# Optimal Control for Generalized Network-Flow Problems

Abhishek Sinha<sup>10</sup> and Eytan Modiano, Fellow, IEEE

Abstract—We consider the problem of throughput-optimal packet dissemination, in the presence of an arbitrary mix of unicast, broadcast, multicast, and anycast traffic, in an arbitrary wireless network. We propose an online dynamic policy, called Universal Max-Weight (UMW), which solves the problem efficiently. To the best of our knowledge, UMW is the first known throughput-optimal policy of such versatility in the context of generalized network flow problems. Conceptually, the UMW policy is derived by relaxing the precedence constraints associated with multi-hop routing and then solving a min-cost routing and max-weight scheduling problem on a virtual network of queues. When specialized to the unicast setting, the UMW policy yields a throughput-optimal cycle-free routing and link scheduling policy. This is in contrast with the well-known throughput-optimal backpressure (BP) policy which allows for packet cycling, resulting in excessive latency. Extensive simulation results show that the proposed UMW policy incurs a substantially smaller delay as compared with the BP policy. The proof of throughput-optimality of the UMW policy combines ideas from the stochastic Lyapunov theory with a sample path argument from adversarial queueing theory and may be of independent theoretical interest.

*Index Terms*— Throughput-optimal policies, generalized flows, queueing theory.

#### I. INTRODUCTION

THE Generalized Network Flow problem involves efficient transportation of messages, generated at the source node(s), to a set of designated destination node(s) over a multi-hop network. Depending on the number of destination nodes associated with each source node, the problem is known either as *unicast* (single destination node), *broadcast* (all nodes are destination nodes), *multicast* (some nodes are destination nodes) or *anycast* (several choices for a single destination node). Over the last few decades, a tremendous amount of research effort has been directed to address each of the above problems in different networking contexts. However, despite the increasingly diverse mix of internet traffic, to the best of our knowledge, there exists no universal solution to the general problem, only isolated solutions that do not interoperate and

Manuscript received December 21, 2016; revised June 6, 2017 and November 17, 2017; accepted December 12, 2017; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor D. Leith. Date of publication December 29, 2017; date of current version February 14, 2018. This work was supported in part by the NSF under Grant CNS-1217048 and Grant CNS-1524317 and in part by the DARPA I2O and Raytheon BBN Technologies under Contract HROO II-1 5-C-0097. Part of the paper appeared in the proceedings of INFOCOM, 2017, IEEE [1]. (*Corresponding author: Abhishek Sinha.*)

The authors are with the Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA 02139 USA (e-mail: sinhaa@mit.edu; modiano@mit.edu).

Digital Object Identifier 10.1109/TNET.2017.2783846

are often suboptimal. In this paper, we provide the first such universal solution: A throughput optimal dynamic control policy for the generalized network flow problem.

We start with a brief discussion of the above networking problems and then survey the relevant literature.

In the Broadcast problem, packets generated at a source need to be distributed among all nodes in the network. In the classic paper of Edmonds [2], the broadcast capacity of a wired network is derived and an algorithm is proposed to compute the maximum number of edge-disjoint spanning trees, which together achieve the maximum broadcast throughput. The algorithm in [2] is combinatorial in nature and does not have a wireless counterpart, with associated interference-free edge activations. Following Edmonds' work, a variety of different broadcast algorithms have been proposed in the literature, each one targeted to optimize different metrics such as delay [3], energy consumption [4] and fault-tolerance [5]. In the context of optimizing throughput, [6] proposes a randomized broadcast policy, which is optimal for wired networks. However, extending this algorithm to the wireless setting proves to be difficult [7]. Sinha et al. [8] propose an optimal broadcast algorithm for a wireless network, albeit with exponential complexity. In a recent series of papers [9], [10], a simple throughput-optimal broadcast algorithm has been proposed for wireless networks with an underlying DAG topology. However, this algorithm does not extend to non-DAG networks.

The Multicast problem is a generalization of the broadcast problem, in which the packets generated at source nodes needs to be efficiently distributed to a subset of nodes in the network. In its combinatorial version, the multicast problem reduces to finding the maximum number of edge-disjoint trees, spanning the source node and destination nodes. This problem is known as the Steiner Tree Packing problem, which is NP-hard [11]. Numerous algorithms have been proposed in the literature for solving the multicast problem. In [12] and [13], back-pressure type algorithms are proposed for multicasting over wired and wireless networks respectively. These algorithms forward packets over a set of pre-computed distribution trees and are limited to the throughput obtainable by these trees. Moreover, computing and maintaining these trees is impractical in large and time-varying networks. We note that because of the need for packet duplications, the Multicast and Broadcast problems do not satisfy standard flow conservation constraints, and thus the design of throughput-optimal algorithms is non-trivial.

The **Unicast** problem involves a single source and a single destination. The celebrated Back-Pressure (BP) algorithm [14]

1063-6692 © 2017 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications\_standards/publications/rights/index.html for more information.

was proposed for the unicast problem. In this algorithm, the routing and scheduling decisions are taken based on local queue length differences. As a result, BP explores all possible paths for routing and usually takes a long time for convergence, resulting in considerable latency, especially in lightly loaded networks. Subsequently, a number of refinements have been proposed to improve the delay characteristics of the BP algorithm. In [15] BP is combined with hop length based shortest path routing for faster route discovery, and [16] proposes a second order algorithm using the Hessian matrix to improve delay.

The **Anycast** problem involves routing from a single source to *any one* of the several given destinations. Anycast is increasingly used in Content-Distribution Networks (CDNs) for optimally distributing geo-replicated contents [17].

Our proposed solution uses a *virtual network of queues* - one virtual queue per link in the network. We solve the routing problem dynamically using a simple "weighted-shortest-route" computation on the virtual network and using the corresponding route on the physical network. Optimal link scheduling is performed by a max-weight computation, also in the virtual network, and then using the resulting activation in the physical network. The overall algorithm is dynamic, cycle-free, and solves the generalized routing and scheduling problem optimally (i.e., maximally stable or throughput optimal). In addition to this, the proposed **UMW** policy has the following advantages:

- Generalized Solution: Unlike the BP policy, which solves only the unicast problem, the proposed UMW policy efficiently addresses all of the aforementioned network flow problems in both wired and wireless networks in a very general setting.
- 2) Delay Reduction: Although the celebrated BP policy is throughput-optimal, its average delay performance is known to be poor due to the occurrence of packetcycling in the network [15], [18]. In our proposed UMW policy, each packet traverses a dynamically selected *acyclic* route, which drastically reduces the average latency.
- 3) State-Complexity Reduction: Unlike the BP policy, which maintains *per-flow* queues at each node, the proposed UMW policy maintains only a virtual queue counter and a priority queue per link, *irrespective of the number and type of flows in the network*. This reduces the amount of data-structure overhead that needs to be maintained for efficient operation (c.f. [18]).

Recently, the techniques introduced in this paper have been utilized in [19] to solve the problem of wireless broadcast with point-to-multipoint transmissions.

The rest of the paper is organized as follows: In Section II we discuss the basic system model and formulate the problem. In Section III we give a brief overview of the proposed **UMW** policy. Section IV discusses the structure and dynamics of the virtual queues, on which **UMW** is based. In Section V we prove its stability in a multi-hop physical network. In Section VII we propose a distributed heuristic policy derived from the UMW policy. Section VIII provides extensive simulation results, comparing **UMW** with other



Fig. 1. A wireless network and its two maximal feasible link activations under the primary interference constraint. (a) A wireless network. (b) Activation vector s1. (c) Activation vector s2.

competing policies. Section IX concludes the paper with a few directions for further research.

#### II. SYSTEM MODEL AND PROBLEM FORMULATION

## A. Network Model

We consider a wireless network with arbitrary topology, represented by the graph  $\mathcal{G}(V, E)$ . The network consists of |V| = n nodes and |E| = m links. Time is slotted. A link, if activated, can transmit one packet per slot. Due to wireless interference constraints, only certain subsets of links may be activated together at any slot. The set of all admissible link activations is known as the activation set and is denoted by  $\mathcal{M} \subseteq 2^E$ . We do not impose any restriction on the structure of the activation set  $\mathcal{M}$ . As an example, in the case of *node*exclusive or primary interference constraint [20], the activation set  $\mathcal{M}_{primary}$  consists of the set of all matchings [21] in the graph  $\mathcal{G}(V, E)$ . Wired networks are a special case of the above model, where the activation set  $\mathcal{M}_{wired} = 2^E$ . In other words, in wired networks, packets can be transmitted over all links simultaneously. See Figure 1 for an example of a wireless network with primary interference constraints.

For simplicity in exposition, in the following, we assume that the network topology is static, with bounded arrivals per slot. However, the proposed policy and its analysis apply even if the network is time-varying. Moreover, an attractive feature of our policy is that time-varying networks *do not incur* any extra computational overhead in implementation. The performance of the proposed policy in time-varying networks is evaluated numerically in Section VIII.

## B. Traffic Model

In this paper, we consider the *Generalized Network Flow* problem, where incoming packets at a source node are to be distributed among an arbitrary set of destination nodes in a multi-hop fashion. Formally, the set of all distinct classes of incoming traffic is denoted by C. A class c traffic is identified by its source node  $s^{(c)} \in V$  and the set of its required destination nodes  $\mathcal{D}^{(c)} \subseteq V$ . As explained below, by varying the structure of the destination set  $\mathcal{D}^{(c)}$  of class *c*, this general framework yields the following four fundamental flow problems as special cases:

- UNICAST: All class c packets, arriving at a source node  $s^{(c)}$ , are required to be delivered to a single destination node  $\mathcal{D}^{(c)} = \{t^{(c)}\}.$
- **BROADCAST**: All class c packets, arriving at a source node  $s^{(c)}$ , are required to be delivered to all nodes in the network, i.e.,  $\mathcal{D}^{(c)} = V$ .
- **MULTICAST:** All class c packets, arriving at a source node  $s^{(c)}$ , are required to be delivered to a proper subset of nodes  $\mathcal{D}^{(c)} = \{t_1^{(c)}, t_2^{(c)}, \dots, t_k^{(c)}\} \subseteq V$ .
- ANYCAST: A Packet of class c, arriving at a source node s<sup>(c)</sup>, is required to be delivered to any one of a given set of k nodes D<sup>(c)</sup> = t<sub>1</sub><sup>(c)</sup> ⊕ t<sub>2</sub><sup>(c)</sup> ⊕ ... ⊕ t<sub>k</sub><sup>(c)</sup>. Thus the anycast problem is similar to the unicast problem, with all destinations forming a single super destination node.

Arrivals are i.i.d. in every slot, with  $A^{(c)}(t)$  packets from class c arriving at the source node  $s^{(c)}$  at slot t. The mean rate of arrival for class c is  $\mathbb{E}A^{(c)}(t) = \lambda^{(c)}$ . The arrival rate to the network is characterized by the vector  $\lambda = \{\lambda^{(c)}, c \in C\}$ . The total number of external packet arrivals to the entire network at any slot t is assumed to be bounded by a finite number  $A_{\text{max}}$ .

## C. Policy-Space

An admissible policy  $\pi$  for the generalized network flow problem executes the following two actions at every slot t:

- LINK ACTIVATIONS: Activating a subset of interferencefree links s(t) from the activation set  $\mathcal{M}$ .
- PACKET DUPLICATIONS AND FORWARDING: Possibly duplicating<sup>1</sup> and forwarding packets over the activated links. Due to the link capacity constraint, at most one packet may be transmitted over an active link per slot.

The set of all admissible policies is denoted by  $\Pi$ . The set  $\Pi$  is unconstrained otherwise and includes policies which may use all past and future packet arrival information.

A policy  $\pi \in \Pi$  is said to support an arrival rate vector  $\lambda$  if, under the action of the policy  $\pi$ , the destination nodes of any class c receive distinct class c packets at the rate  $\lambda^{(c)}$ ,  $c \in C$ . Formally, let  $R^{(c)}(t)$  denote the number of distinct class-c packets received in common by all destination nodes  $i \in \mathcal{D}^{(c)}$ , a under the action of the policy  $\pi$ , up to time t.

Definition 1: [Policy Supporting Rate Vector  $\lambda$ ]: A policy  $\pi \in \Pi$  is said to support an arrival rate vector  $\lambda$  if

$$\liminf_{t \to \infty} \frac{R^{(c)}(t)}{t} = \lambda^{(c)}, \quad \forall c \in \mathcal{C}, \text{ w.p.1}$$
(1)

The network-layer capacity region  $\Lambda(\mathcal{G}, \mathcal{C})^3$  is defined to be the set of all supportable rates, i.e.,

$$\Lambda(\mathcal{G},\mathcal{C}) \stackrel{\text{def}}{=} \{ \lambda \in \mathbb{R}^{|\mathcal{C}|}_+ : \exists \pi \in \Pi \text{ supporting } \lambda \}$$
(2)

Clearly, the set  $\Lambda(\mathcal{G}, \mathcal{C})$  is *convex* (using the usual *time-sharing* argument). A policy  $\pi^* \in \Pi$ , which supports any arrival rate  $\lambda$  in the interior of the capacity region  $\Lambda(\mathcal{G}, \mathcal{C})$ , is called a *throughput-optimal* policy.

As a side note, unlike [14], which establishes positive recurrence of the Markovian queue-lengths under the BP policy, Definition (1) is concerned with rate stability only. The reason for this choice will become clear in the next section where we will see that unlike BP, the stochastic dynamics of the physical queues are non-Markovian under the proposed UMW policy.

#### D. Admissible Routes of Packets

We will design a throughput-optimal policy, which delivers a packet p to any node in the network *at most* once.<sup>4</sup> This immediately implies that the set of all admissible routes  $\mathcal{T}^{(c)}$ for packets of any class c, in general, comprises of trees rooted at the corresponding source node  $s^{(c)}$ . In particular, depending on the type of class c traffic, the topology of the admissible routes  $\mathcal{T}^{(c)}$  takes the following special forms:

- UNICAST TRAFFIC:  $T^{(c)} = \text{set of all } s^{(c)} t^{(c)}$ paths in the graph  $\mathcal{G}$ .
- **BROADCAST TRAFFIC**:  $T^{(c)}$  = set of all spanning trees in the graph  $\mathcal{G}$ , rooted at  $s^{(c)}$ .
- MULTICAST TRAFFIC:  $\mathcal{T}^{(c)} = \text{set of all Steiner}$ trees [11] in  $\mathcal{G}$ , rooted at  $s^{(c)}$  and spanning the vertices  $\mathcal{D}^{(c)} = \{t_1^{(c)}, t_2^{(c)}, \dots, t_k^{(c)}\}.$
- ANYCAST TRAFFIC:  $\mathcal{T}^{(c)}$  = union of all  $s^{(c)} t_i^{(c)}$  paths in the graph  $\mathcal{G}$ ,  $i = 1, 2, \dots, k$ .

## E. Characterization of the Network-Layer Capacity Region

Consider any arrival vector  $\lambda \in \Lambda(\mathcal{G}, \mathcal{C})$ . By definition, there exists an admissible policy  $\pi \in \Pi$ , which supports the arrival rate  $\lambda$  by means of storing, duplicating and forwarding packets efficiently. Taking time-averages over the actions of the policy  $\pi$ , it is clear that there exists a *randomized* flowdecomposition and scheduling policy to route the packets such that none of the edges in the network are overloaded. Indeed, in the following theorem, we show that for every  $\lambda \in \Lambda(\mathcal{G}, \mathcal{C})$ , there exist non-negative scalars  $\{\lambda_i^{(c)}\}$ , indexed by the admissible routes  $T_i^{(c)} \in \mathcal{T}^{(c)}$  and a convex combination of the link activation vectors  $\overline{\mu} \in \operatorname{conv}(\mathcal{M})$  such

 $<sup>^{1}</sup>$ In order to transmit a packet over multiple downstream links (*e.g.* in Broadcast or Multicast), the sender must duplicate the packet and send the copies to the respective downstream link buffers.

<sup>&</sup>lt;sup>2</sup>To be precise, the *super*-destination node in case of Anycast.

<sup>&</sup>lt;sup>3</sup>Note that, Network-layer capacity region is, in general (e.g. multicast), different from the Information-Theoretic capacity region [22].

<sup>&</sup>lt;sup>4</sup>This should be contrasted with the popular throughput-optimal unicast policy *Back-Pressure* [14], which does not satisfy this constraint and may deliver the same packet to a node multiple times, thus potentially degrading its delay performance.

that,

$$\lambda^{(c)} = \sum_{T^{(c)} \in \mathcal{T}^{(c)}} \lambda_i^{(c)}, \quad \forall c \in \mathcal{C}$$
(3)

$$\lambda_e \stackrel{\text{(def.)}}{=} \sum_{\substack{(i,c): e \in T^{(c)}: T^{(c)} \in \mathcal{T}^{(c)}}} \lambda_i^{(c)} \leq \overline{\mu}_e, \quad \forall e \in E.$$
(4)

Eqn. (3) denotes decomposition of the average incoming flows into different admissible routes and Eqn. (4) denotes the fact that none of the edges in the network are overloaded, i.e. arrival rate of packets to any edge e under the policy  $\pi$ is *at most* the rate allocated by the policy  $\pi$  to the edge e to serve packets.

To state the result precisely, define the set  $\overline{\Lambda}$  to be the set of all arrival vectors  $\lambda \in \mathbb{R}^{|\mathcal{C}|}_+$ , for which there exists a randomized activation vector  $\mu \in \operatorname{conv}(\mathcal{M})$  and a non-negative flow decomposition  $\{\lambda_i^{(c)}\}$ , such that Eqns. (3) and (4) are satisfied. We have the following theorem:

Theorem 1: The network-layer capacity region  $\Lambda(\mathcal{G}, \mathcal{C})$  is characterized by the set  $\overline{\Lambda}$ , up to its boundary, *i.e.*,

 $\operatorname{int}(\overline{\Lambda}) \subseteq \Lambda(\mathcal{G}, \mathcal{C}) \subseteq \overline{\Lambda},$ 

where int(S) denotes the interior of the Euclidean set S.

Proof of Theorem 1 consists of two parts: converse and achievability. Proof of the *converse* is given in Appendix A, where we show that all supportable arrival rates must belong to the set  $\overline{\Lambda}$ . The main result of this paper, as developed in the subsequent sections, is the construction of an efficient admissible policy, called Universal Max-Weight (UMW), which achieves *any* arrival rate in the interior of the set  $\overline{\Lambda}$ .

## III. OVERVIEW OF THE UMW POLICY

In this section, we present a brief overview of our throughput-optimal UMW policy, designed and analyzed in the subsequent sections. Central to the UMW policy is a global state vector called virtual queues Q(t), used for packet routing and link activations. Each component of the virtual queues is updated at every slot according to a onehop queueing (Lindley) recursion, corresponding to a *relaxed* network, described in detail in section IV. Unlike the wellknown Back-Pressure algorithm for the unicast problem [14], in which packet routing decisions are made hop-by-hop using physical queue lengths Q(t), the UMW policy prescribes an admissible route to each incoming packet immediately upon its arrival (dynamic source routing). This route selection decision is dynamically made by solving a suitable min-cost routing problem (e.g., shortest path, MST etc.) at the source with edge costs given by the current virtual-queue vector Q(t). Link activation decisions at each slot are made by a Max-Weight algorithm with link-weights set equal to Q(t). Having fixed the routing and activation policy as above, in section V we design a packet scheduling algorithm for the physical network, which efficiently resolves contention among multiple packets that wait to cross the same (active) edge at the same slot. We show that the overall policy is throughput-optimal. One significantly new feature of our algorithm is that it is entirely oblivious to the length of the physical queues of the network and utilizes the auxiliary virtual-queue state variables for stabilizing the former.

Our proof of throughput-optimality of **UMW** leverages ideas from *deterministic* adversarial queueing theory and combines it effectively with the *stochastic* Lyapunov-drift based techniques and may be of independent theoretical interest.

# IV. GLOBAL VIRTUAL QUEUES: STRUCTURES, Algorithms, and Stability

Here we introduce the notion of *virtual queues*,<sup>5</sup> which is obtained by *relaxing* the dynamics of the physical queues of the network in the following intuitive fashion.

## A. Precedence Constraints

In a multi-hop network, if a packet p is being routed along the path  $T = l_1 - l_2 - \ldots - l_k$ , where  $l_i \in E$  is the  $i^{\text{th}}$  link on its path, then by the principle of causality, the packet pcannot be physically transmitted over the  $j^{\text{th}}$  link  $l_j$  if it has not already been transmitted by the first j-1 links  $l_1, l_2, \ldots, l_{j-1}$ . This constraint is known as the *precedence constraint* in the network scheduling literature [24]. In the following, we make a radical departure by relaxing this constraint to obtain a simpler single-hop virtual system, which will play a key role in designing our policy and its optimality analysis.

# B. The Virtual Queue Process $\{\tilde{Q}(t)\}_{t>1}$

The Virtual queue process  $\tilde{\mathbf{Q}}(t) = (\tilde{Q}_e(t), e \in E)$ is an |E| = m dimensional controlled stochastic process, imitating a fictitious queueing network without the precedence constraints. In particular, when a packet p of class c arrives at the source node  $s^{(c)}$ , a dynamic policy  $\pi$  prescribes a suitable route  $T^{(c)}(t) \in \mathcal{T}^{(c)}$  to the packet. Denoting the set of all edges in the route  $T^{(c)}(t)$  by  $\{l_1, l_2, \ldots, l_k\}$ , this incoming packet induces a virtual arrival simultaneously at each of the virtual queues  $(\tilde{Q}_{l_i}), i = 1, 2, \ldots, k$ , right upon its arrival to the source. Since the virtual network is assumed to be relaxed with no precedence constraints, any packet present in the virtual queue is eligible for service. See Figure 2 for an illustration.

The (controlled) service process allocated to the virtual queues is denoted by  $\{\mu^{\pi}(t)\}_{t\geq 1}$ . We require the service process to satisfy the same activation constraints as in the original system, i.e.,  $\mu^{\pi}(t) \in \mathcal{M}, \forall t \geq 1$ .

Let  $A_e^{\pi}(t)$  be the total number of virtual packet arrivals (from all classes) to the queue  $\tilde{Q}_e$  at time t under the action of the policy  $\pi$ , i.e.,

$$A_e^{\pi}(t) = \sum_{c \in \mathcal{C}} A^{(c)}(t) \mathbb{1}\left(e \in T^{(c)}(t)\right), \quad \forall e \in E.$$
(5)

Hence, we have the following one-step evolution (Lindley recursion) of the virtual queue process  $\{\tilde{Q}_e(t)\}_{t\geq 1}$ :

$$\tilde{Q}_{e}(t+1) = \left(\tilde{Q}_{e}(t) + A_{e}^{\pi}(t) - \mu_{e}^{\pi}(t)\right)^{+}, \ \forall e \in E,$$
 (6)

<sup>5</sup>Note that our notion of *virtual queues* is completely different from and unrelated to the notion of *shadow-queues* proposed earlier in [13] and [18], and *virtual queues* proposed in [23].



Fig. 2. Illustration of the virtual queue system for the four-node network  $\mathcal{G}$ . Upon arrival, the incoming packet p, belonging to a unicast session from node 1 to 4, is prescribed a path  $\mathcal{T}_p = \{\{1,2\},\{2,3\},\{3,4\}\}$ . Relaxing the precedence constraints, the packet p is counted as an arrival to the virtual queues  $\tilde{Q}_{12}$  and  $\tilde{Q}_{23}$  and  $\tilde{Q}_{34}$  simultaneously at the *same slot*. In the physical system, the packet p may take a while before reaching any edge in its path, depending on the control policy.

We emphasize that  $A_e^{\pi}(t)$  is a function of the routing tree  $T^{(c)}(t)$  that the policy chooses at time t, from the set of all admissible routes  $\mathcal{T}^{(c)}$ . This is discussed in the following.

## C. Dynamic Control and Stability of the Virtual Queues

Next, we design a dynamic routing and link activation policy for the virtual network, which stabilizes the virtual queue process  $\{\tilde{Q}(t)\}_{t\geq 1}$ , for all arrival rate-vectors  $\lambda \in \operatorname{int}(\overline{\Lambda})$ . This policy is obtained by minimizing the one-step drift of a quadratic Lyapunov-function of the *virtual queue lengths* (as opposed to the real queue lengths used in the Back-Pressure policy [14]). In the following section, we will show that when this dynamic policy is used in conjunction with a suitable packet scheduling policy in the physical network, the overall policy is throughput-optimal for the physical network.

To derive a stabilizing policy for the virtual network, consider a quadratic Lyapunov function  $L(\tilde{Q}(t))$  defined in terms of the virtual queue lengths:

$$L(\tilde{\pmb{Q}}(t)) = \sum_{e \in E} \tilde{Q}_e^2(t)$$

From the one-step dynamics of the virtual queues (6), we have:

$$\begin{split} \tilde{Q}_e(t+1)^2 &\leq (\tilde{Q}_e(t) - \mu_e^{\pi}(t) + A_e^{\pi}(t))^2 \\ &= \tilde{Q}_e^2(t) + (A_e^{\pi}(t))^2 + (\mu_e^{\pi}(t))^2 + 2\tilde{Q}_e(t)A_e^{\pi}(t) \\ &- 2\tilde{Q}_e(t)\mu_e^{\pi}(t) - 2\mu_e^{\pi}(t)A_e^{\pi}(t) \end{split}$$

Since  $\mu_e^{\pi}(t) \ge 0$  and  $A_e^{\pi}(t) \ge 0$ , we have

$$\begin{split} \tilde{Q}_e^2(t+1) - \tilde{Q}_e^2(t) &\leq (A_e^{\pi}(t))^2 + (\mu_e^{\pi}(t))^2 \\ &+ 2\tilde{Q}_e(t)A_e^{\pi}(t) - 2\tilde{Q}_e(t)\mu_e^{\pi}(t) \end{split}$$

Hence, the one-step Lyapunov drift  $\Delta^{\pi}(t)$ , conditional on the current virtual queue lengths  $\tilde{Q}(t)$ , under the operation of

any admissible Markovian policy  $\pi \in \Pi$  is upper-bounded by

$$\Delta^{\pi}(t) \stackrel{\text{def}}{=} \mathbb{E} \left( L(\tilde{\boldsymbol{Q}}(t+1)) - L(\tilde{\boldsymbol{Q}}(t)) | \tilde{\boldsymbol{Q}}(t) \right)$$

$$\leq B + 2 \sum_{e \in E} \tilde{Q}_{e}(t) \mathbb{E} \left( A_{e}^{\pi}(t) | \tilde{\boldsymbol{Q}}(t) \right)$$

$$- 2 \sum_{e \in E} \tilde{Q}_{e}(t) \mathbb{E} \left( \mu_{e}^{\pi}(t) | \tilde{\boldsymbol{Q}}(t) \right)$$
(7)

where B is a constant, bounded by  $\sum_e \mathbb{E}(A_e^{\pi}(t))^2 + \mathbb{E}(\mu_e^{\pi}(t))^2) \leq n^2 A_{\max}^2 + m.$ 

The upper-bound on the drift, given by (7), holds good for any admissible policy in the virtual network. In particular, by minimizing the upper-bound pointwise, and exploiting the separable nature of the objective, we derive the following decoupled dynamic routing and link activation policy for the virtual network:

Dynamic Routing Policy: The drift-minimizing route for each class c, over the set of all admissible routes, is selected by minimizing the following cost function, appearing in the middle of Eqn. (7)

$$\mathsf{RoutingCost}^\pi \equiv \sum_{e \in E} \tilde{Q}_e(t) A_e^\pi(t),$$

where we remind the reader that  $A_e^{\pi}(t)$  denotes the routing policy-dependent arrivals to the virtual queue corresponding to the link *e* at time *t*.

Using Eqn. (5), we may rewrite the objective-function as

$$\operatorname{RoutingCost}^{\pi} = \sum_{c \in \mathcal{C}} A^{(c)}(t) \bigg( \sum_{e \in E} \tilde{Q}_e(t) \mathbb{1} \big( e \in T^{(c)}(t) \big) \bigg)$$
(8)

Using the separability of the objective (8), the above optimization problem decomposes into following min-cost routeselection problem  $T_{opt}^{(c)}(t)$  for each class c:

$$T_{\text{opt}}^{(c)}(t) \in \operatorname*{arg\,min}_{T^{(c)} \in \mathcal{T}^{(c)}} \left( \sum_{e \in E} \tilde{Q}_e(t) \mathbb{1}\left(e \in T^{(c)}\right) \right)$$
(9)

Depending on the type of flow of class c, the route-selection problem (9) is equivalent to one of the following wellknown combinatorial problems on the graph  $\mathcal{G}$ , with its edges weighted by the virtual queue length vector  $\tilde{Q}$ :

- UNICAST TRAFFIC:  $T_{opt}^{(c)}(t) =$  The shortest  $s^{(c)} t^{(c)}$  path in the weighted-graph  $\mathcal{G}$ .
- **BROADCAST TRAFFIC:**  $T_{opt}^{(\hat{c})}(t) =$  The minimum weight spanning tree rooted at the source  $s^{(c)}$ , in the weighted-graph  $\mathcal{G}$ .
- **MULTICAST TRAFFIC:**  $T_{opt}^{(c)}(t) =$  The minimum weight Steiner tree rooted at the source  $s^{(c)}$  and spanning the destinations  $\mathcal{D}^{(c)} = \{t_1^{(c)}, t_2^{(c)}, \dots, t_k^{(c)}\}$ , in the weighted-graph  $\mathcal{G}$ .
- ANYCAST TRAFFIC:  $T_{opt}^{(c)}(t) =$  The shortest of the k shortest  $s^{(c)} t_i^{(c)}$  paths, i = 1, 2, ..., k in the weighted-graph  $\mathcal{G}$ .

Thus, the routes are selected according to a *dynamic source routing* policy [25]. Apart from the minimum weight Steiner tree problem for the multicast traffic (which is NP-hard with several known efficient approximation algorithms [26]), all of the above routing problems on the *weighted* virtual graph may be solved efficiently using standard algorithms [27].

Dynamic Link Activation Policy: A feasible link activation schedule  $\mu^*(t) \in \mathcal{M}$  is dynamically chosen at each slot by maximizing the last term in the upper-bound of the drift-expression (7), given as follows:

$$\boldsymbol{\mu}^{*}(t) \in \operatorname*{arg\,max}_{\boldsymbol{\mu} \in \mathcal{M}} \left( \sum_{e \in E} \tilde{Q}_{e}(t) \mu_{e} \right)$$
(10)

This is the well-known max-weight scheduling policy, which can be solved efficiently under various interference models (e.g., *Primary* or node-exclusive model [28]).

In solving the above routing and scheduling problems, we tacitly made the assumption that the virtual queue vector  $\tilde{Q}(t)$  is globally known in each slot. We will discuss practical distributed implementation of our algorithm in section VII. Next, we establish stability of the virtual queues under the above policy, which will be instrumental for proving

throughput-optimality of the overall UMW policy:

Theorem 2: Under the above dynamic routing and link scheduling policy, the virtual queue process  $\{\tilde{Q}(t)\}_{t\geq 0}$  is strongly stable for any arrival rate  $\lambda \in int(\overline{\Lambda})$ , i.e.,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{e \in E} \mathbb{E}(\tilde{Q}_e(t)) < \infty.$$

*Proof:* Consider an arrival rate vector  $\lambda \in int(\overline{\Lambda})$ . Thus, from Eqns. (3) and (4), it follows that there exists a scalar  $\epsilon > 0$  and a vector  $\mu \in conv(\mathcal{M})$ , such that we can decompose the total arrival for each class  $c \in C$  into a finite number of routes, such that

$$\lambda_e \stackrel{\text{(def.)}}{=} \sum_{\substack{(i,c):e \in T_i^{(c)}, T_i^{(c)} \in \mathcal{T}^{(c)}}} \lambda_i^{(c)} \le \mu_e - \epsilon, \quad \forall e \in E$$
(11)

We can express  $\mu$  as,

$$\boldsymbol{\mu} = \sum_{i} p_i \boldsymbol{s}_i, \tag{12}$$

for some activation vectors  $s_i \in \mathcal{M}, \forall i$  and some probability distribution p.

Now consider the following auxiliary stationary randomized routing and link activation policy **RAND**  $\in \Pi$  for the virtual queue system  $\{\tilde{Q}(t)\}$ , which will be useful in our proof. The randomized policy **RAND** randomly selects the activation vector  $s_j$  with probability  $p_j$  and routes the incoming packet of class c along the route  $T_i^{(c)} \in \mathcal{T}^{(c)}$ , with probability  $\frac{\lambda_i^{(c)}}{\lambda^{(c)}}, \forall i, c$ . Hence, the total expected arrival rate to the virtual queue  $\tilde{Q}_e$ at time slot t, due to the action of the stationary randomized policy **RAND** is given by

$$\mathbb{E}A_e^{\mathbf{RAND}}(t) = \lambda_e = \sum_{(i,c):e \in T_i^{(c)}, T_i^{(c)} \in \mathcal{T}^{(c)}} \lambda_i^{(c)}, \quad \forall e \in E$$
(13)

and the expected total service rate to the virtual server for the queue  $\tilde{Q}_e$  is given by

$$\mathbb{E}\mu_e^{\mathbf{RAND}}(t) = \sum_i p_i \boldsymbol{s}_i(e) = \mu_e \tag{14}$$

Since our Max-Weight policy, **UMW**, minimizes the RHS of the drift expression in Eqn. (7) from the set of all feasible policies  $\Pi$ , we can write

$$\Delta^{\mathbf{UMW}}(t) \leq B + 2 \sum_{e \in E} \tilde{Q}_e(t) \mathbb{E} \left( A_e^{\mathbf{RAND}}(t) | \tilde{\boldsymbol{Q}}(t) \right) \\ - 2 \sum_{e \in E} \tilde{Q}_e(t) \mathbb{E} \left( \mu_e^{\mathbf{RAND}}(t) | \tilde{\boldsymbol{Q}}(t) \right) \\ \stackrel{(a)}{=} B + 2 \sum_{e \in E} \tilde{Q}_e(t) \left( \mathbb{E} A_e^{\mathbf{RAND}}(t) - \mathbb{E} \mu_e^{\mathbf{RAND}}(t) \right) \\ \stackrel{(b)}{=} B + 2 \sum_{e \in E} \tilde{Q}_e(t) \left( \lambda_e - \mu_e \right) \\ \stackrel{(c)}{\leq} B - 2\epsilon \sum_{e \in E} \tilde{Q}_e(t), \tag{15}$$

where (a) follows from the fact that the randomized policy **RAND** is memoryless and hence, independent of the virtual queues  $\tilde{Q}(t)$ , (b) follows from Eqns. (13) and (14) and finally (c) follows from Eqn. (11).

Taking expectation of both sides w.r.t. the virtual queue lengths  $\tilde{Q}(t)$ , we bound the expected drift at slot t as

$$\mathbb{E}L\big(\tilde{\boldsymbol{Q}}(t+1)\big) - \mathbb{E}L\big(\tilde{\boldsymbol{Q}}(t)\big) \le B - 2\epsilon \sum_{e \in E} \mathbb{E}(\tilde{Q}_e(t)) \quad (16)$$

Summing Eqn. (16) from t = 0 to T - 1 and remembering that  $L(\tilde{Q}(T)) \ge 0$  and  $L(\tilde{Q}(0)) = 0$ , we conclude that

$$\frac{1}{T} \sum_{t=0}^{T-1} \sum_{e \in E} \mathbb{E}(\tilde{Q}_e(t)) \le \frac{B}{2\epsilon}$$
(17)

 $\square$ 

Taking lim sup of both sides proves the claim.

As a consequence of the strong stability of the virtual queues  $\{\tilde{Q}_e(t), e \in E\}$ , we have the following sample-path result, which will be the key to our subsequent analysis:

Lemma 1: Under the action of the above policy, we have for any  $\lambda \in int(\bar{\Lambda})$ :

$$\lim_{t \to \infty} \frac{Q_e(t)}{t} = 0, \quad \forall e \in E, \text{ w.p. 1.}$$

In other words, the virtual queues are rate-stable [29].

Proof: See Appendix C.

The sample path result of Lemma 1 may be interpreted as follows: For any given realization  $\omega$  of the underlying sample space  $\Omega$ , define the function

$$F(\omega, t) = \max_{e \in E} Q_e(\omega, t).$$

Note that, for any  $t \in \mathbb{Z}_+$ , due to the boundedness of arrivals per slot, the function  $F(\omega, t)$  is well-defined and finite. In view of this, Lemma (1) states that under the action of the UMW policy,  $F(\omega, t) = o(t)$  almost surely.<sup>6</sup> This result will be used in our sample pathwise stability analysis of the physical queueing process  $\{Q(t)\}_{t>0}$ .

## D. Consequence of the Stability of the Virtual Queues

It is apparent from the virtual queue evolution equation (6), that the stability of the virtual queues under the **UMW** policy implies that the arrival rate at each virtual queue is *at most* the service rate offered to it under the **UMW** routing and scheduling policy. In other words, *effective load* of each edge *e* in the virtual system is at most unity. This is a necessary condition for stability of the physical queues when the same routing and link activation policy is used for the multi-hop physical network. In the following, we make the notion of "effective load" mathematically precise.

*Skorokhod Mapping:* Iterating on the system equation (6), we obtain the following well-known discrete time Skorokhod-Map representation [30] of the virtual queue dynamics

$$\tilde{Q}_e(t) = \left(\sup_{1 \le \tau \le t} \left(A_e^{\pi}(t-\tau,t) - S_e^{\pi}(t-\tau,t)\right)\right)^+, \quad (18)$$

where  $A_e^{\pi}(t_1, t_2) \stackrel{\text{def}}{=} \sum_{\tau=t_1}^{t_2-1} A_e^{\pi}(\tau)$ , is the total number of arrivals to the virtual queue  $\tilde{Q}_e$  in the time interval  $[t_1, t_2)$  and  $S_e^{\pi}(t_1, t_2) \stackrel{\text{def}}{=} \sum_{\tau=t_1}^{t_2-1} \mu_e^{\pi}(\tau)$ , is the total amount of service allocated to the virtual queue  $\tilde{Q}_e$  in the interval  $[t_1, t_2)$ . For reference, we provide a proof of Eqn. (18) in Appendix B.

Combining Equation (18) with Lemma 1, we conclude that under the **UMW** policy, *almost surely* for any sample path  $\omega \in \Omega$ , for each edge  $e \in E$  and any  $t_0 < t$ , we have

$$A_e(\omega; t_0, t) \le S_e(\omega; t_0, t) + F(\omega, t), \tag{19}$$

where  $F(\omega, t) = o(t)$ .

Implications for the Physical Network: Note that, every packet arrival to a virtual queue  $\tilde{Q}_e$  at time t corresponds to a packet in the physical network, that will eventually cross the edge e. Hence, the loading condition (19) implies that under the **UMW** policy, the total number of packets injected during any time interval  $(t_0, t]$ , willing to cross the edge e, is less than the total amount of service allocated to the edge e in that time interval up to an additive term of o(t). Thus informally, the "effective load" of any edge  $e \in E$  is at most unity.

By utilizing the sample-path result in Eqn. (19), in the following section we show that there exists a simple packet scheduling scheme for the physical network, which guarantees the stability of the physical queues, and consequently, throughput-optimality.

# V. OPTIMAL CONTROL OF THE PHYSICAL NETWORK

With the help of the virtual queue structure as defined above, we next focus our attention on designing a throughputoptimal control policy for the multi-hop physical network. As discussed in Section II, a control policy for the physical network consists of three components, namely (1) Routing, (2) Link activations and (3) Packet scheduling. In the proposed

$${}^{6}g(t) = o(t)$$
 if  $\lim_{t \to \infty} \frac{g(t)}{t} = 0.$ 

UMW policy, the (1) Routing and (2) Link activations for the physical network is done exactly in the same way as in the virtual network, based on the current values of the virtual queue state variables  $\tilde{Q}(t)$ , described in Section IV-C. It should be noted that, in the particular case of wireless networks, it is possible that a particular edge with positive virtual queue length is scheduled for transmission at a slot, even though the edge does not have any packet to transmit in its physical queue. The surprising fact, that follows from Theorem 4 is that this kind of wasted transmissions are rare and it *does not affect the throughput*.

There exist many possibilities for the third component, namely the packet scheduler, which efficiently resolves contention when multiple packets attempt to cross an active edge *e* at the same time-slot *t*. Popular choices for the packet scheduler include FIFO, LIFO etc. In this paper, we focus on a particular scheduling policy which has its origin in the context of *adversarial queueing theory* [31]. In particular, we extend the *Nearest To Origin* (NTO) policy to the generalized network flow setting, where a packet may be duplicated. This policy was proposed in [32] in the context of wired networks for the unicast problem. We appropriately extend this policy for use in generalized flow problems, including multicast, broadcast, and anycast, even in wireless networks. Our proposed scheduling policy is called *Extended* NTO (**ENTO**) and is defined as follows:

Definition 2 (Extended NTO): If multiple packets attempt to cross an active edge e at the same time slot t, the Extended Nearest To Origin (ENTO) policy gives priority to the packet which has traversed the least number of hops along its path from its origin up to the edge e.

The Extended NTO policy may be easily implemented by maintaining a priority queue [27] for each edge. The initial priority of each incoming packet at the source is set to zero. Upon transmission by an edge, the priority of a transmitted packet is decreased by one. The transmitted packet is then copied into the next-hop priority queue(s) (if any) according to its assigned route. See Figure 3 for an illustration. The pseudo code for the full UMW algorithm is provided in Algorithm 1.

We next state the following theorem which proves the stability of the physical queues under the **ENTO** policy:

Theorem 3: Under the action of the UMW policy with **ENTO** packet scheduling, the physical queues are ratestable [29] for any arrival vector  $\lambda \in int(\overline{\Lambda})$ , i.e.,

$$\lim_{t \to \infty} \frac{\sum_{e \in E} Q_e(t)}{t} = 0, \quad \text{w.p. 1}$$

*Proof:* This theorem is proved by extending the argument of Gamarnik [32] and combining it with the sample path loading condition in Eqn. (19). See Appendix D for the detailed argument.  $\Box$ 



Fig. 3. A schematic diagram showing the scheduling policy ENTO in action. The packets  $p_1$  and  $p_2$  originate from the sources  $S_1$  and  $S_2$ . Part of their assigned routes is shown in blue and red respectively. The packets contend for crossing the active edge  $e_3$  at the same time slot. According to the ENTO policy, the packet  $p_2$  has higher priority (having crossed a single edge  $e_4$  from its source) than  $p_1$  (having crossed two edges  $e_1$  and  $e_2$  from its source) for crossing the edge  $e_5$ , this edge does not fall in the path connecting the source  $S_1$  to the edge  $e_3$  and hence does not enter into priority calculations.

Algorithm 1 Universal Max-Weight Algorithm (UMW) at Slot t for the Generalized Flow Problem in a Wireless Network

- **Require:** Graph  $\mathcal{G}(V, E)$ , Virtual queue lengths  $\{\tilde{Q}_e(t), e \in E\}$  at the slot t.
- 1: **[Edge-Weight Assignment]** Assign each edge of the graph  $e \in E$  a weight  $W_e(t)$  equal to  $\tilde{Q}_e(t)$ , i.e.

$$\boldsymbol{W}(t) \leftarrow \tilde{\boldsymbol{Q}}(t)$$

- 2: [Route Assignment] Compute a Minimum Weight Route  $T^{(c)}(t) \in \mathcal{T}^{(c)}(t)$  for a class c incoming packet in the weighted graph  $\mathcal{G}(V, E)$ , according to Eqn. (9).
- 3: [Link Activation] Choose the activation  $\mu(t)$  from the set of all feasible activations  $\mathcal{M}$ , which maximizes the total activated link-weights, i.e.

$$\boldsymbol{\mu}(t) \gets \arg\max_{\boldsymbol{s} \in \mathcal{M}} \boldsymbol{s} \cdot \boldsymbol{W}(t)$$

- 4: **[Packet Forwarding]** Forward physical packets from the physical queues over the activated links according to the **ENTO** scheduling policy.
- 5: [Virtual Queue Counter Update] Update the virtual queues assuming a precedence-relaxed system, *i.e.*,

$$\tilde{Q}_e(t+1) \leftarrow \left(\tilde{Q}_e(t) + A_e(t) - \mu_e(t)\right)^+, \ \forall e \in E$$

As a direct consequence of Theorem 3, we have the main result of this paper:

Theorem 4: The UMW policy is throughput-optimal.

*Proof:* For any class  $c \in C$ , the number of packets  $R^{(c)}(t)$ , received by all nodes  $i \in \mathcal{D}^{(c)}$  may be bounded

as follows:

$$A^{(c)}(0,t) - \sum_{e \in E} Q_e(t) \stackrel{(*)}{\leq} R^{(c)}(t) \leq A^{(c)}(0,t), \quad (20)$$

where the lower-bound (\*) follows from the simple observation that if a packet p of class c has not reached all destination nodes  $\mathcal{D}^{(c)}$ , then at least one copy of it must be present in some physical queue.

Dividing both sides of Eqn. (20) by t, taking limits and using SLLN and Theorem 3, we conclude that w.p. 1

$$\lim_{t \to \infty} \frac{R^{(c)}(t)}{t} = \lambda^{(c)}$$

Hence from the definition (1), we conclude that UMW is throughput-optimal.  $\hfill \Box$ 

# VI. EXTENSION OF THE UMW POLICY TO TIME-VARYING NETWORKS

The proposed UMW policy may be readily extended to time-varying networks. In particular, we consider a simple model of a time-varying network, introduced earlier in our paper [10], where each link can be in two states at any slot - ON and OFF. Hence, the state of the network at slot t may be denoted by the binary vector  $\boldsymbol{\sigma}(t) \in \{0, 1\}^m$ , where

$$\boldsymbol{\sigma}(e,t) = \begin{cases} 1, & \text{if } e \text{ is ON at slot } t \\ 0, & \text{otherwise.} \end{cases}$$

The link-state process  $\{\sigma(\tau)\}_{\tau>1}$  is assumed to be evolving according to a stationary ergodic process. The routing and scheduling action of a feasible policy at slot t may depend on the observed network states  $\{\sigma(\tau)\}_0^t$ . A link e can be activated at slot t only if it is ON at that slot, i.e.,  $\sigma(e, t) = 1$ . Carrying out similar virtual queue construction and analysis, the one step conditional drift expression in Eqn. 7 is modified as follows

$$\Delta^{\pi}(t) \stackrel{\text{def}}{=} \mathbb{E} \left( L(\tilde{\boldsymbol{Q}}(t+1)) - L(\tilde{\boldsymbol{Q}}(t)) | \tilde{\boldsymbol{Q}}(t), \boldsymbol{\sigma}(t) \right)$$

$$\leq B + 2 \sum_{e \in E} \tilde{Q}_e(t) \mathbb{E} \left( A_e^{\pi}(t) | \tilde{\boldsymbol{Q}}(t), \boldsymbol{\sigma}(t) \right)$$

$$- 2 \sum_{e \in E} \tilde{Q}_e(t) \mathbb{E} \left( \mu_e^{\pi}(t) | \tilde{\boldsymbol{Q}}(t), \boldsymbol{\sigma}(t) \right)$$
(21)

Note that, the routing policy, which minimizes the RHS of Eqn. (21) is the same as Eqn. (9). Moreover, given the virtual queue lengths, it *does not depend* on the current network state  $\sigma(t)$ . The link scheduling policy, however, depends on the current network state. In particular, we activate a set of ON links which are feasible and have the maximum total weight, i.e.,

$$\boldsymbol{\mu}^{*}(t) \in \operatorname*{arg\,max}_{\boldsymbol{\mu} \in \mathcal{M}} \left( \sum_{e \in E} \tilde{Q}_{e}(t) \boldsymbol{\sigma}(e, t) \mu_{e} \right)$$
(22)

Using similar analysis as in [10] and following the analysis of the UMW policy, it can be shown Theorem 2 and Theorem 4 holds and hence, the policy is throughput-optimal. A simulation result for the performance of the UMW policy in time-varying network for broadcast traffic has been provided in Section VIII C.

## VII. DISTRIBUTED IMPLEMENTATION

The UMW policy in its original form, as given in Algorithm 1, is centralized in nature. This is because the sources need to know the topology of the network and the current value of the virtual queues  $\tilde{Q}(t)$  to solve the shortest route and the Max-Weight problems at steps (2) and (3) of the algorithm. With the advent of Software Defined Networking (SDN) technology, where the logically centralized control plane is separated from the forwarding elements [33], [34], the optimal centralized UMW policy may indeed be favorable in some practical setting. In fact, centralization of control plane functionalities are explicitly favored over the distributed schemes for enhanced network performance [35].

Nevertheless, it is also possible to devise a distributed version of the proposed policy. Although the topology of the network may be obtained efficiently by topology discovery algorithms [36], keeping track of the virtual queue evolution (Eqn. (6)) is subtle. Note that, in the special case where all packets arrive only at a single source node, no information exchange is necessary and the virtual queue updates (Step 5) may be implemented at the source locally. In the general case with multiple sources, it is necessary to periodically exchange packet arrival information among the sources to implement Step 5 exactly. To circumvent this issue, we propose the following class of heuristic **UMW** policies:

Heuristic UMW: Assign the edge weights to be the Physical queue lengths Q(t), instead of the virtual queue lengths  $\tilde{Q}(t)$ , in *either* step (2) or step (3) or *both* in the original UMW Algorithm 1.

Routing based on physical queue lengths still requires the exchange of queue length information. However, this can be efficiently implemented using the standard distributed Bellman-Ford (for shortest path routing), or the distributed MST algorithm by Gallager *et al.* [37]. The simulation results in section VIII-B show that the heuristic policy works well in practice and its delay performance is substantially better than the virtual queue based optimal **UMW** policy in wireless networks. The affirmative empirical results from the simulation section immediately prompt us to make the following conjecture:

Conjecture 1: The Heuristic UMW policy is throughput-optimal.

## VIII. NUMERICAL SIMULATION

# A. Delay Improvement Compared to the Back Pressure Policy - The Unicast Setting

To empirically demonstrate the superior performance of the UMW policy over the Back-Pressure class of policies in the unicast setting, we consider the wired network shown in Figure 4 and implement the following policies: (1) UMW (opt), (2) UMW (heuristic), (3) Backpressure (original [14]), and (4) Shortest-Path based Backpressure [15].



Fig. 4. The wired network topology used for unicast simulation.



Fig. 5. Comparison of time-averaged queue-lengths under the **BP** (original and shortest-path based [15]) and UMW (optimal and heuristic) policies in the unicast setting of Fig. 4. In terms of performance, we have UMW (opt.)> UMW (heu.) > BP (SP-based [15]) > BP (original).

All links are assumed to have a unit capacity. We consider two concurrent unicast sessions with source-destination pairs given by  $(s_1 = 1, t_1 = 8)$  and  $(s_2 = 5, t_2 = 2)$  respectively. It is easy to see that Max-Flow $(s_1 \rightarrow t_1) = 2$  and Max-Flow $(s_2 \rightarrow t_2) = 1$  and there exist mutually disjoint paths to achieve the optimal rate-pair  $(\lambda_1, \lambda_2) = (2, 1)$ . Assuming Poisson arrivals at the sources  $s_1$  and  $s_2$  with intensities  $\lambda_1 = 2\rho$  and  $\lambda_2 = \rho$ ,  $0 \le \rho \le 1$ , where  $\rho$  denotes the "load factor", Figure 5 shows a plot of the total average queue lengths as a function of the load factor  $\rho$  under the operation of the four policies considered above.

From the plot, we conclude that both the optimal and heuristic UMW policies *uniformly* outperform the **BP** (original) and SP-based BP policy in terms of average queue lengths, and hence (by Little's Law), end-to-end delay. The primary reason being, the **BP** class of policies, in principle, explores all possible paths to route packets to their destinations. The UMW policy, on the other hand, transmits all packets along "optimal" acyclic routes. This results in substantial reduction in latency.

# B. Using the Heuristic UMW Policy for Improved Latency in the Wireless Networks - The Broadcast Setting

Next, we empirically demonstrate that the heuristic UMW policy that uses physical queue lengths Q(t) (instead of virtual queues  $\tilde{Q}(t)$  as in the optimal UMW policy) not only achieves the full broadcast capacity but yields better delay performance in this particular wireless network. As discussed earlier, the heuristic policy is practically easier to implement in a distributed fashion. We simulate a  $3 \times 3$  wireless grid network shown in Figure 6, with *primary* interference constraints [20].



Fig. 6. The wireless topology used for broadcast simulation.



Fig. 7. Comparison of the Avg. Queue lengths as a function of the arrival rate for the optimal (in blue) and the heuristic (in red) UMW Policy for the grid network in Figure 6 in the **broadcast** setting.

The broadcast capacity of the network is known to be  $\lambda^* = \frac{2}{5}$  [8]. The ENTO policy is used for packet scheduling. The average queue length is plotted in Figure 7 as a function of the packet arrival rate  $\lambda$  under the operation of the (a) UMW (optimal) and (b) UMW (heuristic) policies. The plot shows that the heuristic policy results in much smaller queue lengths than the optimal policy. The reason being that physical queues capture the network congestion "more accurately" for proper link activations.

# C. Performance of the Optimal and Heuristic UMW Policy in Time-Varying Networks - The Broadcast Setting

In this section, we take a closer look at the wireless grid network in Figure 6 by numerically evaluating the broadcasting performance of the proposed policies, when the network is time-varying. In particular, we assume that at each slot a link is ON with probability  $p_{ON}$ , and is OFF w.p.  $1 - p_{ON}$ , independent of everything else. Packets can be transmitted only over the ON links at a given slot. Using similar analysis that we did for the static network, it can be easily shown that the proposed UMW policy remains throughput-optimal when the Max-Weight link activation at each slot is done with respect to the ON links at that slot. The packet routing policy remains the same as in the original UMW policy. The performance of the optimal and heuristic UMW policy is shown in Figure 8 for two different values of the parameter  $p_{ON}$ . It can be seen from the plot that the heuristic policy incurs substantially smaller queue lengths, compared to the optimal policy, especially in the low-load regime. Also, from the nearly identical vertical asymptotes in the queue length vs arrival rate



Fig. 8. Comparison of the time-averaged total queue lengths under the optimal (solid line) and heuristic (dashed line) UMW policy in the time-varying grid network (with parameter  $p_{ON}$ ), for the broadcast problem.

plots, we conclude that the heuristic policy is also throughputoptimal in this case.

# IX. CONCLUSION

In this paper, we have proposed a new, efficient and throughput-optimal policy, named Universal Max-Weight (UMW), for the Generalized Network Flow problem. The UMW policy can simultaneously handle a mix of Unicast, Broadcast, Multicast, and Anycast traffic in arbitrary networks and is empirically shown to have superior performance compared to the existing policies. The next step would be to investigate whether the UMW policy still retains its optimality when implemented with physical queue lengths, instead of the virtual queue lengths. An affirmative answer to this question would imply a more efficient implementation of the policy.

## APPENDIX

## A. Proof of Converse of Theorem 1

**Proof:** Consider any admissible arrival rate vector  $\lambda \in \Lambda(\mathcal{G}, \mathcal{C})$ . By definition, there exists an admissible policy  $\pi \in \Pi$  which supports the arrival vector  $\lambda$  in the sense of Eqn. (1). Without any loss of generality, we may assume the policy  $\pi$  to be stationary and the associated DTMC to be ergodic. Let  $A_i^{(c)}(t)$  denote the total number of packets from class c that have finished their routing along the route  $T_i^{(c)} \in \mathcal{T}^{(c)}$  up to time t. Note that, each packet is routed along one admissible route only. Hence, if the total number of arrival to the source  $s^{(c)}$  of class c up to time t is denoted by the random variable  $A^{(c)}(t)$ , we have

$$A^{(c)}(t) \stackrel{(a)}{\geq} \sum_{T_i^{(c)} \in \mathcal{T}^{(c)}} A_i^{(c)}(t) \stackrel{(b)}{=} R^{(c)}(t).$$
(23)

In the above, the inequality (a) follows from the observation that any packet p which has finished its routing along some route  $T_i^{(c)} \in \mathcal{T}^{(c)}$  by the time t, must have arrived at the source by the time t. The equality (b) follows from the observation that any packet p which has finished its routing by time t along some route  $T_i^{(c)} \in \mathcal{T}^{(c)}$ , has reached all of the destination nodes  $\mathcal{D}^{(c)}$  of class c by time t and vice versa.

Dividing both sides of equation (23) by t and taking limit as  $t \to \infty$ , we have w.p. 1

$$\lambda^{(c)} \stackrel{(d)}{=} \lim_{t \to \infty} \frac{A^{(c)}(t)}{t} \ge \liminf_{t \to \infty} \frac{1}{t} \sum_{\substack{T_i^{(c)} \in \mathcal{T}^{(c)}}} A_i^{(c)}(t)$$
$$= \liminf_{t \to \infty} \frac{R^{(c)}(t)}{t}$$
$$\stackrel{(f)}{=} \lambda^{(c)},$$

where equality (d) follows from the SLLN, and equality (f) follows from the Definition (1).

From the above inequalities, we conclude that w.p. 1

$$\lim_{t \to \infty} \frac{1}{t} \sum_{T_i^{(c)} \in \mathcal{T}^{(c)}} A_i^{(c)}(t) = \lambda^{(c)}, \quad \forall c \in \mathcal{C}$$
(24)

Now we use the fact that the policy  $\pi$  is stationary and the associated DTMC is *ergodic*. Thus the time-average limits exist and they are constant *a.s.*. For all  $T_i^{(c)} \in \mathcal{T}^c$  and  $c \in \mathcal{C}$ , define

$$\lambda_i^{(c)} \stackrel{\text{def}}{=} \lim_{t \to \infty} \frac{1}{t} A_i^{(c)}(t) \tag{25}$$

Hence, from the above, we get

$$\lambda^{(c)} = \sum_{T_i^{(c)} \in \mathcal{T}^{(c)}} \lambda_i^{(c)}$$
(26)

Now consider any edge  $e \in E$  in the graph  $\mathcal{G}$ . Since the variable  $A_i^{(c)}(t)$  denotes the total number of packets from class c, that have *completely* traversed along the tree  $T_i^{(c)}$ , the following inequality holds good for any time t

$$\sum_{(i,c):e \in T_i^{(c)}, T_i^{(c)} \in \mathcal{T}^{(c)}} A_i^{(c)}(t) \le \sum_{\tau=1}^t \mu_e(\tau),$$
(27)

where the left-hand side denotes a lower-bound on the number of packets that have crossed the edge e and the right-hand side denotes the amount of service that has been provided to the edge e up to time t by the policy  $\pi$ .

Dividing both sides by t and taking limits of both side, and noting that the limit on the left-hand side exists w.p. 1, we have

$$\sum_{(i,c):e\in T_i^{(c)}, T_i^{(c)}\in\mathcal{T}^c} \lambda_i^{(c)} \le \overline{\mu}_e,$$
(28)

where  $\overline{\mu} = \lim_{t\to\infty} \frac{1}{t} \sum_{\tau=1}^{t} \mu(\tau)$ . Since  $\mu(\tau) \in \mathcal{M}, \forall \tau$  and the set  $\operatorname{conv}(\mathcal{M})$  is closed, we conclude that  $\overline{\mu} \in \operatorname{conv}(\mathcal{M})$ . Eqns. (26) and (28) concludes the proof of the theorem.  $\Box$ 

## B. Proof of the Skorokhod Map Representation in Eqn. (18)

*Proof:* From the dynamics of the virtual queues in Eqn. (6), we have for any  $t \ge 1$ 

$$\tilde{Q}_e(t) \ge \tilde{Q}_e(t-1) + A_e(t-1) - \mu_e(t-1).$$
 (29)

Iterating (29)  $\tau$  times  $1 \le \tau \le t$ , we obtain

$$\ddot{Q}_e(t) \ge \dot{Q}_e(t-\tau) + A_e(t-\tau,t) - S_e(t-\tau,t),$$

where  $A_e(t_1, t_2) = \sum_{\tau=t_1}^{t_2-1} A_e(\tau)$  and  $S_e(t_1, t_2) = \sum_{\tau=t_1}^{t_2-1} \mu_e(\tau)$ , as defined before. Since each of the virtual-queue components are non-negative at all times (viz. (6)), we have  $\tilde{Q}_e(t-\tau) \geq 0$ . Thus,

$$\tilde{Q}_e(t) \ge A_e(t-\tau, t) - S_e(t-\tau, t).$$

Since the above holds for any time  $1 \le \tau \le t$  and the queues are always non-negative, we obtain

$$\tilde{Q}_e(t) \ge \left(\sup_{1 \le \tau \le t} \left(A_e(t-\tau,t) - S_e(t-\tau,t)\right)\right)_+ \quad (30)$$

To show that Eqn. (30) holds with equality, we consider two cases.

Case I ( $\hat{Q}_e(t) = 0$ ): Since the RHS of Eqn. (30) is non-negative, we immediately obtain equality throughout in Eqn (30).

Case II ( $\hat{Q}_e(t) > 0$ ): Consider the latest time  $t - \tau', 1 \le \tau' \le t$ , prior to t, at which  $\tilde{Q}_e(t - \tau') = 0$ . Such a time  $t - \tau'$  exists because we assumed the system to start with empty queues at time t = 0. Hence  $Q_e(z) > 0$  throughout the time interval  $z \in [t - \tau' + 1, t]$ . As a result, in this time interval the system dynamics for the virtual-queues (6) takes the following form

$$\tilde{Q}_e(z) = \tilde{Q}_e(z-1) + A_e(z-1) - \mu_e(z-1),$$

Iterating the above recursion in the interval  $z \in [t - \tau' + 1, t]$ , we obtain

$$\tilde{Q}_e(t) = A_e(t - \tau', t) - S_e(t - \tau', t)$$
(31)

We conclude the proof upon combining Eqns. (30) and (31).  $\Box$ 

## C. Proof of Lemma 1

*Proof:* We will establish this result by appealing to the *Strong Stability Theorem* (Theorem 2.8) of [29]. For this, we first consider an associated system  $\{\hat{Q}(t)\}_{t\geq 0}$  with a slightly different queueing recursion, as considered in [29] (Eqn. 2.1, pp-15). For a given sequence  $\{A(t), \mu(t)\}_{t\geq 0}$ , define the following recursion for all  $e \in E$ ,

$$\hat{Q}_e(t+1) = (\hat{Q}_e(t) - \mu_e(t))_+ + A_e(t), 
\hat{Q}_e(0) = 0.$$
(32)

Recall the dynamics of the virtual queues (Eqn. (6)):

$$Q_e(t+1) = (Q_e(t) + A_e(t) - \mu_e(t))^+,$$
  

$$\tilde{Q}_e(0) = 0.$$
(33)

The following proposition is easy to establish. *Proposition 5: For all*  $e \in E$ 

$$A_{\max} + \tilde{Q}_e(t) \stackrel{(*)}{\geq} \hat{Q}_e(t) \stackrel{(**)}{\geq} \tilde{Q}_e(t), \quad \forall t \ge 0.$$
(34)  
sing expectation throughout the first inequality (\*) of

Taking expectation throughout the first inequality (\*) of Eqn. (34) for any  $e \in E$ , we have for each  $t \ge 0$ 

$$\mathbb{E}(Q_e(t)) \le \mathbb{E}(Q_e(t)) + A_{\max}$$

Thus,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}(\hat{Q}_e(t)) \leq \limsup_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}(\tilde{Q}_e(t)) + A_{\max}$$

$$\stackrel{(a)}{<} \infty,$$

where (a) follows from the strong stability of the virtual queues under **UMW**. This shows that, the associated queue process  $\{\hat{Q}(t)\}_{t>0}$  is also strongly stable under **UMW**.

Since the total external arrival  $A(t) = \sum_{e} A_e(t)$  at slot t is assumed to be bounded w.p. 1, applying Theorem 2.8, part (b) of [29], we conclude that for any  $e \in E$ 

$$\lim_{t \to \infty} \frac{\hat{Q}_e(t)}{t} = 0, \quad \text{w.p. 1}$$

Using the second inequality (\*\*) of Proposition 5 and the non-negativity of the virtual queues, we conclude that for any  $e\in E$ 

$$\lim_{t \to \infty} \frac{\hat{Q}_e(t)}{t} = 0, \quad \text{w.p. 1}$$

Finally, using the union bound, we conclude that

$$\lim_{t \to \infty} \frac{\tilde{Q}_e(t)}{t} = 0, \quad \forall e \in E \text{ w.p. 1.}$$

#### D. Proof of Theorem 3

Throughout this proof, we will fix a sample point  $\omega \in \Omega$ , giving rise to a sample path satisfying the condition (19). All random processes<sup>7</sup> will be evaluated at this sample path. For the sake of notational simplicity, we will drop the argument  $\omega$  for evaluating any random variable X at the sample point  $\omega$ , e.g., the deterministic sample-path  $X(\omega, t)$  will be simply denoted by X(t). We now establish a simple analytical result which will be useful in the main proof of the theorem:

Lemma 2: Consider a non-negative function  $\{F(t), t \geq 1\}$  defined on the set of natural numbers, such that F(t) = o(t). Define  $M(t) = \sup_{0 \leq \tau \leq t} F(\tau)$ . Then 1. M(t) is non-decreasing in t. 2. M(t) = o(t)

*Proof:* That M(t) is non-decreasing follows directly from the definition of  $M(t) = \sup_{0 \le \tau \le t} F(t)$ . We now prove the claim (2).

Case I (The Function F(t) Is Bounded): In this case, the function M(t) is also bounded and the claim follows immediately.

*Case II (The Function* F(t) *Is Unbounded):* Define the subsequence  $\{r_k\}_{k\geq 1}$ , corresponding to the time of maximums

of the function M(t) up to time t. Formally the sequence  $\{r_k\}_{k>1}$  is defined recursively as follows,

$$r_1 = 1 \tag{35}$$

$$r_k = \{\min t > r_{k-1} : F(t) > \max_{\tau < t-1} F(\tau)\}$$
(36)

Since the function F(t) is assumed to be unbounded, we have  $r_k \to \infty$  as  $k \to \infty$ . In the literature [39], the sequence  $\{r_k\}$  is also known as the sequence of *records* of the function F(t). With this definition, for any  $t \ge 1$  and for  $r_k \le t$  corresponding to the latest record up to time t, we readily have  $M(t) = F(r_k)$ . Hence,

$$\frac{M(t)}{t} = \frac{F(r_k)}{t} \stackrel{(a)}{\leq} \frac{F(r_k)}{r_k},\tag{37}$$

where Eqn. (a) follows from the fact that  $r_k \leq t$ . Thus for any sequence of natural numbers  $\{t_i\}_1^\infty$ , we have a corresponding sequence  $\{r_{k_i}\}_{i=1}^\infty$  such that for each *i*, we have

$$\frac{M(t_i)}{t_i} = \frac{F(r_{k_i})}{t} \stackrel{(a)}{\leq} \frac{F(r_{k_i})}{r_{k_i}}$$

This implies,

 $\square$ 

$$\limsup_{t \to \infty} \frac{M(t)}{t} \le \limsup_{t \to \infty} \frac{F(t)}{t} \stackrel{(b)}{=} 0, \tag{38}$$

where Eqn (b) follows from our hypothesis on the function F(t). Also since  $M(t) \ge F(t)$ , from Eqn. (38) we conclude that

$$\lim_{t \to \infty} \frac{M(t)}{t} = 0 \tag{39}$$

As a direct consequence of Lemma 2 and the property of the sample-point  $\omega$  under consideration, we have:

$$A_e(t_0, t) \le S_e(t_0, t) + M(t), \quad \forall e \in E, \quad \forall t_0 \le t, \qquad (40)$$

for some non-decreasing non-negative function M(t) = o(t). Equipped with Eqn. (40), we return to the proof of the Theorem 3.

Proposition 6: ENTO is rate-stable.

*Proof:* We generalize the argument by Gamarnik [32] to prove the proposition. Recall that we are analyzing the time-evolution of a fixed sample point  $\omega \in \Omega$ , satisfying Eqn. (40).

Let  $R_e(0)$  denote the total number of packets waiting to cross the edge e at time t = 0. Also, let  $R_k(t)$  denote the total number of packets at time t, which are *exactly* k hops away from their respective sources. Such packets will be called "layer k" packets in the sequel. If a packet is duplicated along its assigned route T (which is, in general, a tree), each copy of the packet is counted separately in the variable  $R_k(t)$ , *i.e.*,

$$R_{k}(t) = \sum_{T \in \mathcal{T}} R_{(e_{k}^{T}, T)}(t), \qquad (41)$$

<sup>&</sup>lt;sup>7</sup>Recall that, a discrete-time integer-valued random process  $X(\omega; t)$  is a measurable map from the sample space  $\Omega$  to the set of all integersequences  $\mathbb{Z}^{\infty}$  [38], i.e.,  $X: \Omega \to \mathbb{Z}^{\infty}$ .

where the variable  $R_{(e,T)}(t)$  denotes the number of packets following the routing tree T, that are waiting to cross the edge  $e \in T$  at time t. The edge  $e_k^T$  is an edge located  $k^{\text{th}}$  hop away from the source in the tree T. If there are more than one such edge (because the tree T has more than one branch), we include all these edges in the summation (41). We show by induction that  $R_k(t)$  is *almost surely* bounded by a function, which is o(t).

Base Step k = 0: Fix an edge e and time t. Let  $t_0 \leq t$  be the largest time at which no packets of layer 0 (packets which have not crossed any edge yet) were waiting to cross e. If no such time exists, set  $t_0 = 0$ . Hence, the total number of layer 0 packets waiting to cross the edge e at time  $t_0$  is at most  $Q_e(0)$ . During the time interval  $[t_0, t]$ , as a consequence of the **UMW** control policy (40), at most  $S_e(t_0, t) + M(t)$  external packets, that want to cross the edge e in future, have been admitted to the network. Also, by the choice of the time  $t_0$ , the edge ewas always having packets to transmit during the entire time interval  $[t_0, t]$ . Since **ENTO** scheduling policy is followed, layer 0 packets have priority over all other packets. Hence, it follows that the total number of packets at the edge e at time t satisfies

$$\sum_{T:e \in e_0^T} R_{(e,T)}(t) \le R_e(0) + S_e(t_0,t) + M(t) - S_e(t_0,t)$$
$$\le R_e(0) + M(t)$$
(42)

As a result, we have 
$$R_0(t) \leq \sum_e R_e(0) + |E|M(t)$$
, for all   
t. Let  $B_0(t) \stackrel{\text{def}}{=} \sum_e R_e(0) + |E|M(t)$ . Since  $M(t) = o(t)$ , we have  $B_0(t) = o(t)$ . Note that, since  $M(t)$  is monotonically non-decreasing by definition, so is  $B_0(t)$ .

Induction Step: Suppose that, for some monotonically nondecreasing functions  $B_j(t) = o(t), j = 0, 1, 2, ..., k - 1$ , we have  $R_j(t) \leq B_j(t)$ , for all time t. We next show that  $R_k(t) \leq B_k(t)$  for all t, where  $B_k(t) = o(t)$ .

Again, fix an edge e and an arbitrary time t. Let  $t_0 \leq t$ denote the largest time before t, such that there was no layer kpacket waiting to cross the edge e. Set  $t_0 = 0$  if no such time exists. Hence, the edge e was always having packets to transmit during the time interval  $[t_0, t]$  (packets in layer k or lower). The layer k packets that wait to cross edge eat time t are composed only of a subset of packets which were in layers  $0 \leq j \leq k-1$  at time  $t_0$  or packets that arrived during the time interval  $[t_0, t]$  and have edge e as one of their  $k^{th}$  edge on the route followed. By our induction assumption, the first group of packets has a size bounded by  $\sum_{j=0}^{k-1} B_j(t_0) \leq \sum_{j=0}^{k-1} B_j(t)$ , where we have used the fact (from our previous induction step) that the functions  $B_j(\cdot)$ 's are monotonically non-decreasing. The size of the second group of packets is given by  $\sum_{T:e \in e_k^T} A_T(t_0, t)$ . We next estimate the number of layer k packets that crossed the edge eduring the time interval  $[t_0, t]$ . Since ENTO policy is used, layer k packets were not processed only when there were packets in layers up to k-1 that wanted by  $\sum_{j=0}^{k-1} B_j(t_0) \le \sum_{j=0}^{k-1} B_j(t)$ , which denotes the total possible number of packets in layers up to k-1 at time  $t_0$ , plus  $\sum_{j=0}^{k-1} \sum_{T:e \in e_j^T} A_T(t_0, t)$ , which is the number of new packets that arrived in the interval  $[t_0, t]$ 

and intend to cross the edge e within first k - 1 hops. Thus, we conclude that at least

$$\max\left\{0, S_e(t_0, t) - \sum_{j=0}^{k-1} B_j(t) - \sum_{j=0}^{k-1} \sum_{T: e \in e_j^T} A_T(t_0, t)\right\}$$
(43)

packets of layer k crossed e during the time interval  $[t_0, t]$ . Hence,

$$\sum_{T:e \in e_k^T} R_{(e,T)}(t)$$

$$\leq \sum_{j=0}^{k-1} B_j(t) + \sum_{T:e \in e_k^T} A_T(t_0, t)$$

$$- \left(S_e(t_0, t) - \sum_{j=0}^{k-1} B_j(t) - \sum_{j=0}^{k-1} \sum_{T:e \in e_j^T} A_T(t_0, t)\right)$$

$$= 2\sum_{j=0}^{k-1} B_j(t) + \sum_{j=0}^k \sum_{T:e \in e_j^T} A_T(t_0, t) - S_e(t_0, t)$$

$$\stackrel{(a)}{\leq} 2\sum_{j=0}^{k-1} B_j(t) + M(t),$$

where Eqn. (a) follows from the arrival bound (40). Hence, the total number of layer k packets at time t is bounded by

$$R_k(t) \le 2|E| \sum_{j=0}^{k-1} B_j(t) + M(t)|E|$$
(44)

Define  $B_k(t)$  to be the RHS of the above equation, i.e.

$$B_k(t) \stackrel{\text{(def)}}{=} 2|E| \sum_{j=0}^{k-1} B_j(t) + M(t)|E|$$
(45)

Using our induction assumption and Eqn. (45), we conclude that  $B_k(t) = o(t)$  and it is monotonically non-decreasing. This completes the induction step.

To conclude the proof of the proposition, notice that the total size of the physical queues at time t may be written as

$$\sum_{e \in E} Q_e(t) = \sum_{k=1}^{n-1} R_k(t)$$
(46)

Since the previous inductive argument shows that for all k, we have  $R_k(t) \leq B_k(t)$  where  $B_k(t) = o(t)$  a.s., we conclude

$$\lim_{t \to \infty} \frac{\sum_{e \in E} Q_e(t)}{t} = 0, \quad \text{w.p. 1}, \tag{47}$$

This implies that the physical queues are rate stable [29], jointly under the operation of UMW and ENTO.  $\Box$ 

#### References

- A. Sinha and E. Modiano, "Optimal control for generalized networkflow problems," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, May 2017, pp. 1–9. [Online]. Available: http://bit.ly/infocom17Sinha
- [2] R. Rustin, *Combinatorial Algorithms*. Toronto, ON, Canada: Algorithmics Press, 1973.
- [3] A. Czumaj and W. Rytter, "Broadcasting algorithms in radio networks with unknown topology," in *Proc. 44th Annu. IEEE Symp. Found. Comput. Sci.*, Oct. 2003, pp. 492–501.

- [4] J. Widmer, C. Fragouli, and J.-Y. Le Boudec, "Low-complexity energyefficient broadcasting in wireless ad-hoc networks using network coding," in *Proc. 1st Workshop Netw. Coding, Theory, Appl.*, 2005, pp. 1–6.
- [5] E. Kranakis, D. Krizanc, and A. Pelc, "Fault-tolerant broadcasting in radio networks," J. Algorithms, vol. 39, no. 1, pp. 47–67, 2001.
- [6] L. Massoulie, A. Twigg, C. Gkantsidis, and P. Rodriguez, "Randomized decentralized broadcasting algorithms," in *Proc. 26th IEEE Int. Conf. Comput. Commun. (INFOCOM)*, May 2007, pp. 1073–1081.
- [7] D. Towsley and A. Twigg, "Rate-optimal decentralized broadcasting: The wireless case," in *Proc. ACITA*, 2008, pp. 323–333.
  [8] A. Sinha, G. Paschos, and E. Modiano, "Throughput-optimal multi-
- [8] A. Sinha, G. Paschos, and E. Modiano, "Throughput-optimal multihop broadcast algorithms," in *Proc. 17th ACM Int. Symp. Mobile Ad Hoc Netw. Comput. (MobiHoc)*, New York, NY, USA, 2016, pp. 51–60. [Online]. Available: http://doi.acm.org/10.1145/2942358.2942390
- [9] A. Sinha, G. Paschos, C.-P. Li, and E. Modiano, "Throughput-optimal broadcast on directed acyclic graphs," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, Apr. 2015, pp. 1248–1256.
- [10] A. Sinha, L. Tassiulas, and E. Modiano, "Throughput-optimal broadcast in wireless networks with dynamic topology," in *Proc. 17th ACM Int. Symp. Mobile Ad Hoc Netw. Comput. (MobiHoc)*, New York, NY, USA, 2016, pp. 21–30. [Online]. Available: http://doi.acm.org/10. 1145/2942358.2942389
- [11] K. Jain, M. Mahdian, and M. R. Salavatipour, "Packing Steiner trees," in Proc. 14th Annu. ACM-SIAM Symp. Discrete Algorithms, 2003, pp. 266–274.
- [12] S. Sarkar and L. Tassiulas, "A framework for routing and congestion control for multicast information flows," *IEEE Trans. Inf. Theory*, vol. 48, no. 10, pp. 2690–2708, Oct. 2002.
- [13] L. Bui, R. Srikant, and A. Stolyar, "Optimal resource allocation for multicast sessions in multi-hop wireless networks," *Phil. Trans. Roy. Soc. London A, Math., Phys. Eng. Sci.*, vol. 366, no. 1872, pp. 2059–2074, 2008.
- [14] L. Tassiulas and A. Ephremides, "Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks," *IEEE Trans. Autom. Control*, vol. 37, no. 12, pp. 1936–1948, Dec. 1992.
- [15] L. Ying, S. Shakkottai, A. Reddy, and S. Liu, "On combining shortestpath and back-pressure routing over multihop wireless networks," *IEEE/ACM Trans. Netw.*, vol. 19, no. 3, pp. 841–854, Jun. 2011.
- [16] M. Zargham, A. Ribeiro, and A. Jadbabaie, "Accelerated backpressure algorithm," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2013, pp. 2269–2275.
- [17] A. Sinha, P. Mani, J. Liu, A. Flavel, and D. A. Maltz, "Distributed load management in anycast-based CDNs," in *Proc. 53rd Annu. Allerton Conf. Commun., Control, Comput. (Allerton)*, Sep./Oct. 2015, pp. 74–82.
- [18] L. Bui, R. Srikant, and A. Stolyar, "Novel architectures and algorithms for delay reduction in back-pressure scheduling and routing," in *Proc. IEEE INFOCOM*, Apr. 2009, pp. 2936–2940.
- [19] A. Sinha and E. Modiano, "Throughput-optimal broadcast in wireless networks with point-to-multipoint transmissions," in *Proc. 18th ACM Int. Symp. Mobile Ad Hoc Netw. Comput. (Mobihoc)*, New York, NY, USA, 2017, pp. 3:1–3:10. [Online]. Available: http://doi.acm.org/10. 1145/3084041.3084064
- [20] C. Joo, X. Lin, and N. B. Shroff, "Greedy maximal matching: Performance limits for arbitrary network graphs under the node-exclusive interference model," *IEEE Trans. Autom. Control*, vol. 54, no. 12, pp. 2734–2744, Dec. 2009.
- [21] D. B. West *et al.*, Introduction to Graph Theory, vol. 2. Upper Saddle River, NJ, USA: Prentice-Hall, 2001.
- [22] R. Ahlswede, N. Cai, S.-Y. Li, and R. W. Yeung, "Network information flow," *IEEE Trans. Inf. Theory*, vol. 46, no. 4, pp. 1204–1216, Jul. 2000.
- [23] M. J. Neely, "Energy optimal control for time-varying wireless networks," *IEEE Trans. Inf. Theory*, vol. 52, no. 7, pp. 2915–2934, Jul. 2006.
- [24] J. K. Lenstra and A. H. G. R. Kan, "Complexity of scheduling under precedence constraints," *Oper. Res.*, vol. 26, no. 1, pp. 22–35, 1978.
- [25] D. B. Johnson and D. A. Maltz, "Dynamic source routing in ad hoc wireless networks," in *Mobile Computing*. New York, NY, USA: Springer, 1996, pp. 153–181.
- [26] J. Byrka, F. Grandoni, T. Rothvoß, and L. Sanità, "An improved LPbased approximation for Steiner tree," in *Proc. 42nd ACM Symp. Theory Comput.*, 2010, pp. 583–592.
- [27] T. H. Cormen, Introduction to Algorithms. Cambridge, MA, USA: MIT Press, 2009.
- [28] L. X. Bui, S. Sanghavi, and R. Srikant, "Distributed link scheduling with constant overhead," *IEEE/ACM Trans. Netw.*, vol. 17, no. 5, pp. 1467–1480, Oct. 2009.

- [29] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems," *Synthesis Lectures Commun. Netw.*, vol. 3, no. 1, pp. 1–211, 2010.
- [30] S. Meyn, Control Techniques for Complex Networks. Cambridge, U.K.: Cambridge Univ. Press, 2008.
- [31] M. Andrews *et al.*, "Universal-stability results and performance bounds for greedy contention-resolution protocols," *J. ACM*, vol. 48, no. 1, pp. 39–69, 2001.
- [32] D. Gamarnik, "Stability of adaptive and non-adaptive packet routing policies in adversarial queueing networks," in *Proc. 31st Annu. ACM Symp. Theory Comput.*, 1999, pp. 206–214.
- [33] C. E. Rothenberg et al., "Revisiting routing control platforms with the eyes and muscles of software-defined networking," in Proc. 1st Workshop Hot Topics Softw. Defined Netw. (HotSDN), New York, NY, USA, 2012, pp. 13–18. [Online]. Available: http://doi.acm.org/10. 1145/2342441.2342445
- [34] N. Feamster, H. Balakrishnan, J. Rexford, A. Shaikh, and J. van der Merwe, "The case for separating routing from routers," in *Proc. ACM SIGCOMM Workshop Future Directions Netw. Archit. (FDNA)*, New York, NY, USA, 2004, pp. 5–12. [Online]. Available: http://doi.acm.org/10.1145/1016707.1016709
- [35] R. Jain and S. Paul, "Network virtualization and software defined networking for cloud computing: A survey," *IEEE Commun. Mag.*, vol. 51, no. 11, pp. 24–31, Nov. 2013.
- [36] R. Chandra, C. Fetzer, and K. Hogstedt, "A mesh-based robust topology discovery algorithm for hybrid wireless networks," in *Proc. AD-HOC Netw. Wireless*, 2002, pp. 1–25.
- [37] R. G. Gallager, P. A. Humblet, and P. M. Spira, "A distributed algorithm for minimum-weight spanning trees," ACM Trans. Program. Lang. Syst., vol. 5, no. 1, pp. 66–77, 1983.
- [38] R. Durrett, Probability: Theory and Examples. Cambridge, U.K.: Cambridge Univ. Press, 2010.
- [39] N. Glick, "Breaking records and breaking boards," Amer. Math. Monthly, vol. 85, no. 1, pp. 2–26, 1978.



Abhishek Sinha received the B.E. degree in electronics and telecommunication engineering from Jadavpur University, Kolkata, India, in 2010, the M.E. degree in telecommunication engineering from the Indian Institute of Science, Bangalore, India, in 2012, and the Ph.D. degree from the Massachusetts Institute of Technology in 2017, where he is involved in the Laboratory for Information and Decision Systems. He is currently a Senior Engineer with Qualcomm Research, San Diego, CA, USA. His research interests include network control,

information theory, optimization, and applied probability. He was a recipient of several awards, including the Best Paper Award in ACM MobiHoc 2016, the Prof. Jnansaran Chatterjee Memorial Gold Medal and the T.P. Saha Memorial Gold Centered Silver Medal from Jadavpur University, and Jagadis Bose National Science Talent Search Scholarship.



**Eytan Modiano** (F'12) received the B.S. degree in electrical engineering and computer science from the University of Connecticut, Storrs, CT, USA, in 1986, and the M.S. and Ph.D. degrees in electrical engineering from the University of Maryland, College Park, MD, USA, in 1989 and 1992, respectively. He was a Naval Research Laboratory Fellow from 1987 to 1992 and a National Research Council Post-Doctoral Fellow from 1992 to 1993. He was with the MIT Lincoln Laboratory from 1993 to 1999. Since 1999, he has been a Faculty Member with MIT,

where he is currently a Professor and an Associate Department Head with the Department of Aeronautics and Astronautics, and the Associate Director of the Laboratory for Information and Decision Systems.

His research interests include communication networks and protocols with emphasis on satellite, wireless, and optical networks. He is an Associate Fellow of the AIAA, and served on the IEEE Fellows Committee. He was a co-recipient of the MobiHoc 2016 Best Paper Award, the Wiopt 2013 Best Paper Award, and the Sigmetrics 2006 Best Paper Award. He was the Technical Program Co-Chair for the IEEE Wiopt 2006, the IEEE Infocom 2007, the ACM MobiHoc 2007, and the DRCN 2015. He is Editor-in-Chief of the IEEE/ACM TRANSACTIONS ON NETWORKING, and served as an Associate Editor of the IEEE TRANSACTIONS ON NETWORKING.