

A Cross-Layer Backpressure Architecture for Wireless Multihop Networks

Rafael Laufer, *Member, IEEE, ACM*, Theodoros Salonidis, *Member, IEEE, ACM*,
Henrik Lundgren, *Member, IEEE, ACM*, and Pascal Le Guyadec

Abstract—Contemporary wireless multihop networks operate much below their capacity due to the poor coordination among transmitting nodes. In this paper, we present XPRESS, a cross-layer backpressure architecture designed to reach the capacity of wireless multihop networks. Instead of a collection of poorly coordinated wireless routers, XPRESS turns a mesh network into a wireless switch. Transmissions over the network are scheduled using a throughput-optimal backpressure algorithm. Realizing this theoretical concept entails several challenges, which we identify and address with a cross-layer design and implementation on top of our wireless hardware platform. In contrast to previous work, we implement and evaluate backpressure scheduling over a TDMA MAC protocol, as it was originally proposed in theory. Our experiments in an indoor testbed show that XPRESS can yield up to 128% throughput gains over 802.11.

Index Terms—Backpressure scheduling and routing, design and implementation, wireless multihop networks.

I. INTRODUCTION

EXISTING networks are designed in layers, where protocols operate independently at each layer of the network stack. This approach provides flexibility with a modular design and standardization, but it may result in severe performance degradation when these protocols do not cooperate well. This is usually the case of wireless multihop networks, where noise and interference at lower layers affect the routing and congestion control performed at upper layers.

A common approach to address these performance issues is then to modify a single layer of the protocol stack while keeping other layers intact. Cross-layer architectures offer a radical alternative by advocating cooperation among the multiple layers of the protocol stack. At the core of these architectures is the backpressure scheduling algorithm [1], which, in theory, achieves the network capacity.

Manuscript received May 25, 2012; revised November 08, 2012; accepted February 17, 2013; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor L. Qiu. Date of publication April 01, 2013; date of current version April 14, 2014. This work was supported by the European Community Seventh Framework Programme under the OPNEX project (FP7-ICT-224218). A preliminary version of this paper was published in the Proceedings of the 17th ACM International Conference on Mobile Computing and Networking (MobiCom), Las Vegas, NV, USA, September 19–23, 2011.

R. Laufer is with Bell Laboratories, Alcatel-Lucent, Holmdel, NJ 07733 USA (e-mail: rafael.laufer@alcatel-lucent.com).

T. Salonidis is with the IBM T. J. Watson Research Center, Yorktown Heights, NY 10598 USA (e-mail: tsaloni@us.ibm.com).

H. Lundgren and P. Le Guyadec are with the Technicolor Paris Research and Innovation Center, 92443 Issy les Moulineaux Cedex, France (e-mail: henrik.lundgren@technicolor.com; pascal.leguyadec@technicolor.com).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNET.2013.2249592

Translating this theoretical concept into a fully practical system, however, entails several challenges, mainly due to its idealized assumptions. In essence, backpressure assumes a globally synchronized time-slotted MAC protocol as well as a centralized controller that both computes and disseminates a schedule (i.e., a set of links allowed to transmit) for each time slot. Moreover, the schedule computation requires the global knowledge of both per-flow queue backlogs and network state (i.e., link quality and link interference pattern), which therefore must be measured at the wireless nodes and provided to the controller in a timely manner. Recently, practical backpressure systems that relax some of these assumptions and approximate the backpressure algorithm on top of the IEEE 802.11 MAC protocol have been proposed [2]–[4]. Even though these approaches have shown performance benefits, the step toward a practical system implementing optimal backpressure scheduling has not yet been made.

In this paper, we present XPRESS, a throughput-optimal backpressure architecture for wireless multihop networks. In XPRESS, a mesh network is transformed into a wireless switch, where packet routing and scheduling decisions are made by a backpressure scheduler. XPRESS is composed of a central controller, which performs backpressure scheduling based on the measured wireless network state, and also of the wireless nodes, which periodically provide the network measurements and execute the computed schedule using a cross-layer protocol stack. The implementation of XPRESS on our 802.11 platform resulted in novel techniques that overcome the above-mentioned backpressure challenges and have a wider applicability to the design of centralized multihop wireless systems. Our contributions are as follows.

First, the XPRESS cross-layer stack gracefully integrates the transport, network, and MAC layers. In order to achieve synergy among these layers on our programmable 802.11 platform, we implement: 1) a congestion control scheme to ensure the network operates within the capacity region; 2) a coordination mechanism between network-layer flow queues and MAC-layer link queues, which enables per-link queue implementation on memory-constrained wireless interfaces; and 3) a multihop TDMA MAC protocol that not only ensures global synchronization among nodes, but also enables coordinated transmissions within slot boundaries according to the exact backpressure schedule.

Second, we find that a TDMA MAC on top of 802.11 PHY results in binary interference patterns. This relaxes the channel-state estimation requirement from exact values to binary state (i.e., delivery ratios are close to 0% or 100%), which allows us to efficiently find the link transmission sets and accurately

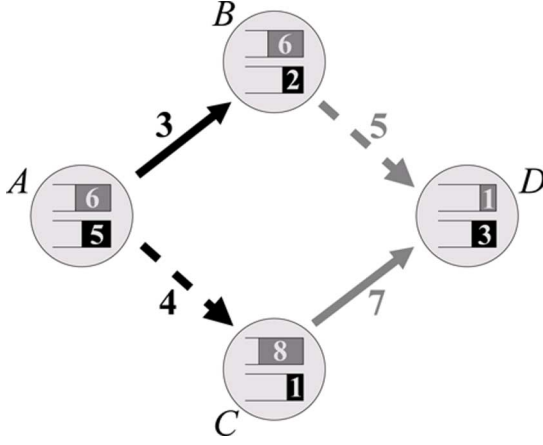


Fig. 1. Backpressure scheduling in a network with two flows, black and gray, from A to D . Links in sets $\{(A, B), (C, D)\}$ (continuous) and $\{(A, C), (B, D)\}$ (dashed) can be scheduled in the same slot.

estimate the queue backlogs of the network. Complementary techniques, such as multislot schedule computations and speculative scheduling, reduce the protocol overhead further and provide the scheduler at the controller with a longer computation time budget at the expense of outdated, but still accurate, network state.

Third, we propose an interference estimation mechanism with only $O(N)$ measurement complexity that allows the backpressure scheduler to determine at the TDMA frame timescale the links which can transmit without interference. The mechanism uses received signal strength (RSS) as well as an adaptive technique based on packet loss to cope with the RSS measurement limitations of 802.11.

Fourth, our evaluations in an indoor testbed show that XPRESS provides close to perfect fairness in small-scale centralized WLANs and 63%–128% throughput gains over 802.11 in multihop mesh configurations. We also show that XPRESS accurately emulates the optimal backpressure schedule and delivers relatively low delays when operating close to capacity.

Finally, we provide an analysis of the communication and computation overhead of XPRESS and identify different system design choices and limitations.

II. BACKPRESSURE SCHEDULING

The backpressure algorithm was introduced in [1] as a scheduling policy that maximizes the throughput of wireless multihop networks. Assuming slotted time, the basic idea of backpressure scheduling is to select the “best” set of noninterfering links for transmission at each slot. We now describe this idea in a 4-node network with two flows, black and gray, from node A to D , depicted in Fig. 1. Each node maintains a separate queue for each flow. For each queue, the number of backlogged packets is shown. Assume that we have two link sets, $\{(A, B), (C, D)\}$ and $\{(A, C), (B, D)\}$, shown as continuous and dashed lines, respectively. The links in each set do not interfere and can transmit in the same time slot.

The scheduler executes the following three steps at each slot. First, for each link, it finds the flow with the maximum differential queue backlog. For example, for link (A, B) , the gray flow

has a difference of 0 packets and the black flow has a difference of 3 packets. The maximum value is then assigned as the weight of the link (see Fig. 1). Second, the scheduler selects the set of noninterfering links with the maximum sum of weights for transmission. This requires to compute the sum of link weights for each possible set. In the example, set $\{(A, B), (C, D)\}$ sums to $3 + 7 = 10$ and set $\{(A, C), (B, D)\}$ sums to $4 + 5 = 9$. The scheduler then selects the set with the maximum sum of weights, i.e., $\{(A, B), (C, D)\}$, to transmit at this slot. Finally, packets from the selected flows are transmitted on the selected links, i.e., black flow on link (A, B) and gray flow on link (C, D) . The same computation is then performed at every slot.

A. Backpressure Algorithm

More formally, the backpressure scheduling algorithm consists of the following steps executed for each time slot.

Flow Scheduling and Routing: For each link (i, j) , select the flow f_{ij}^* with the maximum queue differential backlog

$$f_{ij}^* = \arg \max_{f \in F} (q_i^f - q_j^f) \quad (1)$$

where q_i^f and q_j^f are the queue backlogs for flow f at nodes i and j , respectively, and F is the set of flows. The maximization in (1) implicitly performs routing by selecting the link (i, j) that each flow may use during the slot. The weight w_{ij} of each link is then selected as the weight of flow f_{ij}^*

$$w_{ij} = \max_{f \in F} (q_i^f - q_j^f). \quad (2)$$

Link Scheduling: Select the optimal link capacity vector $\mu^* = (\mu_{ij}^*)$ that satisfies

$$\mu^* = \arg \max_{\mu \in \Lambda} \sum_{(i,j)} \mu_{ij} w_{ij} \quad (3)$$

where $\mu = (\mu_{ij})$ are the link capacity vectors. The capacity μ_{ij} of each link (i, j) is the maximum rate in bits per second that the link can transmit subject to the channel state and the interference due to the other links in the vector. The set of all feasible link capacity vectors define the capacity region Λ .

Transmission: During the time slot, a selected link (i, j) transmits a packet from flow f_{ij}^* using rate μ_{ij}^* .

B. Congestion Control

The backpressure algorithm is throughput-optimal when the flow rates are within the capacity region. This issue can be addressed by combining the backpressure algorithm with the network utility maximization (NUM) framework, originally proposed for wireline networks [5]. This framework leads to a simple distributed congestion control algorithm where the source s of each flow f adjusts the flow rate x_f as

$$x_f^* = \arg \max_{x_f \geq 0} U_f(x_f) - x_f q_s^f = U_f'^{-1}(q_s^f) \quad (4)$$

where q_s^f is the queue backlog for flow f at the source s and $U_f'^{-1}(q_s^f)$ is the inverse of the first derivative of the utility function $U_f(\cdot)$ at the point q_s^f . In [6]–[8], it is proven that the congestion control scheme of (4) regulates the flow rates to be within

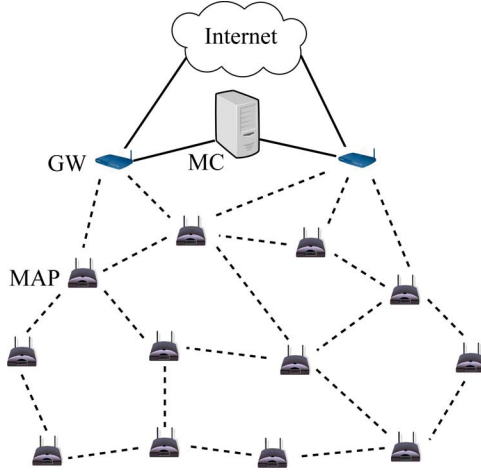


Fig. 2. XPRESS architecture, composed of MAPs to provide wireless coverage to mobile clients, GWs to provide Internet connectivity, and an MC for wireless scheduling.

the capacity region and cooperates with the backpressure scheduler to maximize throughput.

III. XPRESS DESIGN

This section presents the XPRESS system, a cross-layer backpressure architecture for wireless multihop networks. To our knowledge, XPRESS is the first system to implement backpressure scheduling over a time-slotted MAC, as it was originally proposed in theory. We first provide a high-level system overview, and then we detail the data plane and control plane designs. Finally, we describe the design of our backpressure scheduler with speculative scheduling.

A. Overview

In XPRESS, the wireless network is composed of several mesh access points (MAPs), a few gateways (GWs), and a mesh controller (MC), as depicted in Fig. 2. We use the term “node” to refer to a mesh node that can be either an MAP or a GW. The MAPs provide wireless connectivity to mobile clients and also operate as wireless routers, interconnecting with each other in a multihop fashion to forward user traffic. Mobile clients communicate with MAPs over a different channel, and thus are not required to run the XPRESS protocol stack. The GWs are connected to both the wireless network and the wired infrastructure and provide a bridge between the two. The MC is responsible for the coordination of the wireless transmissions in the network, and it is analogous to a switching control module. In our design, the MC is deployed in a dedicated node in the wired infrastructure and connects to the gateways through high-speed links. In an alternative design, the MC could be implemented within one of the gateways, if necessary.

At a high level, the operation of XPRESS can be described as follows. XPRESS runs a slotted MAC protocol, where a sequence of slots are organized into *frames*. For each slot, XPRESS selects a set of noninterfering links to transmit based on the flow queue lengths and the network state. Each node thus maintains per-flow queues and monitors adjacent links to estimate interference and losses. The queue lengths and network monitoring results are periodically transmitted to the

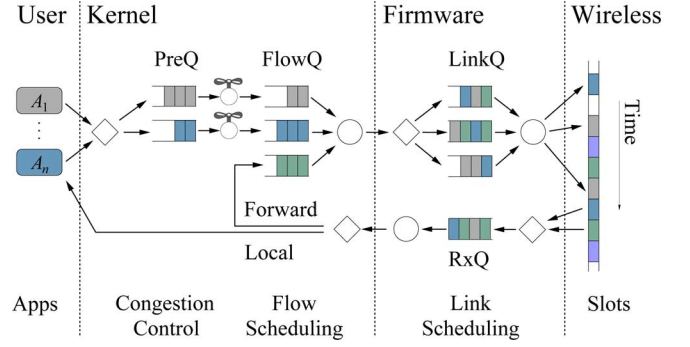


Fig. 3. Data plane at XPRESS nodes. Diamonds are packet classifiers, while circles are packet schedulers. Rate control and flow scheduling occur at the OS kernel, whereas link scheduling occurs at the network card firmware.

MC over an uplink control channel. Upon reception of this information, the MC updates its local topology and interference databases and runs the backpressure scheduler to calculate the throughput-optimal schedule for multiple upcoming slots, (i.e., a frame). The MC then disseminates the computed schedule to the nodes over a downlink control channel. The nodes in turn apply the new schedule for transmissions in the slots of the next frame. This cycle then repeats periodically.

B. Data Plane

The XPRESS data plane spans the transport, network, and MAC layers of the protocol stack, as depicted in Fig. 3. The transport and network layers implement congestion control and flow scheduling, respectively. The MAC layer implements link scheduling and a TDMA MAC protocol. The organization of these modules into host OS kernel and network interface card firmware depends on the architecture used. For convenience, Fig. 3 depicts this organization on our testbed hardware devices (cf. Section VI), where the full MAC firmware resides on the wireless cards, while the upper layers reside in the host OS kernel. In the figure, diamonds represent packet classifiers, while circles represent packet schedulers. The data flow from left to right is outgoing packets originating from the applications to the wireless medium; the data flow in the opposite direction is incoming packets that are routed or delivered to the applications. Packets in the slotted wireless medium (far right), which are neither incoming nor outgoing, represent transmissions between two other nodes in the network.

XPRESS Flows: In XPRESS, a flow is defined at the IP layer by its source and destination mesh nodes. Our design is general and can easily accommodate other flow definitions. However, compared to the usual 5-tuple flow definition of source and destination IP addresses, source and destination transport ports, and transport protocol (i.e., TCP or UDP), this design decision reduces processing and communication overhead in XPRESS at the expense of flow granularity.

Congestion Control and Flow Scheduling: Locally originated packets first pass through a flow classifier, represented by the leftmost diamond in Fig. 3. Each flow has two individual queues, namely, a pre-queue (PreQ) and a flow queue (FlowQ). After classification, packets are inserted into the PreQ and must pass through the congestion controller, represented by the

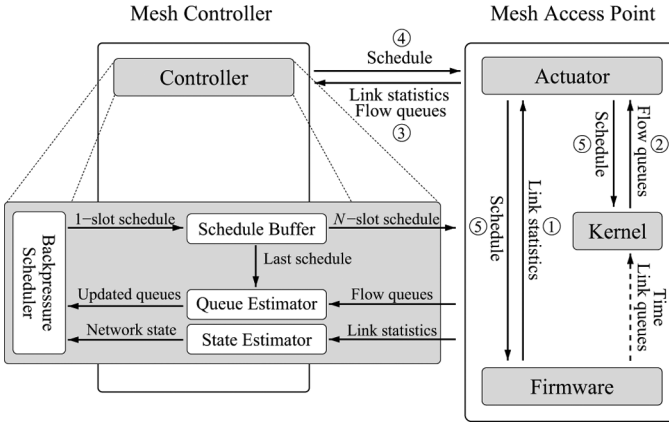


Fig. 4. XPRESS control plane and its system components. Numbered arrows indicate messages exchanged every frame, whereas the dashed arrow represents an internal message exchanged every slot between the network card firmware and the OS kernel of the MAP.

faucet handle in the figure. Congestion control is performed according to (4) and depends only on the length of the local FlowQ. A longer FlowQ reduces the allowed input rate, while a shorter FlowQ allows a higher rate. After congestion control, packets enter the FlowQ and wait to be scheduled. The kernel-space packet scheduler is synchronized with the slotted MAC with respect to time and link queue state. Just shortly before a scheduled transmission slot starts, the kernel scheduler dequeues a packet from the scheduled FlowQ and sends it down to the firmware link queue for transmission.

Link Scheduling: The MAC protocol keeps an individual queue for each neighbor in order to enable *link scheduling*, which allows a higher spatial reuse than *node scheduling*. As packets dequeued from the FlowQ arrive at the link-level packet classifier, they are classified according to the destination MAC address and inserted into the appropriate link queues (LinkQ). The slotted MAC, realized by a TDMA MAC protocol (cf. Section IV-C), maintains network-wide node synchronization and ensures that transmissions occur strictly within slot boundaries. When a transmission slot starts, the MAC protocol dequeues a packet from the scheduled LinkQ and transmits it over the air. If the transmission fails and the retransmission limit is not reached, the packet remains in the appropriate LinkQ until the next slot for the same neighbor.

Packet Reception and Forwarding: Once a packet is received, it is first filtered based on the destination MAC address and then inserted into a single receive queue (RxQ) at the firmware. The packet is delivered to the network layer at the kernel, where it is routed and tagged for local delivery or forwarding. In the latter case, the packet is inserted into the respective FlowQ and waits to be scheduled, just like a locally generated packet after passing congestion control.

C. Control Plane

The different system components must exchange control information to coordinate the network-wide transmissions. Fig. 4 depicts the XPRESS control plane with its different components at the MC and an MAP. The uplink control channel is represented by the arrows labeled from 1 to 3, while the downlink

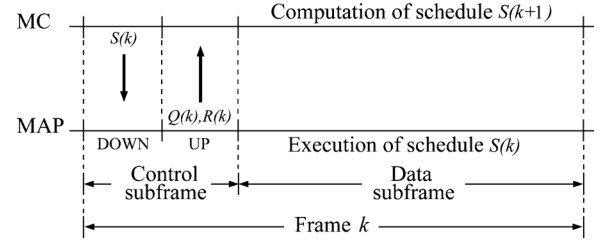


Fig. 5. Queue sampling and schedule computation in XPRESS. FlowQ lengths $Q(k)$ and rates $R(k)$ provided at the start of frame k are used to estimate $Q(k+1)$ and compute the schedule $S(k+1)$ of the next frame $k+1$.

control channel is represented by arrows 4 and 5. There is also an internal control channel (dashed arrow) between the MAC and the kernel packet scheduler, to synchronize them with respect to time and LinkQ state.

Uplink Channel: Each node runs an actuator application that communicates with the controller application at the MC. A schedule computation cycle begins at the start of a new frame. At this point, the TDMA MAC protocol notifies the actuator about the new frame and piggybacks in the same message the link statistics collected during the last frame, such as received signal strength and delivery ratio. This is the Step 1 in Fig. 4. The actuator then collects the FlowQ lengths (Step 2), combines all information, and sends it to the MC using the uplink channel (Step 3).

Mesh Controller (MC): The zoom of the MC in Fig. 4 shows the different steps taken for the calculation of a TDMA frame schedule composed of N slots. First, the link statistics are used to estimate the network state, namely, the interference relations and link loss rates. We explain this step in Section V. The FlowQ lengths are also collected and fed into a backpressure scheduler. The scheduler then uses the FlowQ lengths and network state to compute the network schedule, using the backpressure algorithm (cf. Section II) and our speculative scheduling technique, which we detail in Section III-D.

Downlink Channel: When the schedule computation is finished, the MC disseminates the new schedule using the downlink control channel (Step 4). The actuator receives this packet and forwards the new schedule both to the OS kernel as well as to the MAC layer (Step 5). The TDMA MAC starts using this new schedule for data transmission in the next frame. Packets will then be dequeued from the FlowQs to the LinkQs in accordance with this new schedule.

D. Backpressure Scheduler

XPRESS introduces a speculative scheduling technique to reduce scheduling overhead. This technique computes a schedule for a group of slots on a TDMA frame basis and performs the optimal backpressure computation for all slots in the frame based on speculated network queue state.

Fig. 5 depicts the speculative scheduling operation. At the start of TDMA frame k , the MC computes a schedule $S(k+1)$ for all data slots of the *next* frame $k+1$. This approach provides the MC with a time budget of one TDMA frame to perform optimal computations. However, this comes at the expense of uncertainty due to changes in the queue backlogs during frame k . These changes are due to: 1) incoming packets at

the source node of each flow; 2) wireless losses; and 3) the schedule $S(k)$ executed by the nodes during the slots of frame k . The MC addresses this uncertainty using the schedule $S(k)$ it computed during the previous frame $k - 1$, as well as the FlowQ lengths $Q(k)$ and source rate estimates $R(k)$ provided by the mesh nodes at the start of frame k . Additionally, as explained in Section V, the scheduler uses only links of high packet delivery ratios, which reduces the uncertainty of wireless losses.

First, the MC computes an estimate $\hat{Q}(k + 1)$ of the queue backlogs at the start of frame $k + 1$. The MC initially updates $Q(k)$ by adding $R(k)$ to the FlowQ of each source node as an estimate of the number of incoming packets at frame k . Then, the MC locally emulates the transmissions of schedule $S(k)$ on this updated queue state—that is, for each slot, the scheduled FlowQs of transmitters are decremented, and the FlowQs of receivers are incremented until $\hat{Q}(k + 1)$ is obtained.

Given $\hat{Q}(k + 1)$, the MC computes the schedule of the first slot of frame $k + 1$ with the backpressure algorithm (Section II). After this computation, the MC updates the scheduled FlowQs, creating a new queue estimate for the second slot (see Fig. 4). The rate $R(k)$ is also used to update the queue estimate of each slot according to the input rate at each source. With the new queue estimate, the backpressure algorithm computation is repeated. This process continues for the following slots, until the schedule $S(k + 1)$ of all data slots of frame $k + 1$ are calculated. After the computation, the MC transmits this schedule to the mesh nodes for execution during frame $k + 1$. In Section VI-D, we experimentally validate our speculative scheduling technique and show that the estimated queue backlog \hat{Q} closely follows the actual backlog Q .

IV. XPRESS IMPLEMENTATION

The XPRESS design is general and can be realized on a wide range of platforms. In this section, we describe the main components of our cross-layer implementation in the Linux OS and the firmware of our WiFi cards. We follow a top-down approach and describe these components in the order of the outgoing data path in Fig. 3.

A. Congestion Control

Congestion control is performed only at the source node of each flow by adjusting the flow input rate in accordance with (4). More precisely, the source rate x_f of each flow f is continuously adjusted to the optimal rate $x_f^* = U_f'^{-1}(q_s^f)$ for the flows to remain within the capacity region. In XPRESS, we use $U_f(x_f) = K \log(x_f)$ as the utility function, where K is a constant parameter defined later in Section VI. The logarithmic function allows a good tradeoff between fairness and efficiency in wireless networks [9]. The maximum allowed rate of each flow is then periodically adjusted to $x_f^* = K/q_s^f$, where q_s^f is the length of the FlowQ of flow f at the source s .

We implement this congestion control in the Linux kernel between the PreQ and FlowQ. Locally generated outgoing packets are intercepted using the Linux Netfilter local-out hook. These packets are classified by flow and put into their respective PreQ. A per-flow token-bucket rate controller adjusts the rate at which they are allowed to enter the FlowQ. This rate

control is performed by periodically inserting tokens into the bucket according to the optimal rate x_f^* . When there are enough tokens in the bucket, a packet is dequeued from the PreQ and sent to the respective FlowQ.

B. Queues and Scheduler

Flow Queues: Outgoing packets are intercepted using the Netfilter post-routing hook in the Linux kernel. Intercepted incoming packets that have been routed, and thus are ready to be forwarded, are classified and put into the corresponding FlowQ. We pass the FlowQ backlog information to the actuator module through the Linux `/proc` interface. The actuator in turn forwards this information over the uplink control channel to the MC for schedule computation.

Backpressure Scheduler: In our speculative scheduling, the schedule for each slot is computed at the MC using the backpressure algorithm (cf. Section II). In Section V, we show that this can be reduced to a maximum weight independent set (MWIS) computation in our system due to binary interference. The MWIS computation is an NP-hard problem in theory and can be computationally intensive in practice. Our C++ implementation is based on an algorithm for enumerating maximal independent sets at the beginning of each frame [10]. We then find the MWIS using a linear search over the independent sets. For efficiency, these sets could be stored in a heap structure keyed by their weights. At each slot, queue lengths change, which triggers a heap update. After the update, the new MWIS can be found as the root of the heap.

The MC runs on a Linux PC with a quad-core 2.7 GHz AMD Opteron computer, with 16 GB RAM. In our testbed, the execution time of our MWIS implementation on the MC never exceeded 100 μ s. In Section VII, we show that the MWIS problem can be solved efficiently for practical network sizes.

Kernel Packet Scheduler: The packet scheduler, in turn, is implemented in the MAPs as a kernel thread that waits for the computed schedule. The schedule contains information about which flows to transmit in each slot, as well as the next hop to be used by each flow. Based on this schedule and the instantaneous LinkQ lengths (see below), the packet scheduler decides when to dequeue packets from the FlowQs and send them to the LinkQs.

Link Queues: We implement LinkQs as individual neighbor queues at the firmware level to allow link scheduling. Individual neighbor queues are needed in link scheduling to prevent head-of-line blocking caused by a single common transmission queue. This problem may occur if we have wireless losses and the packet at the head of the queue is destined for a different neighbor than the one assigned for the slot, resulting in no packet transmitted during that slot.

Cross-Layer Queue Synchronization: Implementing link scheduling at the firmware is challenging due to the limited memory of wireless cards. As an example, our platform can only hold about a dozen full-size packets in the firmware at any time. As a result, the kernel and the firmware must be tightly synchronized with respect to both memory utilization and time to avoid *memory-exhaustion blocking* and slot underutilization. Memory-exhaustion blocking occurs if a given LinkQ has accumulated too many packets, leaving no memory space for

other queues. Once the firmware memory is full, no more packets can be sent from the FlowQs to *any* LinkQ. At the same time, the kernel must send a sufficient number of packets to populate the different LinkQs of the upcoming slots. In this case, transmission slots will remain unused if the scheduled LinkQ is empty.

We implement a cross-layer queue synchronization to address these issues. The firmware periodically advertises its LinkQ lengths as well as the current time slot to the kernel. In our implementation, this occurs at every slot for a tight synchronization. The kernel uses this information to send each packet to the firmware just two slots before its actual transmission, unless the corresponding LinkQ already has a sufficient number of packets in its queue.

C. TDMA MAC Protocol

We implement a time-slotted MAC protocol in the firmware of our 802.11a/b/g cards and disable some inherent CSMA functionality, including carrier sensing, backoff, RTS/CTS, and NAV. This gives us full control of transmission times and allow us to divide the MAC operation into frames and time slots. In our implementation, each TDMA frame is composed of a group of slots, which are divided into a control subframe (CS) followed by a data subframe (DS). We adopt our multihop TDMA slot synchronization protocol in [11] to synchronize MAPs over the wireless channel during the CS. After the control slots, nodes are synchronized, and the rest of the frame is composed of data slots. Each data slot is assigned to a particular set of noninterfering links based on the schedule computed by the MC. The MAC packet scheduler receives this schedule from the actuator, and then takes care of transmitting packets according to the schedule and within slot boundaries. The actuator notifies the MAC about the new schedule as soon as it receives it from the MC. In our case, a new schedule is applied on every frame.

V. INTERFERENCE ESTIMATION

We now introduce the design of our interference estimation technique to provide the backpressure scheduler with the link transmission sets and the corresponding capacities of their links. The capacity μ_{ij} of each link (i, j) is estimated on a TDMA frame timescale as $P_{ij} \times R_{ij}$, where P_{ij} is the packet delivery ratio (PDR) and R_{ij} is the PHY rate of link (i, j) during the TDMA frame.

Finding the link transmission sets and their capacities is a challenge because each link capacity depends on both the channel condition and the interference created by the other links in the set. A direct approach would be to enumerate and schedule each link set in the same slot, and then measure the PDR of their links. In a network with L links and N_r PHY rates, this requires $O(N_r \times 2^L)$ measurements during each TDMA frame, which is prohibitive.

In order to understand how interference manifests in our TDMA system, we perform the following experiment for all link pairs that do not share a node in our testbed.¹ The links of each pair simultaneously transmit backlogged UDP traffic for

¹In our current implementation, nodes have a single radio; therefore, links that share a common node always interfere.

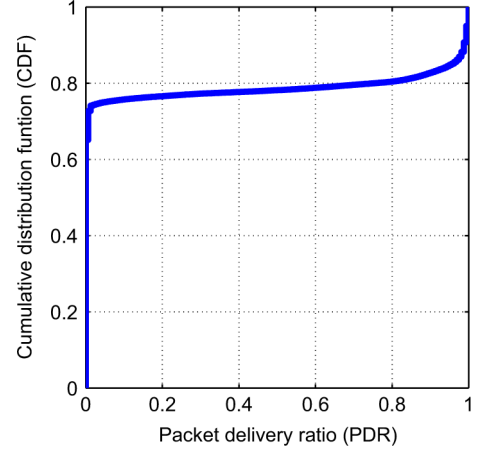


Fig. 6. The cumulative distribution function (CDF) of the TDMA frame PDR of links transmitting backlogged in pairs at 24 Mbps using the TDMA MAC protocol.

1 min using broadcast packets at 24 Mbps PHY rate. During this time, the receivers measure the RSS and PDR values at each TDMA frame. Fig. 6 depicts the cumulative distribution function (CDF) of the measured PDRs. We observe that interference creates a binary pattern where only 10% of links have intermediate PDR values between 0% and 90%. The dominance of the binary pattern is due to the clock synchronization of the TDMA protocol and the way 802.11a MAC decodes packets. According to our measurements, the TDMA protocol keeps the node clocks synchronized within $4 \mu\text{s}$. Thus, interfering transmissions overlap with packet reception at the very beginning of the 802.11a MAC preamble, where the capture effect is weak [12]. This is independent of the PHY data rate because it occurs at the beginning of the MAC preamble, where the lowest PHY data rate is used.

The binary interference phenomenon has important implications on our system. First, it is enough to estimate “low” or “high” PDR rather than exact values. Second, the link transmission sets consist only of links with “high” PDR and can be identified as the independent sets of a binary conflict graph. This bypasses exhaustive enumeration and measurement on all 2^L link combinations to determine the link transmission sets. Finally, the link scheduling problem of the backpressure algorithm, i.e., the selection in (3), reduces from an exhaustive search over all link capacity vectors to an MWIS problem in the conflict graph.

A. Interference Measurement

Nodes in XPRESS measure the RSS of each neighbor from the synchronization beacons they transmit in the control subframe of each TDMA frame. Each node has an assigned slot in this subframe, at which the node transmits its beacon using the lowest 802.11a PHY rate (i.e., 6 Mb/s) without interference. This results in one RSS measurement per frame for each link. Based on these measurements, we estimate the signal-to-interference ratio (SIR) of a link (i, j) under the interference of a transmitting node k as the difference $S_{ij} - S_{kj}$, between S_{ij} , the RSS of link (i, j) , and S_{kj} , the RSS of link (k, j) , both in dBm. If the SIR of the link exceeds a threshold, which depends on the PHY rate, the link PDR is estimated “high.” Since RSS

values are independent of PHY rate, this approach can estimate the PDR of all transmission sets for all PHY rates and reduce the number of measurement complexity from $O(N_R \times 2^L)$ to $O(N)$.

Despite its advantages, this RSS approach might be far from perfect for three reasons. First, the RSS values reported by 802.11 wireless cards can be highly inaccurate due to the type of hardware, poor calibration, environmental conditions, location, temperature, multipath effects, and external interference. Second, it relies on the RSS of decoded packets and hence cannot detect *hidden interferers* that are within interference range, but not within communication range. Finally, this approach only analyzes link pairs and thus does not consider the potential *additive interference* resulting from more than one transmitter.

To address these limitations, we use PDR measurements to complement the previous RSS technique. During the data subframe of each TDMA frame, the MAC at each node i estimates the PDR of link (i, j) as the fraction of successful MAC transmissions. At the end of the TDMA frame, the MAC reports these measurements to the MC. This information is then used by the MC to update the conflict graph at each TDMA frame, as described next.

B. Conflict Graph Update

The MC uses a conflict graph to represent interference in the network. A vertex v_{ij} in the conflict graph corresponds to the link (i, j) in the network graph. An edge between vertices v_{ij} and v_{kl} denotes interference between links (i, j) and (k, l) in either the DATA or ACK directions. The vertex independent sets in the conflict graph correspond to the link transmission sets.

The conflict graph update mechanism is executed at each TDMA frame, after the MC receives the RSS and PDR measurements. The conflict graph construction occurs in two stages. The first stage uses the RSS measurements, and the second stage refines it further with the PDR measurements.

First Stage: Recall that S_{ij} is the RSS of link (i, j) measured at node j once per frame. If S_{ij} is not measured, the transmission from i to j is estimated outside of communication and interference range. First, the MC creates a vertex v_{ij} in the conflict graph if $S_{ij} \geq \gamma_j^R$, the RX sensitivity threshold of receiver j at PHY rate R . Then, the MC adds an edge between each pair of vertices v_{ij} and v_{kl} in the conflict graph if either they share a common node or if the SIR of DATA or ACK directions is less than the SIR receiver threshold δ_j^R at PHY rate R .

Second Stage: For each link (i, j) selected by the first stage, the MC checks its reported PDR value P_{ij} . If $P_{ij} \geq 90\%$, the link remains in the conflict graph. Otherwise, the MC finds the link transmission set I_{ij} where this link failed the most by inspecting the schedule of the previous TDMA frame. Among the links in I_{ij} , a few were outside the interference range, and others were hidden interferers. To identify them, the MC first builds a connectivity graph where 1-hop neighbors are defined as the nodes from which a beacon is received, 2-hop neighbors are the 1-hop neighbors of these neighbors, and so on. To reduce the search space, we assume that the interferers of link (i, j) are those transmitters in I_{ij} , which are 2-hop neighbors of either

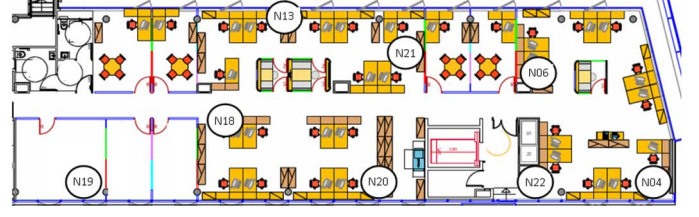


Fig. 7. Our wireless indoor testbed ($40 \times 8 \text{ m}^2$).

i or j . For each such node k , the algorithm adds an edge between vertex v_{ij} and vertex v_{kl} in the conflict graph, indicating that these two links will be scheduled in different slots in the next TDMA frame. If link (i, j) fails again in the same set I_{ij} during the next TDMA frame, the 3-hop neighbors of i and j are considered. While the link fails, this process continues in subsequent frames until all interferers of the link have been identified.

VI. TESTBED EVALUATION

In this section, we evaluate XPRESS in an indoor testbed. We first describe our testbed as well as the experimental methodology (Section VI-A), followed by results in single-hop (Section VI-B) and multihop scenarios (Section VI-C). Finally, we validate the queue length estimations in our speculative scheduling and show that XPRESS achieves the network capacity (Section VI-D).

A. Experimental Methodology

Fig. 7 shows our 8-node testbed deployed in one floor of our premises. This is a typical office environment with cubicles and glass meeting booths. Each node is a PC with 1.6 GHz CPU, 512 MB RAM, and 80 GB HD for storing measurement data. Each node runs the XPRESS software in the Linux kernel 2.6.28 and has a Technicolor WiFi card with a proprietary firmware that runs our multihop TDMA MAC protocol detailed in [11]. The WiFi cards are set to operate at 5 GHz in the 802.11a frequency band. The MC node is a dedicated Linux PC with a quad-core 2.7 GHz AMD Opteron computer, with 16 GB RAM. The execution time of our MWIS implementation on the MC did not exceed $100 \mu\text{s}$, and thus XPRESS supports backpressure computation for time slots as short as $100 \mu\text{s}$.

In our testbed, the TDMA control subframe consists of eight slots of $625 \mu\text{s}$ each, resulting in a total of 5 ms. In these slots, each node broadcasts its clock synchronization beacon at 6 Mbps PHY rate to other mesh nodes in the testbed, based on the protocol in [11]. The TDMA data subframe consists of 152 slots of $625 \mu\text{s}$ each, for a total of 95 ms. This results in a TDMA frame of 100 ms and a control overhead of 5%.

In our current prototype, multihop clock synchronization is performed using wireless transmissions in the TDMA control subframe, but other control information is exchanged using a reliable Ethernet connection. This provides a controlled testing environment, which is important to reproduce and to validate our experimental results. In our evaluation, however, we take the actual control overhead into account by making each slot in the TDMA control subframe larger than the duration required by the synchronization beacon. According to our analysis in Section VII, the 5-ms control subframe provides enough time to

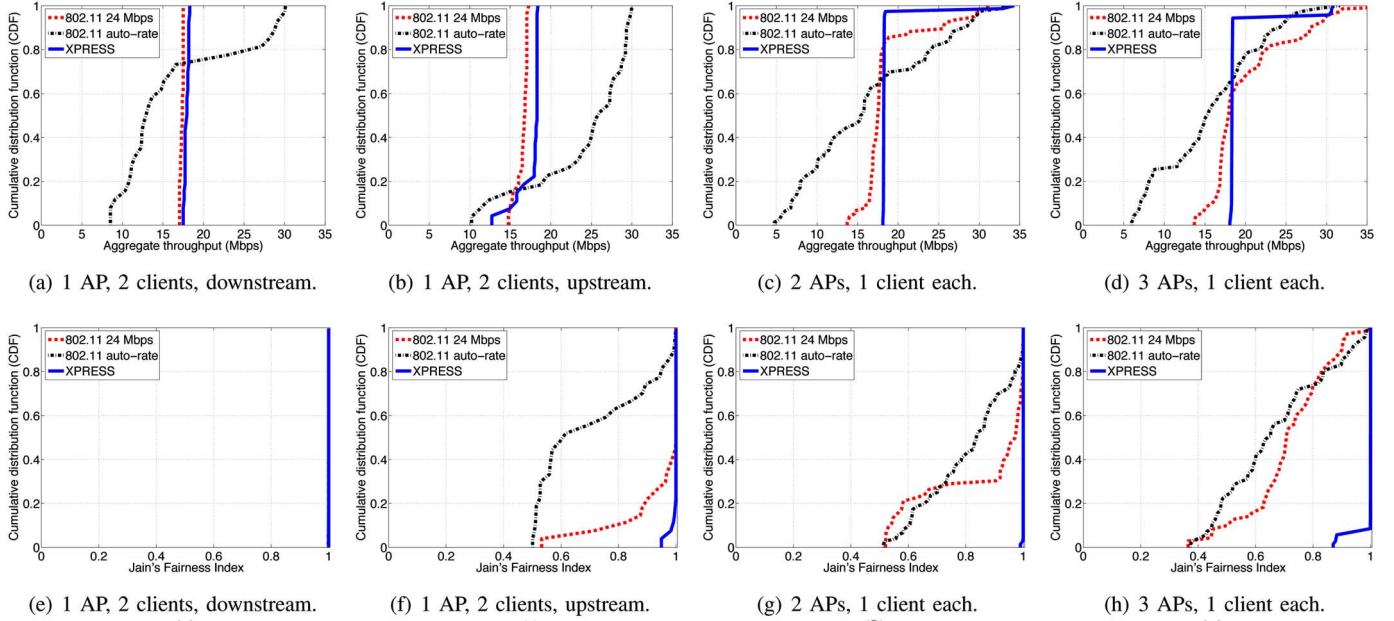


Fig. 8. (a)–(d) The throughput in different centralized WLAN scenarios. (e)–(h) The Jain's Fairness Index for the corresponding scenarios.

transmit all uplink and downlink control information at 6 Mbps in our 8-node testbed. As a result, the schedule information could in practice be included in the wireless synchronization beacons, which have shown to be highly reliable in our testbed.

Due to our dense deployment, we take the following steps to enable multihop routing and spatial reuse. First, we use the minimum transmission power of our WiFi platform (12 dBm). Second, we experimentally determine 24 Mbps as the highest PHY rate that can create multihop topologies without severely compromising link reliability. These power and PHY rate settings result in a total of 20 single links with a PDR higher than 90%, and in 20 link pairs with a PDR higher than 90% for each link when transmitting simultaneously. This provides enough spatial reuse for our testbed size and density.

In our experiments, we compare XPRESS to 802.11 DCF. For XPRESS, we fix the PHY rate for the data subframe to 24 Mbps. At this rate, the SIR threshold is defined as 17 dB based on the profile of RSS versus PDR collected for each link. The profile was constructed by offline measurements of links transmitting alone and in pairs. The 17 dB was the highest threshold that enabled a PDR greater than 90% in most cases. We set the upper limit of XPRESS FlowQs to 400 packets. The congestion control parameter K was experimentally determined as 100, a value that allowed a single link to achieve its full capacity for our logarithmic utility function. For 802.11 DCF, we use both a fixed 24 Mbps PHY rate and the automatic PHY rate adaptation scheme of our WiFi card (noted as *auto-rate* hereafter). We also disable the RTS/CTS exchange because it creates a high overhead for 802.11 DCF at the relatively high rate of 24 Mbps.

In order to maintain repeatability across different test runs, we select a 5-GHz channel free of external interference and set the MAC retransmission limit to 7 for both XPRESS and 802.11. We use *Iperf* [13] to generate UDP traffic with 1470-B payload packets and measure throughput as the goodput received at the flow destination.

B. Single-Hop WLAN Performance

In WLAN scenarios, XPRESS can operate as a centralized WLAN architecture, such as a wireless home network with a few access points (APs) controlled by the home gateway or an enterprise WLAN with lightweight APs and a central controller.

We consider four different WLAN scenarios: 1) one AP with two clients and only downstream traffic from the AP to the clients; 2) one AP with two clients and only upstream traffic from the clients to the AP; 3) two different APs, each with one client; and 4) three different APs, each with one client. For each scenario, we considered all possible combinations of nodes in our testbed. These results are representative for other scenarios that we tested.

Fig. 8(a)–(d) depicts the CDF of the aggregate throughput for XPRESS and 802.11. The figures show that XPRESS yields roughly 18 Mbps for most configurations in both scenarios. This is the maximum UDP throughput of a single link at 24 Mbps PHY rate, hence XPRESS scheduled links in sequence. Fig. 8(e)–(h) depicts the CDF of the Jain's Fairness Index (JFI) for the same scenarios and shows that the JFI of XPRESS is always close to unity, demonstrating almost perfect fairness.

We see that the aggregate XPRESS throughput is spread below and above the 802.11 throughput. When XPRESS performs sequential scheduling, 802.11 may achieve a higher throughput in two cases: when it performs successful parallel scheduling and/or when 802.11 auto-rate increases the PHY rate above 24 Mbps. For example, in Fig. 8(a) and (b), XPRESS achieves a higher throughput than 24-Mbps 802.11 in the majority of the configurations, with auto-rate outperforming XPRESS in 25% of the configurations in Fig. 8(a) and 85% in Fig. 8(b). In Fig. 8(c), the parallel scheduling in 24-Mbps 802.11 achieves higher throughput than XPRESS for 17% of the configurations; the higher PHY rates discovered by 802.11

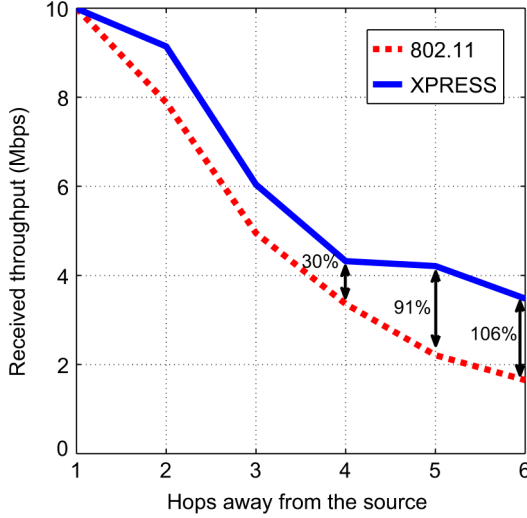


Fig. 9. Throughput of a line topology for both 802.11 and XPRESS for an offered load of 10 Mbps.

auto-rate increased this fraction to 35%. However, for the majority of the configurations, i.e., 83% and 65% respectively, XPRESS achieves much higher throughput, mainly due to successfully coping with hidden terminals and collisions. Fig. 8(d) shows a similar trend.

XPRESS schedules links in parallel in a few configurations. In these cases, XPRESS doubles the throughput, while the best 802.11 results are slightly lower due to 802.11 fairness issues, shown in Fig. 8(f)–(h). In particular, Fig. 8(f) shows that the higher throughput achieved by 802.11 auto-rate in Fig. 8(b) comes at the expense of severe unfairness. In this case, 50% of the configurations had a JFI lower than 0.6, and roughly 30% had a JFI of 0.5, the worst possible. This means that one client monopolized the wireless channel, while the other suffered from starvation. XPRESS, on the other hand, achieves perfect fairness in all these scenarios. The only scenario where 802.11 achieves perfect fairness is in Fig. 8(e) because, in this scenario, a single AP controls the downlink transmissions to both clients.

In summary, XPRESS can provide stable throughput and almost perfect fairness in small-scale WLANs while still achieving higher throughput than 802.11. On the other hand, 802.11 achieves variable throughput, and when 802.11 achieves a higher throughput than XPRESS, it may result in unfairness.

C. Multihop Mesh Performance

We now evaluate XPRESS in multihop mesh scenarios, consisting of line, gateway-centric, and multipath scenarios.

Line Scenario: We consider several single-path/single-flow configurations, ranging from 1 to 6 hops. For each configuration, we restrict the scheduler to use only the links in the path and use a 10-Mbps UDP traffic load. Fig. 9 depicts the average throughput of XPRESS and 802.11 versus the number of hops. The throughput of XPRESS decreases rapidly with the number of hops until 4 hops, where spatial reuse on the path starts being exploited. On the other hand, 802.11 continues to decrease rapidly even after 4 hops. The relative improvement of XPRESS over 802.11 increases with the path length. For instance, for an offered load of 10 Mbps, 802.11 achieves

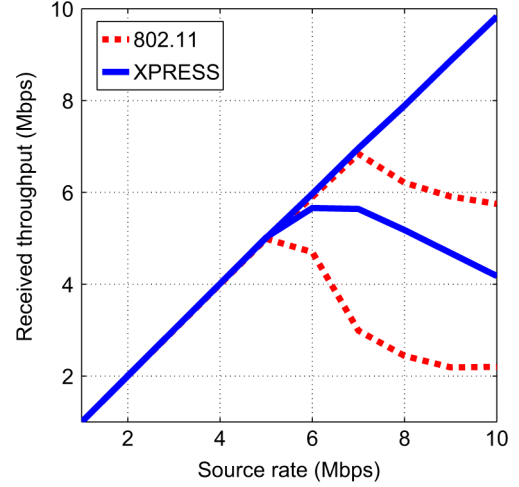


Fig. 10. Throughput of two flows $A \rightarrow B \rightarrow C$ and $B \rightarrow C$ in a GW-centric 2-hop topology.

3.3 Mbps at 4 hops, 2.2 Mbps at 5 hops, and 1.7 Mbps at 6 hops. On the other hand, XPRESS is able to transmit at 4.3 Mbps at 4 hops (30% gain), 4.2 Mbps at 5 hops (91% gain), and 3.5 Mbps at 6 hops (106% gain). This gain is due to the hidden/exposed terminals in 802.11 and the spatial reuse in XPRESS, which starts at 5 hops, where the aggregate throughput is $5 \times 4.2 = 21$ Mbps, as opposed to $4 \times 4.3 = 17.2$ Mbps at 4 hops. These results demonstrate the benefits of XPRESS backpressure scheduling over 802.11 DCF in achieving the full path capacity by an efficient coordination of transmissions.

Gateway-Centric Scenario: We consider a typical mesh scenario of a 2-hop gateway-centric topology with two flows, $A \rightarrow B \rightarrow C$ and $B \rightarrow C$, originating from two MAPs, A and B , toward a gateway C over a single path. Both flows have the same UDP source rate. Fig. 10 presents the per-flow throughput versus the offered load of each flow.

From the figure, we can see that both flows can simultaneously sustain up to 5 Mbps offered load each, for both XPRESS and 802.11. Above this, the aggregate throughput of 802.11 decreases and unfairness occurs at the expense of the 2-hop flow. On the other hand, XPRESS achieves consistently higher per-flow and aggregate throughput than 802.11. We can see that XPRESS serves both flows fairly until 6 Mbps. This is roughly the network capacity, where the 2-hop flow consumes $6 + 6 = 12$ Mbps and the 1-hop flow consumes 6 Mbps, for a total of 18 Mbps. As the source rate increases, we see that the medium time is still equally shared between the two flows. However, the longer flow needs 2 hops to reach the gateway, which reduces its throughput by half. On the other hand, the lack of coordination in 802.11 results in wireless losses due to severe contention, resulting in high throughput loss. At the highest source rate, the aggregate 802.11 throughput is only 7.9 Mbps, while XPRESS achieves 14.0 Mbps (77% gain).

Multipath Scenario: We now investigate the ability of XPRESS to exploit multipath capabilities in our testbed. In these experiments, packets may travel different paths between the same source and destination, depending on the per-slot instantaneous differential queue backlogs. We investigate two different aspects: throughput and end-to-end delay. We set up

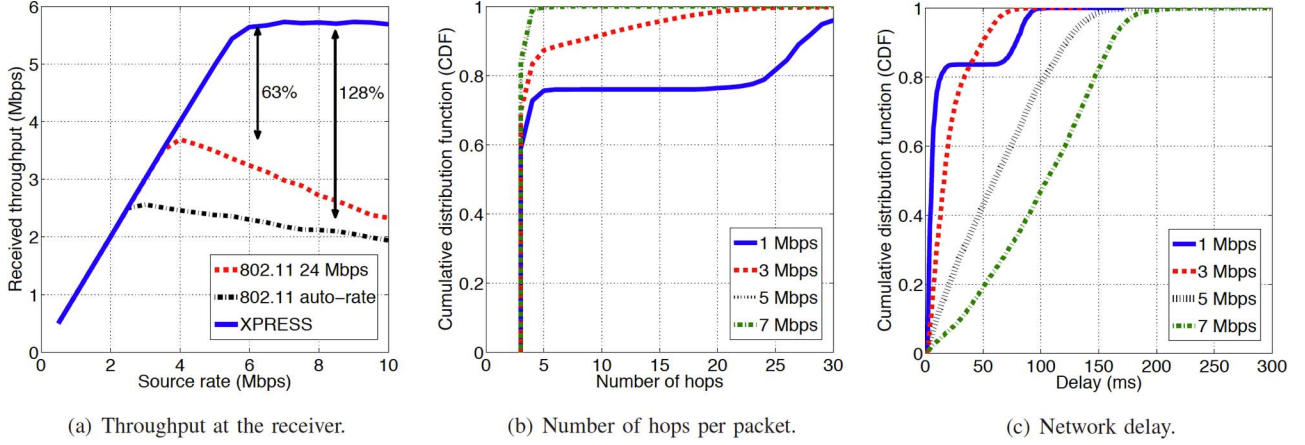


Fig. 11. The (a) throughput, (b) number of hops, and (c) delay for the multipath experiment.

a UDP flow between nodes 4 and 19, the farthest nodes in our testbed, and allow the XPRESS scheduler to use all possible links in the testbed. We compare XPRESS to 802.11 at fixed 24 Mbps PHY rate and auto-rate using the best path computed based on the ETX routing metric, through nodes 4, 22, 21, 18, 19 (cf. Fig. 7).

Fig. 11(a) depicts the received throughput at node 19 versus the input source rate at node 4. The throughput of XPRESS increases linearly with the offered load until 5.5 Mbps, after which it remains stable at the maximum of 5.7 Mbps. On the other hand, 802.11 only reaches a throughput of 3.5 Mbps (63% gain for XPRESS) with a fixed rate of 24 Mbps and 2.5 Mbps (128% gain for XPRESS) with auto-rate, after which throughput declines. The decline in 802.11 at high input rate occurs because of hidden-terminal collisions along the 4-hop path, which trigger packet retransmissions and reduce the end-to-end throughput. XPRESS does not suffer from hidden terminals and is able to sustain the maximum throughput. From the figure, we can also see that 802.11 auto-rate receives less throughput than 802.11 with 24 Mbps under high load. This occurs because the collisions often cause auto-rate to fall back to low PHY rates.

We now investigate the delay properties of XPRESS. In Fig. 11(b), we see the CDF of the path hop count taken by each packet, and in Fig. 11(c) we see the delay measured after the rate controller at the source. Fig. 11(b) shows that, under high loads, almost all packets follow 3-hop or 4-hop paths. As the load decreases, however, an increasing fraction of packets follows longer paths. At the lowest load of 1 Mbps, a fraction of 20% of the packets followed a path between 25 and 30 hops. The reason is that queues are small and the differential backlogs are not effective in path differentiation; this is an inherent property of backpressure scheduling [8]. However, as shown in Fig. 11(c), the delay of the slowest packets under 1 Mbps load does not exceed 100 ms. Thus, despite the long paths taken, once these packets exit the congestion controller and enter the backpressure scheduler, they are delivered to the destination in less than a TDMA frame duration. Delays increase after 5 Mbps load, which is close to the capacity limit of 5.7 Mbps shown in Fig. 11(a). However, for all offered loads beyond the capacity limit, the additional delay is limited to about 50 ms. In this case, the FlowQs are always full, and the scheduler delay

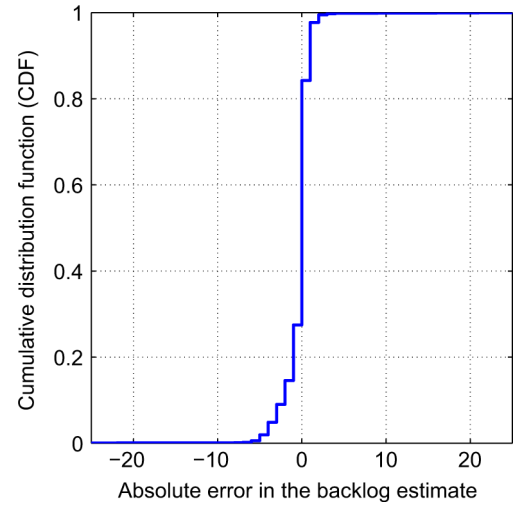


Fig. 12. The CDF of the absolute error in the backlog estimate $\hat{Q} - Q$ for each node at each slot. The estimates are close to the actual queue lengths, indicating that XPRESS achieves the network capacity.

is dominated by the maximum FlowQ buffer sizes. Also, the delays are finite, which indicates that the congestion controller feeds the backpressure scheduler with rates within the network capacity region.

D. Validating Queue Estimations

We now evaluate the queue estimation accuracy of the speculative technique in our scheduler (cf. Section III-D). At each 625- μ s slot in the last 60-s multipath experiment, we collect the estimated FlowQ backlog \hat{Q} and the actual FlowQ backlog Q of each node, obtaining more than 700 000 samples. Fig. 12 shows the CDF of the error $\hat{Q} - Q$ in the queue backlog estimates for all samples during the experiment. We see that the maximum error is only five packets, which is very small compared to the tens or hundreds of packets backlogged at the FlowQs at each node. This validates that XPRESS estimates queue backlogs correctly.

This result has also another important implication for XPRESS. It is well known from theoretical results that throughput optimality is still achieved with imperfect backlog estimates as long as the difference $\hat{Q} - Q$ between each backlog estimate and its actual value is bounded [14, Section 4.7]. The

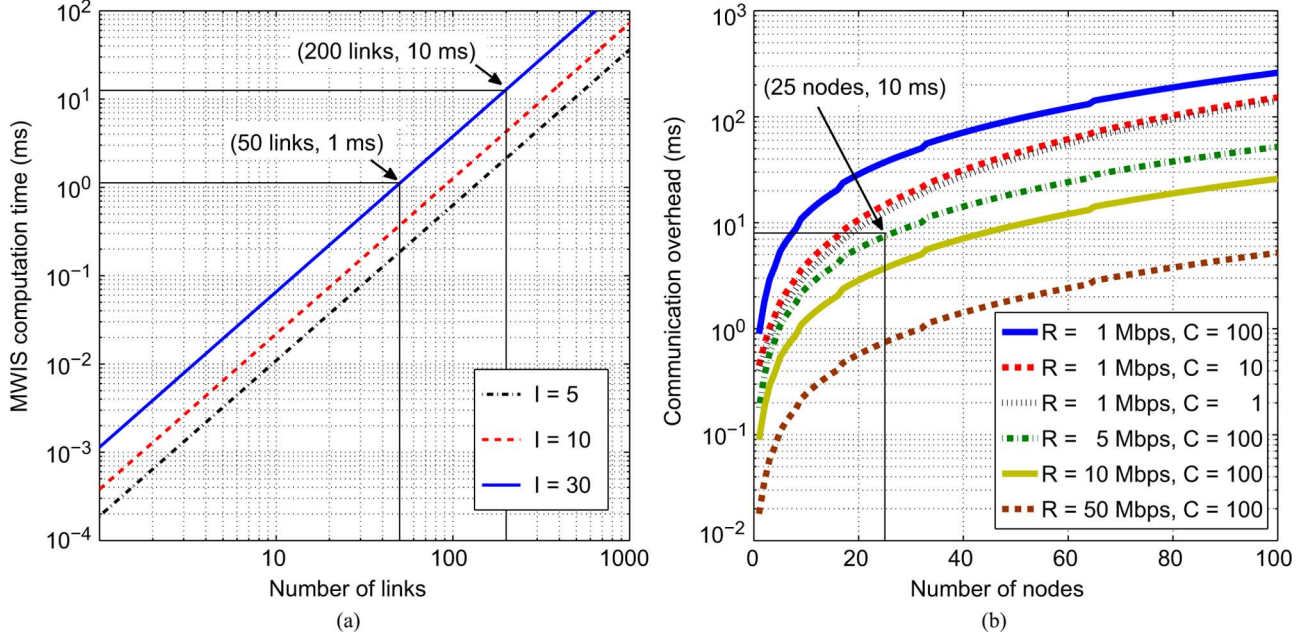


Fig. 13. (a) The mean MWIS computation time as a function of the number of links in the network. The implementation was on CPLEX running on top of an Intel 2.83 GHz E5440 CPU with 16 GB RAM. (b) The per-frame communication overhead as a function of the number of MAPs, for values of C and R . We assume $B = 256$, $D = 10$, $Q = 64$, $N_s = 50$, and $N_p = 10$.

small error shown in Fig. 12 thus supports the claim that, in addition to being significantly higher than the throughput achieved by 802.11, the throughput achieved by XPRESS in Fig. 11(a) is indeed the network capacity.

VII. OVERHEAD ANALYSIS

We now analyze the processing and communication overhead in XPRESS. Our goal is to understand how it scales for larger, yet practical, networks (e.g., 20, 50, 100 nodes) and the corresponding design tradeoffs. We believe that this limit is enough for practical purposes. For example, existing city-wide mesh networks are often split into smaller 10- to 20-node subnetworks, each communicating with a different gateway and operating over a separate channel. This occurs because: 1) multihop flows require several transmissions to reach the destination, which significantly limits the end-to-end throughput in large networks; and 2) the gateway quickly becomes a bottleneck as the number of nodes increases.

Computation Overhead: In theory, the MWIS problem is NP-hard. However, Wang [15] experimentally shows that for practical MWIS problems in wireless multihop networks, the computation time increases polynomially with the number of links and linearly with the average number of interferers of each link.

Fig. 13(a) is derived from the data in [15] and shows the mean MWIS computation time versus the number of links, for average densities I of 5, 10, and 30 interferers per link. With this figure, one can make design decisions about the frequency of the MWIS computation and the degree of multipath routing. For example, consider a network with 25 nodes and 1-ms TDMA slots. Fig. 13(a) shows that the MWIS computation in a very dense ($I = 30$) network of 50 links takes 1 ms, hence it can

run for every time slot. At the same density, the MWIS computation for 200 links takes 10 ms and hence can run once every 10 slots. In this case, the intermediate 10 slots would either reuse the last computed MWIS schedule or resort to heuristics to compute potentially suboptimal schedules. Another alternative is for the MC to prune a few links, which reduces the MWIS computation time and allows per-slot MWIS computations at the cost of less available paths.

Communication Overhead: The communication overhead is the time to transmit the control information on the uplink and downlink control channels (cf. Fig. 5). Assume that these channels can deliver this information at R bps. The downlink channel disseminates the synchronization beacon and C schedules for all N nodes. Each schedule specifies the selected set of links and the corresponding flows allowed to transmit. In the worst case, all nodes are selected for transmission in the same schedule. Thus, each schedule can be encoded with at most $N(\log_2 F + \log_2 D)$ bits, where F is the maximum number of flows and D is the maximum node degree in the network. If B bits are used for the synchronization beacon, an upper bound on the downlink time is

$$T_D = \frac{CN(\log_2 F + \log_2 D) + B}{R}. \quad (5)$$

The uplink overhead consists of N control packets sent by the nodes to the MC. Each packet contains the queue backlogs of all flows at the node, as well as the ID, the RSS, and the PDR of each neighbor. The uplink time T_U is thus

$$T_U = \frac{N}{R} [F \log_2 Q + D(\log_2 N + \log_2 N_s + \log_2 N_p)], \quad (6)$$

where Q is the maximum FlowQ size, N_s is the number of RSS levels measured by the wireless cards, and N_p is the number of PDR levels.

If we assume that there are only flows between the nodes and the gateway (i.e., $F = 2N$), we see from (5) and (6) that, for large network sizes, the uplink overhead $O(NF)$ is higher than the downlink overhead $O(N \log_2 F)$. However, for practical network sizes, they may be comparable depending on C .

Fig. 13(b) plots the total control overhead time $T_U + T_D$ as a function of N for values of R and C , assuming $F = 2N$. With this figure, one can make design decisions between the number of schedules C and the rate R . For instance, when the rate $R = 1$ Mbps, we see that increasing C from 10 to 100 only doubles the total overhead, i.e., same effect as halving R . This holds for any R , and it occurs because when $C = 100$, the uplink and downlink overhead become comparable. For the 25-node XPRESS network previously discussed, assume a requirement of 1-ms TDMA slots, with a control subframe of 10 ms and a data subframe of 100 ms (i.e., 10% overhead). Fig. 13(b) shows that this requirement can be satisfied with $R \geq 5$ Mbps and $C = 100$ schedules.

VIII. RELATED WORK

Backpressure Architectures: The backpressure algorithm was introduced by Tassiulas and Ephremides [1], and since then, a significant effort has been dedicated to distributed approximations that guarantee a fraction of the capacity region [16]. These algorithms maintain the slotted TDMA MAC protocol assumption of the original algorithm, but have not yet been translated to real implementations. On the other hand, recent work applies backpressure concepts to build practical systems. For instance, Wu *et al.* [17] prototype a synchronous and distributed SIR-based maximal matching MAC scheduler on a DSP/FPGA platform. In a different approach, several systems have been built on top of existing MAC protocols, such as 802.11. Akyol *et al.* [2] modify the 802.11 contention window to prioritize links with a higher differential backlog. Radunovic *et al.* [3] enhance the performance of multipath TCP transfers with a simplified backpressure scheme on top of 802.11. In a similar fashion, Aziz *et al.* [18], Ryu *et al.* [19], and Warriar *et al.* [4] approximate backpressure scheduling using prioritization in 802.11. Moeller *et al.* [20] build a backpressure routing architecture over 802.15.4 to enhance data collection in wireless sensor networks. More recently, Bhorkar *et al.* [21] and Choumas *et al.* [22] propose 802.11 backpressure implementations to reduce the end-to-end delay. Most of these works [2]–[4], [18], [21], [22] assume a separate routing protocol.

These approaches were important steps to show practical improvements of backpressure-inspired protocols over regular CSMA/CA. However, it is not clear how close they are to the throughput optimality of XPRESS since they are constrained by the CSMA/CA MAC interference estimation limitations, such as hidden and exposed terminals.

Centralized 802.11 Architectures: Vendors, like Cisco, Aruba, and Meru Networks, offer centralized architectures for enterprise 802.11 wireless LANs. In these architectures, a central controller coordinates downlink transmissions from the APs to the clients. With a similar goal, Liu *et al.* [23] and Srivastava *et al.* [24] propose centralized architectures to schedule AP transmissions. XPRESS shares the centralized

philosophy of these architectures, but also has fundamental differences. For instance, these architectures only schedule single-hop transmissions from the AP to the clients and use heuristics for scheduling. On the other hand, XPRESS controls the entire multihop wireless backbone and implements a throughput-optimal scheduler.

Interference Estimation: Existing techniques for interference estimation can be broadly classified as passive or active. Passive approaches require monitors deployed throughout the wireless network to collect traffic traces, which are later analyzed offline [25], [26]. Active approaches, on the other hand, use the available infrastructure to inject test packets into the network and measure interference. Padhye *et al.* [27] measure interference by comparing the throughput of two links, when transmitting in isolation and in parallel. Although accurate, this technique requires extensive downtime to test all link pairs. Ahmed *et al.* [28] propose an online approach where APs periodically silence their clients and run a quick interference tests. This technique, however, is limited to the downlink of single-hop wireless LANs. For multihop networks, Vutukuru *et al.* [29] propose a passive technique to learn about transmission conflicts. This approach exploits several exposed terminals unused by 802.11, but it does not address hidden interferers. Finally, several works measure RSS to predict PDRs [23], [30], [31]. In addition to RSS measurements, the XPRESS interference estimation design measures PDR during each TDMA frame to detect hidden or additive interferers and refine the network conflict graph.

IX. DISCUSSION

While XPRESS has shown significant improvements over 802.11, there are a few issues that must be dealt with for efficient operation. In this section, we describe these issues and present new research directions inspired by XPRESS.

Imperfect Scheduling: In any practical backpressure scheme, imperfect scheduling could occur due to error-prone schedule transmissions. In the case of not receiving the schedule, the node should fall back to a preconfigured default behavior, i.e., repeat the schedule of the last TDMA frame or refrain from transmitting. However, regardless of the node behavior, the net result in this case is that only a fraction of the capacity region can be achieved [16], with the remaining capacity being lost to either wireless collisions or missed transmission opportunities. Innovative techniques are thus necessary to still guarantee near-optimal performance despite losses in the control channel.

Delay Reduction: In backpressure scheduling, there are no preestablished routes; the route taken by a packet depends on the network congestion. As a result, packets may visit the same node more than once and create loops. This is particularly common in underloaded networks, as seen in Fig. 11(c). New delay reduction techniques [20], [32] as well as an analysis of their throughput–delay tradeoffs are then required to serve both elastic and inelastic traffic.

TCP Support: XPRESS operates at the IP layer, and therefore any transport protocol, including TCP, can run on top of XPRESS. Efficient TCP operation, however, requires handling out-of-order packets, due to the backpressure multipath routing, and supporting variable TCP packet sizes over the XPRESS

TDMA MAC. Out-of-order packets have been addressed with delayed packet delivery in [3], and variable packet sizes could also be addressed with extra implementation effort. Integrating TCP with an optimal backpressure scheduler and understanding their interactions, however, is a challenge for future research.

Rate Adaptation: The XPRESS design applies to both single-rate and multirate systems. For instance, the binary conflict graph representation holds for multiple rates, and the RSS-based interference measurement complexity $O(N)$ is independent of the rate. However, the main challenge in multirate systems is that the MWIS computational complexity increases with the number of rates. We believe that efficient rate adaptation in a centralized multihop network like XPRESS is a challenging topic for future investigation.

Optimal TDMA Frame Length: In XPRESS, the choice of the frame length is dictated by the tradeoff among three factors: 1) scheduling accuracy; 2) fast channel adaptation; and 3) control overhead minimization. Whereas a shorter TDMA frame provides accurate queue backlog estimates and quickly adapts to channel variations, it also incurs a higher control overhead. Our experiments show that a TDMA frame length of 100 ms provides high performance in static indoor environments. Determining the optimal frame length under highly dynamic channel conditions, however, is an interesting avenue for future research.

X. CONCLUSION

We presented the design and implementation of XPRESS, a backpressure architecture for wireless multihop networks. Our design leverages a centralized controller for obtaining throughput-optimal scheduling. In contrast to previous work, we integrated backpressure scheduling with a TDMA MAC protocol to allow precise timing in transmissions. Moreover, we introduced a novel interference estimation technique and an efficient speculative backpressure scheduler. Our results in an indoor testbed confirm that XPRESS achieves the full network capacity and up to 128% gains over 802.11. In our future work, apart from the topics discussed in Section IX, we also intend to evaluate XPRESS in larger networks. We believe our work opens up interesting avenues in wireless network system design, showing that optimal centralized routing and scheduling are feasible for small- to medium-sized wireless multihop networks. For larger networks, we believe that the design primitives of the XPRESS cross-layer protocol stack can lead to new distributed wireless architectures beyond 802.11.

REFERENCES

- [1] 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.
- [2] U. Akyol, M. Andrews, P. Gupta, J. Hobby, I. Saniee, and A. Stolyar, "Joint scheduling and congestion control in mobile ad-hoc networks," in *Proc. IEEE INFOCOM*, Apr. 2008, pp. 619–627.
- [3] B. Radunovic, C. Gkantsidis, D. Gunawardena, and P. Key, "Horizon: Balancing TCP over multiple paths in wireless mesh network," in *Proc. ACM MobiCom*, Sep. 2008, pp. 247–258.
- [4] A. Warrior, S. Janakiraman, S. Ha, and I. Rhee, "DiffQ: Practical differential backlog congestion control for wireless networks," in *Proc. IEEE INFOCOM*, Apr. 2009, pp. 262–270.
- [5] F. Kelly, A. Maulloo, and D. Tan, "Rate control in communication networks: Shadow prices, proportional fairness and stability," *J. Oper. Res. Soc.*, vol. 49, pp. 237–252, 1998.
- [6] L. Chen, S. Low, M. Chiang, and J. Doyle, "Cross-layer congestion control, routing and scheduling design in ad hoc wireless networks," in *Proc. IEEE INFOCOM*, Apr. 2006, pp. 1–12.
- [7] A. Eryilmaz and R. Srikant, "Joint congestion control, routing, and MAC for stability and fairness in wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 8, pp. 1514–1524, Aug. 2006.
- [8] X. Lin and N. B. Shroff, "Joint rate control and scheduling in multihop wireless networks," in *Proc. IEEE CDC*, Dec. 2004, vol. 2, pp. 1484–1489.
- [9] B. Radunovic and J.-Y. L. Boudec, "Rate performance objectives of multihop wireless networks," *IEEE Trans. Mobile Comput.*, vol. 3, no. 4, pp. 334–349, Oct.–Dec. 2004.
- [10] K. Makino and T. Uno, "New algorithms for enumerating all maximal cliques," in *Proc. 9th Scand. Workshop Algor. Theory*, Jul. 2004, pp. 260–272.
- [11] D. Koutsonikolas, T. Salonidis, H. Lundgren, P. LeGyade, C. Hu, and I. Sheriff, "TDM MAC protocol design and implementation for wireless mesh networks," in *Proc. ACM CoNEXT*, Dec. 2008, p. 28.
- [12] J. Lee, J. Ryu, S. Lee, and T. Kwon, "Improved modeling of IEEE 802.11a PHY through fine-grained measurements," *Comput. Netw.*, vol. 54, no. 4, pp. 641–657, Mar. 2009.
- [13] "iperf," 2011 [Online]. Available: <http://sourceforge.net/projects/iperf>
- [14] L. Georgiadis, M. J. Neely, and L. Tassiulas, "Resource allocation and cross-layer control in wireless networks," *Found. Trends Netw.*, vol. 1, no. 1, pp. 1–144, 2006.
- [15] P. Wang, "Throughput optimization of urban wireless mesh network," Ph.D. dissertation, Dept. Elect. Comput. Eng., Univ. Delaware, Newark, DE, USA, 2009.
- [16] X. Lin and N. B. Shroff, "The impact of imperfect scheduling on cross-layer congestion control in wireless networks," *IEEE Trans. Netw.*, vol. 14, no. 2, pp. 302–315, Apr. 2006.
- [17] X. Wu, S. Tavildar, S. Shakkottai, T. Richardson, J. Li, R. Laroia, and A. Jovicic, "FlashLinQ: A synchronous distributed scheduler for peer-to-peer ad hoc networks," in *Proc. Allerton Conf.*, Sep. 2010, pp. 514–521.
- [18] A. Aziz, D. Starobinski, P. Thiran, and A. ElFawal, "EZ-Flow: Removing turbulence in IEEE 802.11 wireless mesh networks without message passing," in *Proc. ACM CoNEXT*, Dec. 2009, pp. 73–84.
- [19] J. Ryu, V. Bhargava, N. Paine, and S. Shakkottai, "Back-pressure routing and rate control for ICNs," in *Proc. ACM MobiCom*, Sep. 2010, pp. 365–376.
- [20] S. Moeller, A. Sridharan, B. Krishnamachari, and O. Gnawali, "Routing without routes: The backpressure collection protocol," in *Proc. IEEE/ACM IPSN*, Apr. 2010, pp. 279–290.
- [21] A. Bhorkar, T. Javidi, and A. Snoren, "Achieving congestion diversity in wireless ad-hoc networks," in *Proc. IEEE INFOCOM Mini-Conf.*, Apr. 2011, pp. 521–525.
- [22] K. Choumas, T. Korakis, I. Koutsopoulos, and L. Tassiulas, "Implementation and end-to-end throughput evaluation of an IEEE 802.11 compliant version of the enhanced-backpressure algorithm," in *Proc. TridentCom*, Jun. 2012, pp. 64–80.
- [23] X. Liu, A. Sheth, M. Kaminsky, K. Papagiannaki, S. Seshan, and P. Steenkiste, "DIRC: Increasing indoor wireless capacity using directional antennas," in *Proc. ACM SIGCOMM*, Aug. 2009, pp. 171–182.
- [24] V. Shrivastava, N. Ahmed, S. Rayanchu, S. Banerjee, S. Keshav, K. Papagiannaki, and A. Mishra, "CENTAUR: Realizing the full potential of centralized WLANs through a hybrid data path," in *Proc. ACM MobiCom*, Sep. 2009, pp. 297–308.
- [25] Y.-C. Cheng, J. Bellardo, P. Benkö, A. Snoeren, G. Voelker, and S. Savage, "Jigsaw: Solving the puzzle of enterprise 802.11 analysis," in *Proc. ACM SIGCOMM*, 2006, pp. 39–50.
- [26] R. Mahajan, M. Rodrig, D. Wetherall, and J. Zahorjan, "Analyzing the MAC-level behavior of wireless networks in the wild," in *Proc. ACM SIGCOMM*, Aug. 2006, pp. 75–86.
- [27] J. Padhye, S. Agarwal, V. Padmanabhan, L. Qiu, A. Rao, and B. Zill, "Estimation of link interference in static multi-hop wireless networks," in *Proc. ACM IMC*, Oct. 2005, p. 28.
- [28] N. Ahmed, U. Ismail, S. Keshav, and K. Papagiannaki, "Online estimation of RF interference," in *Proc. ACM CoNEXT*, Dec. 2008, p. 4.
- [29] M. Vutukuru, K. Jamieson, and H. Balakrishnan, "Harnessing exposed terminals in wireless networks," in *Proc. NSDI*, Apr. 2008, pp. 59–72.
- [30] L. Qiu, Y. Zhang, F. Wang, M. Han, and R. Mahajan, "A general model of wireless interference," in *Proc. ACM MobiCom*, Oct. 2007, pp. 171–182.

- [31] C. Reis, R. Mahajan, M. Rodrig, D. Wetherall, and J. Zahorjan, "Measurement-based models of delivery and interference in static wireless networks," in *Proc. ACM SIGCOMM*, Sep. 2006, pp. 51–62.
- [32] S. S. L. Ying and A. Reddy, "On combining shortest-path and back-pressure routing over multihop wireless networks," in *Proc. IEEE INFOCOM*, Apr. 2009, pp. 1674–1682.



Rafael Laufer (S'09–M'11) received the B.Sc. and M.Sc. degrees in electrical engineering from Universidade Federal do Rio de Janeiro (UFRJ), Rio de Janeiro, Brazil, in 2003 and 2005, respectively, and the Ph.D. degree in computer science from the University of California, Los Angeles (UCLA), Los Angeles, CA, USA, in 2011.

He is currently a Member of the Technical Staff with Bell Laboratories, Alcatel-Lucent, Holmdel, NJ. During his Ph.D. studies, he worked as a Research Intern with Bell Labs, Paris, France; Cisco, San Jose, CA, USA; Technicolor, Paris, France; and EPFL, Lausanne, Switzerland. His major research interests are in wireless networks, network security, cloud computing, and software-defined networks.

Dr. Laufer is a member of the Association for Computing Machinery (ACM). He received the Marconi Society's Young Scholar Award in 2008 in recognition of outstanding academic achievement.



Theodoros Salonidis (S'98–M'04) received the Diploma in electronic and computer engineering from the Technical University of Crete, Chania, Greece, in 1997, and the M.S. and Ph.D. degrees in electrical engineering from the University of Maryland, College Park, MD, USA, in 1999 and 2004, respectively.

He is a Research Staff Member with the IBM T. J. Watson Research Center, Yorktown Heights, NY, USA. He was a Post-doctoral Researcher with Rice University, Houston, TX, USA from 2004 to 2006, and a Researcher with Intel Research, Cambridge, U.K., in 2006 and Thomson/

Technicolor, Paris, France, from 2007 to 2012. From 1999 to 2000, he was a Research Intern with the IBM T. J. Watson Research Center. His current research interests are in performance analysis, design, and implementation of mobile wireless networks and smart grid systems.

Dr. Salonidis is a member of the Association for Computing Machinery (ACM) and the Technical Chamber of Greece.



Henrik Lundgren (M'11) received the Doctor's degree in computer science from Uppsala University, Uppsala, Sweden, in 2005.

He has been a Senior Research Scientist with the Technicolor Research and Innovation Laboratory, Paris, France, since 2006. He was a Post-Doctoral Researcher with the University of California, Santa Barbara, CA, USA, from 2005 to 2006. Prior to his post-graduate studies, he worked from 1996 to 1998 as a Researcher and R&D Engineer in telecommunications with Ellemtel Utvecklings AB (a joint Ericsson and Telia research laboratory), Stockholm, Sweden. His current research interests include mobile and wireless communication systems.

Dr. Lundgren is a member of the Association for Computing Machinery (ACM).



Pascal Le Guyadec received the post-graduate degree in science from the Grenoble University, Grenoble, France, in 1993.

He has worked for several years in R&D for the Legrand Electric Group, France Telecom, and small IT startups. He joined the Advanced VoIP Laboratory, Technicolor, Paris, France, in 2002 to work with the networking aspects of conversational services for the home. Since 2006, he has worked as a Research Engineer with the Technicolor Paris Research & Innovation Center. He is currently in charge of devel-

oping and implementing research measurement tools, deploying testbed infrastructures, and contributing to different research activities.