{V} + \frac{\mathbb{E}\{L(\Theta(0))\}}{Vt} \quad (38)$$

Because the equality constraints (10) hold, we have that $\|\bar{\mathbf{x}}(t) - \bar{\gamma}(t)\| \rightarrow 0$. Taking the lim sup of (38) as $t \rightarrow \infty$ thus yields (28), completing the proof. \square

¹²Left and right partial derivatives exist and are finite for any convex function that is defined over the full space \mathbb{R}^M .

Note that in the special case when there are no auxiliary variables (so that $f(\mathbf{x})$ is linear and $f(\mathbf{x}) = l(\mathbf{x})$), and when all queues are initially empty, the inequality (38) reduces to the following cost guarantee that holds for all time t :

$$f(\bar{\mathbf{x}}(t)) \leq f_{\theta}^* + \epsilon_V + \delta + (B + C)/V$$

APPENDIX B — PROOF OF THE VARIABLE $V(t)$ THEOREM (THEOREM 2)

Proof: (Mean Rate Stability of all Queues) Assume without loss of generality that $\epsilon_U(t) < \epsilon_{max}$, $\epsilon_Z(t) < \sigma$, $\epsilon_Q(t) < \epsilon_{max}$ for all $t \geq t_0$ (else, choose a time \tilde{t}_0 for which this holds). Then, on a single slot t , we can apply the result from the proof of Theorem 1 with $V \triangleq V(t)$ and $\epsilon_x \triangleq \epsilon_x(t)$ (for $x \in \{V, U, Z, Q\}$). Thus, for any time $t \geq t_0$ we have from (36):

$$\mathbb{E}\{L(\Theta(t+1))\} - \mathbb{E}\{L(\Theta(t))\} \leq B + C(t) + V(t)[l_{diff} + \tilde{f}_{diff} + \epsilon_V(t)]$$

where we have neglected the three non-positive terms on the right hand side of (36). Summing the above inequality over $\tau \in \{t_0, \dots, t-1\}$ yields:

$$\frac{\mathbb{E}\{L(\Theta(t))\} - \mathbb{E}\{L(\Theta(t_0))\}}{t - t_0} \leq B + \frac{O(t^{\beta_2+1})}{t - t_0}$$

where we have used the fact that $\sum_{\tau=t_0}^{t-1} C(\tau) \leq O(t^{\beta_1+1})$ and $\sum_{\tau=t_0}^{t-1} V(\tau) \leq O(t^{\beta_2+1})$. Because $L(\Theta(t))$ is a sum of squared queue lengths (for all queues), the above inequality implies that for any queue $Q_l(t)$:

$$\frac{\mathbb{E}\{Q_l(t)^2\}}{t - t_0} \leq B + \frac{O(t^{\beta_2+1})}{t - t_0} + \frac{\mathbb{E}\{L(\Theta(t_0))\}}{t - t_0}$$

Dividing the above inequality by $t - t_0$, taking square roots, and using the fact that $\mathbb{E}\{Q_l(t)^2\} \geq \mathbb{E}\{Q_l(t)\}^2$ yields:

$$\frac{\mathbb{E}\{Q_l(t)\}}{t - t_0} \leq \sqrt{\frac{B}{(t - t_0)} + \frac{O(t^{\beta_2+1})}{(t - t_0)^2} + \frac{\mathbb{E}\{L(\Theta(t_0))\}}{(t - t_0)^2}}$$

Because $\beta_2 + 1 < 2$, the right hand side above converges to 0 as $t \rightarrow \infty$. This holds for all queues $Q_l(t)$, and hence all these queues are *mean rate stable*. Similarly, it holds for all queues $Z_m(t)$ and $U_n(t)$ (for $m \in \tilde{\mathcal{M}}$ and $n \in \{1, \dots, N\}$), and so all these queues are mean rate stable. It follows by Lemma 1 that all inequality constraints (9)-(10) are satisfied. \square

Proof: (Cost Optimality) Again assume (without loss of generality) that $\epsilon_U(t) < \epsilon_{max}$, $\epsilon_Z(t) < \sigma$, $\epsilon_Q(t) < \epsilon_{max}$ for all $t \geq t_0$. We thus have from (37) that:

$$\mathbb{E}\{l(\mathbf{x}(t)) + \tilde{f}(\gamma(t))\} + \frac{\mathbb{E}\{L(\Theta(t+1)) - L(\Theta(t))\}}{V(t)} \leq \frac{B + C(t)}{V(t)} + \epsilon_V(t) + f_{\theta}^* + \alpha(t)(l_{diff} - \tilde{f}_{diff})$$

where $\alpha(t)$ is defined:

$$\alpha(t) \triangleq \max \left[\frac{\epsilon_U(t)}{\epsilon_{max}}, \frac{\epsilon_Z(t)}{\sigma}, \frac{\epsilon_Q(t)}{\epsilon_{max}} \right]$$

and satisfies $\alpha(t) \rightarrow 0$ as $t \rightarrow \infty$. The above holds for all $t \geq t_0$. Summing over $\tau \in \{t_0, \dots, t-1\}$ yields:

$$\begin{aligned} \sum_{\tau=t_0}^{t-1} \mathbb{E} \left\{ l(\mathbf{x}(\tau)) + \tilde{f}(\boldsymbol{\gamma}(\tau)) \right\} + \sum_{\tau=t_0+1}^{t-1} \mathbb{E} \{ L(\boldsymbol{\Theta}(\tau)) \} \left[\frac{1}{V(\tau-1)} - \frac{1}{V(\tau)} \right] + \frac{\mathbb{E} \{ L(\boldsymbol{\Theta}(t)) \}}{V(t-1)} - \frac{\mathbb{E} \{ L(\boldsymbol{\Theta}(t_0)) \}}{V(t_0)} \leq \\ (t-t_0)f_{\theta}^* + \sum_{\tau=t_0}^{t-1} \left[\frac{B+C(\tau)}{V(\tau)} + \epsilon_V(\tau) + \alpha(\tau)(l_{diff} - \tilde{f}_{diff}) \right] \end{aligned}$$

Using non-negativity of $L(\cdot)$ and the fact that $\frac{1}{V(\tau-1)} - \frac{1}{V(\tau)} \geq 0$ (because $V(\tau)$ is non-decreasing), and dividing by $(t-t_0)$ yields:

$$\frac{1}{t-t_0} \sum_{\tau=t_0}^{t-1} \mathbb{E} \left\{ l(\mathbf{x}(\tau)) + \tilde{f}(\boldsymbol{\gamma}(\tau)) \right\} - \frac{\mathbb{E} \{ L(\boldsymbol{\Theta}(t_0)) \}}{(t-t_0)V(t_0)} \leq f_{\theta}^* + \Psi(t) \quad (39)$$

where $\Psi(t)$ is defined:

$$\Psi(t) \triangleq \frac{1}{t-t_0} \sum_{\tau=t_0}^{t-1} \left[\frac{B+C(\tau)}{V(\tau)} + \epsilon_V(\tau) + \alpha(\tau)(l_{diff} - \tilde{f}_{diff}) \right]$$

Note that $C(\tau)/V(\tau) \rightarrow 0$ as $\tau \rightarrow \infty$, and hence $\Psi(t)$ is the time average of a function that converges to 0. We thus have $\Psi(t) \rightarrow 0$ as $t \rightarrow \infty$. By Jensen's inequality applied to the left hand side of (39) we have:

$$l(\bar{\mathbf{x}}(t)) + \tilde{f}(\bar{\boldsymbol{\gamma}}(t)) - \frac{\mathbb{E} \{ L(\boldsymbol{\Theta}(t_0)) \}}{(t-t_0)V(t_0)} \leq f_{\theta}^* + \Psi(t)$$

where $\bar{\mathbf{x}}(t)$ and $\bar{\boldsymbol{\gamma}}(t)$ are time average expectations over the interval $\tau \in \{t_0, \dots, t-1\}$. Because we already know $Z_m(t)$ is mean rate stable for all $m \in \tilde{\mathcal{M}}$, we have that $\|\bar{\boldsymbol{\gamma}}(t) - \bar{\boldsymbol{x}}(t)\| \rightarrow 0$ as $t \rightarrow \infty$ (by Lemma 1), and hence, as in the proof of Theorem 1 (using $f(\mathbf{x}) = l(\mathbf{x}) + \tilde{f}(\tilde{\mathbf{x}})$):

$$\limsup_{t \rightarrow \infty} f(\bar{\mathbf{x}}(t)) \leq f_{\theta}^*$$

Because f_{θ}^* is defined as the infimum cost subject to queue stability,¹³ it can be shown that the lim inf cannot be lower than f_{θ}^* , and so the limit of $f(\bar{\mathbf{x}}(t))$ exists and is equal to the lim sup, proving the result. \square

APPENDIX C — PROOF OF THEOREM 3

To prove Theorem 3, fix time t and define $\Omega(\boldsymbol{\Theta}(t))$ as follows:

$$\begin{aligned} \Omega(\boldsymbol{\Theta}(t)) \triangleq \mathbb{E} \left\{ RHS(t, \boldsymbol{\Theta}(t), \tilde{k}(t), I^{mw}(t), \boldsymbol{\gamma}^{mw}(t)) \mid \boldsymbol{\Theta}(t) \right\} \\ - \mathbb{E} \left\{ RHS(t, \boldsymbol{\Theta}(t), k^{mw}(t), I^{mw}(t), \boldsymbol{\gamma}^{mw}(t)) \mid \boldsymbol{\Theta}(t) \right\} \end{aligned}$$

Now note that because these right-hand sides differ only in terms comprising the $e_k(t)$ expression, we have:

$$\Omega(\boldsymbol{\Theta}(t)) = \mathbb{E} \left\{ e_{\tilde{k}(t)}(t) \mid \boldsymbol{\Theta}(t) \right\} - \min_{k \in \mathcal{K}} [e_k(t)]$$

where the expectation on the right hand side is over the random decision $\tilde{k}(t) = \arg \min_{k \in \mathcal{K}} [\tilde{e}_k(t)]$, which is based on the empirical average $\tilde{e}_k(t)$ formed by the past W random samples. It uses the fact that given a particular

¹³It can be shown that the infimum cost subject to *strong stability* is the same as the infimum cost subject to *mean rate stability*.

(possibly sub-optimal) decision $\tilde{k}(t) \in \mathcal{K}$, the resulting expected max-weight functional (using the exact but unknown distribution function $F_{\tilde{k}(t)}(\omega)$) is $e_{\tilde{k}(t)}$. Now for each $k \in \mathcal{K}$, define $\delta_k(t) \triangleq \tilde{e}_k(t) - e_k(t)$. We thus have:

$$\begin{aligned}
\mathbb{E} \left\{ e_{\tilde{k}(t)}(t) \mid \Theta(t) \right\} &= \mathbb{E} \left\{ \tilde{e}_{\tilde{k}(t)}(t) - \delta_{\tilde{k}(t)}(t) \mid \Theta(t) \right\} \\
&\leq \mathbb{E} \left\{ \tilde{e}_{\tilde{k}(t)}(t) \mid \Theta(t) \right\} + \mathbb{E} \left\{ \max_{k \in \mathcal{K}} [-\delta_k(t)] \mid \Theta(t) \right\} \\
&= \mathbb{E} \left\{ \min_{k \in \mathcal{K}} [\tilde{e}_k(t)] \mid \Theta(t) \right\} + \mathbb{E} \left\{ \max_{k \in \mathcal{K}} [-\delta_k(t)] \mid \Theta(t) \right\} \\
&= \mathbb{E} \left\{ \min_{k \in \mathcal{K}} [e_k(t) + \delta_k(t)] \mid \Theta(t) \right\} + \mathbb{E} \left\{ \max_{k \in \mathcal{K}} [-\delta_k(t)] \mid \Theta(t) \right\} \\
&\leq \min_{k \in \mathcal{K}} [e_k(t)] + \mathbb{E} \left\{ \max_{k \in \mathcal{K}} [\delta_k(t)] + \max_{k \in \mathcal{K}} [-\delta_k(t)] \mid \Theta(t) \right\} \\
&\leq \min_{k \in \mathcal{K}} [e_k(t)] + \sum_{k=1}^K \mathbb{E} \left\{ \max[\delta_k(t), 0] + \max[-\delta_k(t), 0] \mid \Theta(t) \right\} \\
&= \min_{k \in \mathcal{K}} [e_k(t)] + \sum_{k=1}^K \mathbb{E} \left\{ |\delta_k(t)| \mid \Theta(t) \right\}
\end{aligned}$$

It follows that:

$$\Omega(\Theta(t)) \leq \sum_{k=1}^K \mathbb{E} \left\{ |\delta_k(t)| \mid \Theta(t) \right\}$$

Therefore, by iterated expectations we have:

$$\mathbb{E} \left\{ \Omega(\Theta(t)) \right\} \leq \sum_{k=1}^K \mathbb{E} \left\{ |\delta_k(t)| \right\} \quad (40)$$

Note that $\mathbb{E} \left\{ \Omega(\Theta(t)) \right\}$ corresponds to the desired inequality (25), and hence it suffices to bound $\mathbb{E} \left\{ |\delta_k(t)| \right\}$. To this end, for each $k \in \{1, \dots, K\}$, define $T_k(t)$ as the number of timeslots that passed after the W th-latest type k exploration event. Thus, all samples $\omega_w^{(k)}(t)$ occur on type- k explorations events, and are on or after time $t - T_k(t)$. Define $\Theta_k(t) \triangleq \Theta(t - T_k(t))$. We have:

$$|\delta_k(t)| = |\tilde{e}_k(t) - e_k(t)| \leq |\tilde{e}_k(t) - \tilde{e}_k^{prev}(t)| + |\tilde{e}_k^{prev}(t) - e_k^{prev}(t)| + |e_k^{prev}(t) - e_k(t)| \quad (41)$$

where $\tilde{e}_k^{prev}(t)$ and $e_k^{prev}(t)$ are defined using queue lengths from the *previous* time $t - T_k(t)$ as follows:

$$\begin{aligned}
\tilde{e}_k^{prev}(t) &\triangleq \frac{1}{W} \sum_{w=1}^W \min_{I \in \mathcal{I}} [Y_k(I, \omega_w^{(k)}(t), \Theta(t - T_k(t)))] \\
e_k^{prev}(t) &\triangleq \mathbb{E} \left\{ \min_{I \in \mathcal{I}} [Y_k(I, \omega(t), \Theta(t - T_k(t)))] \mid \Theta(t - T_k(t)), k(t) = k \right\}
\end{aligned}$$

where the expectation in the definition of $e_k^{prev}(t)$ is with respect to the independent outcome $\omega(t)$ that has distribution $F_k(\omega)$. Comparing the definition of $e_k^{prev}(t)$ to the definition of $e_k(t)$ in (30), it is clear that they are different only in that they use different queue states (similarly, $\tilde{e}_k(t)$ and $\tilde{e}_k^{prev}(t)$ differ only in that they use different queue states). Because the maximum change in queue size on any single slot is bounded, we have the following lemma.

Lemma 3: For any $k \in \{1, \dots, K\}$, any time t , and regardless of queue backlog $\Theta(t)$, we have:

$$\mathbb{E} \{ |\tilde{e}_k(t) - \tilde{e}_k^{prev}(t)| + |e_k^{prev}(t) - e_k(t)| \} \leq d_1 \mathbb{E} \{ T_k(t) \} \quad (42)$$

where d_1 is a constant that is proportional to the maximum change in any queue over a single slot, and is independent of the current queue sizes and of W and K .

Proof: Define $I_w^{(k)}(t)$ as follows:

$$I_w^{(k)}(t) \triangleq \arg \min_{I \in \mathcal{I}} \left[Y_k(I, \omega_w^{(k)}, \Theta(t - T_k(t))) \right]$$

We have:

$$\begin{aligned} \tilde{e}_k(t) &= \frac{1}{W} \sum_{w=1}^W \min_{I \in \mathcal{I}} \left[Y_k(I, \omega_w^{(k)}(t), \Theta_w^{(k)}(t)) \right] \\ &\leq \frac{1}{W} \sum_{w=1}^W Y_k(I_w^{(k)}(t), \omega_w^{(k)}(t), \Theta_w^{(k)}(t)) \\ &\leq \frac{1}{W} \sum_{w=1}^W Y_k(I_w^{(k)}(t), \omega_w^{(k)}(t), \Theta(t - T_k(t))) + c_1 T_k(t) \\ &= \tilde{e}_k^{prev}(t) + c_1 T_k(t) \end{aligned}$$

where c_1 is a constant that is proportional to the maximum change of any queue value over one slot. With an almost identical argument, it can be shown that $\tilde{e}_k^{prev}(t) \leq \tilde{e}_k(t) + c_2 T_k(t)$, where c_2 is a constant that is proportional to the maximum change of any queue value over one slot. Thus:

$$|\tilde{e}_k(t) - \tilde{e}_k^{prev}(t)| \leq \max[c_1, c_2] T_k(t)$$

Therefore:

$$\mathbb{E} \{ |\tilde{e}_k(t) - \tilde{e}_k^{prev}(t)| \} \leq \max[c_1, c_2] \mathbb{E} \{ T_k(t) \}$$

Similarly, we can show:

$$\mathbb{E} \{ |e_k(t) - e_k^{prev}(t)| \} \leq c_3 \mathbb{E} \{ T_k(t) \}$$

Defining $d_1 \triangleq \max[c_1, c_2] + c_3$ proves the lemma. \square

It now suffices to bound $\mathbb{E} \{ |\tilde{e}_k^{prev}(t) - e_k^{prev}(t)| \}$. For a given $k \in \{1, \dots, K\}$ and a given collection of queue states $\Theta(t - T_k(t))$ at time $t - T_k(t)$, define the following function $Y(\omega)$:

$$Y(\omega) \triangleq \min_{I \in \mathcal{I}} [Y_k(I, \omega, \Theta(t - T_k(t)))] \quad (43)$$

Note that $\tilde{e}_k^{prev}(t)$ is simply an empirical average of the function $Y(\omega)$ over W i.i.d. samples $\omega_w^{(k)}(t)$ (which have distribution $F_k(\omega)$). Note that these values are also *independent of the queue state* $\Theta(t - T_k(t))$, as these samples are taken on or after time $t - T_k(t)$. Further, the value $e_k^{prev}(t)$ is simply an expected value of the random variable $Y(\omega)$ over all outcomes ω that take place with distribution $F_k(\omega)$. Hence we have reduced the problem to a pure ‘‘Law of Large Numbers’’ problem of bounding the expected difference between the exact mean of a random

variable and its empirical average over W i.i.d. samples. Because the queue backlogs $\Theta(t - T_k(t))$ are considered constant in $Y(\omega)$, we can write $Y(\omega)$ in terms of component random variables as follows (using (31)):

$$Y(\omega) = \min_{I \in \mathcal{I}} \left[VY_V(\omega) + \sum_{n=1}^N U_n(t - T_k(t))Y_{U,n}(\omega) + \sum_{m \in \tilde{\mathcal{M}}} Z_m(t - T_k(t))Y_{Z,m}(\omega) - \sum_{l=1}^L Q_l(t - T_k(t))Y_{Q,l}(\omega) \right] \quad (44)$$

where $Y_V(\omega)$, $Y_{U,n}(\omega)$, $Y_{Z,m}(\omega)$, $Y_{Q,l}(\omega)$ are random variables defined as (from (31)):

$$Y_V(\omega) \triangleq l(\hat{x}(k, \omega, I_\omega^*)) \quad (45)$$

$$Y_{U,n}(\omega) \triangleq h_n(\hat{x}(k, \omega, I_\omega^*)) \quad (46)$$

$$Y_{Z,m}(\omega) \triangleq \hat{x}_m(k, \omega, I_\omega^*) \quad (47)$$

$$Y_{Q,l}(\omega) \triangleq \hat{\mu}_l(k, \omega, I_\omega^*) - \hat{a}_l(k, \omega, I_\omega^*) \quad (48)$$

where I_ω^* is the stage-2 control action that achieves the min in (44). Now define \bar{Y}_V , $\bar{Y}_{U,n}$, $\bar{Y}_{Z,m}$, $\bar{Y}_{Q,l}$ as the expectations of the random variables in (45)-(48) over the random variable ω that has distribution $F_k(\omega)$, and define $Y_V^{(W)}$, $Y_{U,n}^{(W)}$, $Y_{Z,m}^{(W)}$, $Y_{Q,l}^{(W)}$ as the corresponding *empirical averages* over the i.i.d. samples $\omega_w^{(k)}$ (for $w \in \{1, \dots, W\}$). We thus have:

$$\begin{aligned} \tilde{e}_k^{prev}(t) - e_k^{prev}(t) &= V(Y_V^{(W)} - \bar{Y}_V) + \sum_{n=1}^N U_n(t - T_k(t))(Y_{U,n}^{(W)} - \bar{Y}_{U,n}) \\ &+ \sum_{m \in \tilde{\mathcal{M}}} Z_m(t - T_k(t))(Y_{Z,m}^{(W)} - \bar{Y}_{Z,m}) + \sum_{l=1}^L Q_l(t - T_k(t))(Y_{Q,l}^{(W)} - \bar{Y}_{Q,l}) \end{aligned}$$

and hence:

$$\begin{aligned} |\tilde{e}_k^{prev}(t) - e_k^{prev}(t)| &\leq V|Y_V^{(W)} - \bar{Y}_V| + \sum_{n=1}^N U_n(t - T_k(t))|Y_{U,n}^{(W)} - \bar{Y}_{U,n}| \\ &+ \sum_{m \in \tilde{\mathcal{M}}} |Z_m(t - T_k(t))||Y_{Z,m}^{(W)} - \bar{Y}_{Z,m}| + \sum_{l=1}^L Q_l(t)|Y_{Q,l}^{(W)} - \bar{Y}_{Q,l}| \end{aligned} \quad (49)$$

We now use the following basic lemma concerning the expected difference between an empirical average and its exact mean:

Lemma 4: Let $\{Y_w\}_{w=1}^\infty$ be an i.i.d. sequence of random variables with a general distribution with finite support, so that there are finite constants y_{min} and y_{max} such that:

$$y_{min} \leq Y_w \leq y_{max} \text{ for all } w \in \{1, 2, \dots\}$$

Define $y_{diff} \triangleq y_{max} - y_{min}$. Define \bar{Y} as the expectation of Y_1 , and define $Y^{(W)}$ as the empirical average over W samples: $Y^{(W)} \triangleq \frac{1}{W} \sum_{w=1}^W Y_w$. Then:

$$\mathbb{E} \left\{ |Y^{(W)} - \bar{Y}| \right\} \leq \frac{y_{diff}}{2\sqrt{W}}$$

Proof: The proof is straightforward and is given in Appendix D for completeness. \square

Because all penalties and cost functions are upper and lower bounded, the random variables in (45)-(48) have finite support, and we define y_{diff}^{max} as the maximum difference in the maximum and minimum possible values over all of the random variables. Using Lemma 4 in (49) yields:

$$\begin{aligned} \mathbb{E} \{ |\tilde{e}_k^{prev}(t) - e_k^{prev}(t)| \mid \Theta(t - T_k(t)), T_k(t) \} &\leq \frac{y_{diff}^{max}}{2\sqrt{W}} \left[V + \sum_{n=1}^N U_n(t - T_k(t)) \right] \\ &\quad + \frac{y_{diff}^{max}}{2\sqrt{W}} \left[\sum_{m \in \tilde{\mathcal{M}}} |Z_m(t - T_k(t))| + \sum_{l=1}^L Q_l(t - T_k(t)) \right] \\ &\leq \frac{y_{diff}^{max}}{2\sqrt{W}} \left[V + \sum_{n=1}^N U_n(t) + \sum_{m \in \tilde{\mathcal{M}}} |Z_m(t)| + \sum_{l=1}^L Q_l(t) \right] \\ &\quad + d_2 T_k(t) \end{aligned}$$

where d_2 is a constant that depends on the maximum change in queue backlog on a given slot. Taking expectations of the above and using the law of iterated expectations yields:

$$\mathbb{E} \{ |\tilde{e}_k^{prev}(t) - e_k^{prev}(t)| \} \leq \frac{y_{diff}^{max}}{2\sqrt{W}} \mathbb{E} \left\{ V + \sum_{n=1}^N U_n(t) + \sum_{m \in \tilde{\mathcal{M}}} |Z_m(t)| + \sum_{l=1}^L Q_l(t) \right\} + d_2 \mathbb{E} \{ T_k(t) \}$$

Using the above inequality with (42), (41) in (40) yields:

$$\mathbb{E} \{ \Omega(\Theta(t)) \} \leq \frac{K y_{diff}^{max}}{2\sqrt{W}} \mathbb{E} \left\{ V + \sum_{n=1}^N U_n(t) + \sum_{m \in \tilde{\mathcal{M}}} |Z_m(t)| + \sum_{l=1}^L Q_l(t) \right\} + c \sum_{k=1}^K \mathbb{E} \{ T_k(t) \} \quad (50)$$

where c is a constant that depends on the maximum possible change in queue backlogs over one slot. The random variable $T_k(t)$ can be viewed as a sum of W geometric random variables (each with mean K/θ), with the possible exception when t is small and some of the past W samples occur during the initialization time $\tau \in \{-WK, -WK + 1, \dots, -1\}$. Therefore, for all t and all k we have:

$$\mathbb{E} \{ T_k(t) \} \leq WK/\theta + WK$$

Then inequality (50) satisfies the condition (25) from Assumption 3 with:

$$\epsilon_V = \epsilon_U = \epsilon_Q = \epsilon_C = \frac{K y_{diff}^{max}}{2\sqrt{W}}, \quad C \triangleq c[WK^2/\theta + WK^2]$$

This completes the proof of Theorem 3.

APPENDIX D — PROOF OF LEMMA 4

Proof: We have:

$$\mathbb{E} \left\{ |Y^{(W)} - \bar{Y}| \right\}^2 \leq \mathbb{E} \left\{ |Y^{(W)} - \bar{Y}|^2 \right\} = \frac{\sigma^2}{W}$$

where σ^2 is the variance of Y_1 . It suffices to bound σ^2 in terms of the constants y_{min} , y_{max} , and y_{diff} . We have:

$$\begin{aligned}\sigma^2 = Var(Y_1) &= Var(Y_1 - y_{min}) \\ &= \mathbb{E}\{(Y_1 - y_{min})^2\} - (\bar{Y} - y_{min})^2 \\ &\leq \mathbb{E}\{(y_{max} - y_{min})(Y_1 - y_{min})\} - (\bar{Y} - y_{min})^2\end{aligned}\tag{51}$$

$$\begin{aligned}&= (\bar{Y} - y_{min})(y_{max} - y_{min} - (\bar{Y} - y_{min})) \\ &= (\bar{Y} - y_{min})(y_{max} - \bar{Y})\end{aligned}\tag{52}$$

where (51) holds because $Y_1 - y_{min} \geq 0$. To compute the final bound on the expression in (52), note that $y_{min} \leq \bar{Y} \leq y_{max}$, and the maximum of the function $f(x) = (x - y_{min})(y_{max} - x)$ over the interval $y_{min} \leq x \leq y_{max}$ is equal to $(y_{max} - y_{min})^2/4$. Thus, $\sigma^2 \leq y_{diff}^2/4$. \square

APPENDIX E — PROOF OF THEOREM 4

The proof of Theorem 3 can be followed in the same way, with the exception that the fixed value W is replaced by the random value $W(t)$ (which may be correlated with queue states). Therefore, repeating the proof in Appendix C, the result of (50) translates to:

$$\mathbb{E}\{\Omega(\Theta(t))\} \leq \frac{Ky_{diff}^{max}}{2} \mathbb{E}\left\{\frac{V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t)}{\sqrt{W(t)}}\right\} + c \sum_{k=1}^K \mathbb{E}\{T_k(t)\}$$

Each term $\mathbb{E}\{T_k(t)\}$ can be bounded by $\hat{W}(t)K/\theta + W_0K$. The final term can thus be bounded as follows:

$$c \sum_{k=1}^K \mathbb{E}\{T_k(t)\} \leq \hat{W}(t)K^2/\theta + W_0K^2$$

where $\hat{W}(t) \triangleq (t+1)^{\beta_1}$. Define $C_1(t) \triangleq \hat{W}(t)K^2/\theta + W_0K^2$.

It is not difficult to show that $W_{rand}(t)$ satisfies:

$$\lim_{t \rightarrow \infty} \frac{W_{rand}(t)}{t} = \frac{\theta}{K} \quad \text{with probability 1}$$

However, $\hat{W}(t)$ increases sub-linearly with t . Therefore, because $W(t) \triangleq \min[\hat{W}(t), W_{rand}(t)]$, we have:

$$\lim_{t \rightarrow \infty} Pr[W(t) \neq \hat{W}(t)] = 0$$

Furthermore, because $W_{rand}(t)$ is simply the min of K delayed renewal processes $W_1(t), \dots, W_K(t)$ (each having i.i.d. geometric inter-arrival times with mean K/θ), we have by the union bound:

$$Pr[W(t) \neq \hat{W}(t)] = Pr[\min[W_1(t), \dots, W_K(t)] < (t+1)^{\beta_1}] \leq K Pr[W_1(t) \leq (t+1)^{\beta_1}]$$

It follows that:

$$\lim_{t \rightarrow \infty} t Pr[W(t) \neq \hat{W}(t)] = 0$$

Therefore:

$$\begin{aligned} & \frac{Ky_{diff}^{max}}{2} \mathbb{E} \left\{ \frac{V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t)}{\sqrt{W(t)}} \right\} \leq \\ & \frac{Ky_{diff}^{max}}{2\sqrt{\hat{W}(t)}} \mathbb{E} \left\{ V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t) \mid W(t) = \hat{W}(t) \right\} Pr[\hat{W}(t) = W(t)] \\ & + \frac{Ky_{diff}^{max}}{2} \mathbb{E} \left\{ V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t) \mid W(t) \neq \hat{W}(t) \right\} Pr[\hat{W}(t) \neq W(t)] \end{aligned}$$

where we have used the fact that $W(t) \geq 1$ always. Adding the (non-negative) conditional expectation to complete the first term on the right hand side yields:

$$\begin{aligned} & \frac{Ky_{diff}^{max}}{2} \mathbb{E} \left\{ \frac{V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t)}{\sqrt{W(t)}} \right\} \leq \\ & \frac{Ky_{diff}^{max}}{2\sqrt{\hat{W}(t)}} \mathbb{E} \left\{ V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t) \right\} \\ & + \frac{Ky_{diff}^{max}}{2} \mathbb{E} \left\{ V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t) \mid W(t) \neq \hat{W}(t) \right\} Pr[\hat{W}(t) \neq W(t)] \\ & \leq \frac{Ky_{diff}^{max}}{2\sqrt{\hat{W}(t)}} \mathbb{E} \left\{ V + \sum_n U_n(t) + \sum_m |Z_m(t)| + \sum_l Q_l(t) \right\} \\ & \quad + \frac{Ky_{diff}^{max} c_0 t}{2} Pr[\hat{W}(t) \neq W(t)] \end{aligned}$$

where c_0 is a constant that is proportional to the maximum change in any queue over one slot. Because $tPr[\hat{W}(t) \neq W(t)] \rightarrow 0$ as $t \rightarrow \infty$, there exists a time t_0 such that for all $t \geq t_0$ we have:

$$\frac{Ky_{diff}^{max} c_0 t}{2} Pr[\hat{W}(t) \neq W(t)] \leq 1$$

We can now define $C(t) \triangleq C_1(t) + 1$ for use in Theorem 2 (note that $C(t) \leq O((t - t_0 + 1)^{\beta_1})$). Further define:

$$\epsilon_x(t) \triangleq \frac{Ky_{diff}^{max}}{2\sqrt{\hat{W}(t)}}$$

for $x \in \{V, U, Z, Q\}$. This satisfies the assumptions of Theorem 2, proving the result.

REFERENCES

- [1] L. Georgiadis, M. J. Neely, and L. Tassiulas. Resource allocation and cross-layer control in wireless networks. *Foundations and Trends in Networking*, vol. 1, no. 1, pp. 1-149, 2006.
- [2] A. Stolyar. Maximizing queueing network utility subject to stability: Greedy primal-dual algorithm. *Queueing Systems*, vol. 50, pp. 401-457, 2005.
- [3] M. J. Neely, E. Modiano, and C. Li. Fairness and optimal stochastic control for heterogeneous networks. *Proc. IEEE INFOCOM*, March 2005.
- [4] M. J. Neely. Energy optimal control for time varying wireless networks. *IEEE Transactions on Information Theory*, vol. 52, no. 7, pp. 2915-2934, July 2006.
- [5] C. Li and M. J. Neely. Energy-optimal scheduling with dynamic channel acquisition in wireless downlinks. *Proc. of 46th IEEE Conf. on Dec. and Control (CDC)*, Dec. 2007.

- [6] A. Gopalan, C. Caramanis, and S. Shakkottai. On wireless scheduling with partial channel-state information. *Allerton Conf. on Comm., Control, and Computing*, Sept. 2007.
- [7] M. J. Neely and R. Urgaonkar. Optimal backpressure routing in wireless networks with multi-receiver diversity. *Ad Hoc Networks (Elsevier)*, 2008, doi:10.1016/j.adhoc.2008.07.009.
- [8] M. J. Neely and R. Urgaonkar. Opportunism, backpressure, and stochastic optimization with the wireless broadcast advantage. *Asilomar Conference on Signals, Systems, and Computers, Pacific Grove, CA*, Oct. 2008.
- [9] R. Agrawal and V. Subramanian. Optimality of certain channel aware scheduling policies. *Proc. 40th Annual Allerton Conference on Communication, Control, and Computing, Monticello, IL*, Oct. 2002.
- [10] P. Chaporkar and A. Proutiere. Optimal joint probing and transmission strategy for maximizing throughput in wireless systems. *IEEE Journal in Selected Areas in Communication*, vol. 26, no. 8, pp. 1546-1555, Oct. 2008.
- [11] T. Javidi, B. Krishnamachari, Q. Zhao, and M. Liu. Optimality of myopic sensing in multi-channel opportunistic access. *IEEE ICC, Beijing, China*, May 2008.
- [12] N. B. Chang and M. Liu. Optimal channel probing and transmission scheduling for opportunistic spectrum access. *Proc. of ACM International Conference on Mobile Computing and Networking (MobiCom)*, Sept. 2007.
- [13] D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, Mass, 1996.
- [14] B. J. Oommen and J. K. Lanctot. Discretized pursuit learning automata. *IEEE Trans. on Systems, Man, and Cybernetics*, vol.20, no.4, pp.931-938, July/August 1990.
- [15] M. A. Haleem and R. Chandramouli. Adaptive stochastic iterative rate selection for wireless channels. *IEEE Communications Letters*, vol. 8, no. 5, pp. 292-294, May 2004.
- [16] D. P. Bertsekas, A. Nedic, and A. E. Ozdaglar. *Convex Analysis and Optimization*. Boston: Athena Scientific, 2003.
- [17] R. Berry and R. Gallager. Communication over fading channels with delay constraints. *IEEE Transactions on Information Theory*, vol. 48, no. 5, pp. 1135-1149, May 2002.
- [18] M. J. Neely. Optimal energy and delay tradeoffs for multi-user wireless downlinks. *IEEE Transactions on Information Theory*, vol. 53, no. 9, pp. 3095-3113, Sept. 2007.
- [19] A. Eryilmaz, R. Srikant, and J. R. Perkins. Stable scheduling policies for fading wireless channels. *IEEE Transactions on Networking*, vol. 13, no. 2, April 2005.
- [20] M. J. Neely. *Dynamic Power Allocation and Routing for Satellite and Wireless Networks with Time Varying Channels*. PhD thesis, Massachusetts Institute of Technology, LIDS, 2003.