

# Aggregating Bandwidth for Multihomed Mobile Collaborative Communities

Puneet Sharma,\* Sung-Ju Lee,\* Jack Brassil,\* and Kang G. Shin†

\*Mobile & Media Systems Lab, Hewlett-Packard Laboratories, Palo Alto, CA 94304

†Department of Electrical Engineering & Computer Science

University of Michigan, Ann Arbor, MI 48109

## Abstract

Multi-homed, mobile wireless computing and communication devices can spontaneously form communities to logically combine and share the bandwidth of each other's wide-area communication links using *inverse multiplexing*. But membership in such a community can be highly dynamic, as devices and their associated WWAN links randomly join and leave the community. We identify the issues and tradeoffs faced in designing a decentralized inverse multiplexing system in this challenging setting, and determine precisely how heterogeneous WWAN links should be characterized, and when they should be added to, or deleted from, the shared pool. We then propose methods of choosing the appropriate channels on which to assign newly-arriving application flows. Using video traffic as a motivating example, we demonstrate how significant performance gains can be realized by adapting allocation of the shared WWAN channels to specific application requirements. Our simulation and experimentation results show that collaborative bandwidth aggregation systems are, indeed, a practical and compelling means of achieving high-speed Internet access for groups of wireless computing devices beyond the reach of public or private access points.

## Index Terms

network communications, wireless communications, mobile computing, mobile communication systems, mobile environments

## I. INTRODUCTION

An increasing number of multi-homed wireless mobile computing devices are being equipped with two distinct types of wireless communication interfaces: a local area network (WLAN) interface such as IEEE 802.11x, and a wide area network (WWAN) interface such as a 2.5G or later generation cellular link. The capabilities of these interfaces differ greatly, most notably with the available WLAN bandwidth exceeding the WWAN's bandwidth by one to three orders of magnitude. For the foreseeable future we anticipate that this bandwidth disparity between local and wide area wireless network connections will remain intact.

Public high-speed Internet connectivity from such devices is now typically achieved by connection via the WLAN interface to an access point which is connected to a high-speed, wired network. It remains unlikely, however, that opportunistic deployment of these access points will ever realize ubiquitous — or even relatively geographically broad — access. Even where access points are densely deployed, seamless roaming between access points remains a technical challenge, and may not serve the business interests of either access point operators, venue owners or service providers. Further, even where access point coverage is rich, the transmission rate of the wired connection — typically 1.5 Mb/s — is limited and shared among a possibly large group of users, and unlikely to increase significantly in transmission speed in the foreseeable future.

To overcome the limited geographic coverage of public access points, we envision an alternative, complementary solution to high-speed Internet access through collaborative resource sharing. A group of wireless, mobile computing and communication devices in close proximity can dynamically form communities interconnected through their compatible high-speed WLAN interfaces; we call these ad hoc communities *Mobile Collaborative Communities* ( $MC^2$ s). Each  $MC^2$  member independently uses its WWAN interface to create a communication *channel* to an inverse multiplexer, and optionally offers to other members (full or partial) access to this channel. The set of participating channels connecting the  $MC^2$  members to the inverse multiplexer can be logically combined with an inverse multiplexing protocol to yield a higher-speed *aggregated channel* than is available from any one of the individual  $MC^2$  members.

The envisioned bandwidth aggregation mechanism is an enabling technology, as illustrated by the following example. A group of train commuters, who do not have Internet access through their WLAN interfaces, could spontaneously form a  $MC^2$ , and all members could receive a video stream delivered at a higher bandwidth — and higher quality — than any one member could receive. Each  $MC^2$  member would also enjoy higher speed, statistically-multiplexed WWAN access, a service often far more desirable than private, but lower-speed access.

Striping data across multiple, parallel communication channels is a conventional communications technique used to improve system performance or reliability in relatively statically-configured disk storage systems [35] and fixed, wired LAN–WAN interconnection systems [12,44]. In stark contrast to the multi-homed device scenario, due to end-device heterogeneity, mobility, and time-varying link transmission characteristics, the system we consider here is highly dynamic, and must be assembled, administered, and maintained in a decentralized fashion. This paper makes following key contributions:

- We present the design of a collaborative bandwidth aggregation architecture that is both practical and readily deployable. The proposed architecture enables third-party bandwidth aggregation service and does not require kernel modification at the end-clients.

- Based on simulations and testbed experiments we have shown that significant performance gains can be realized by application-level striping that adapts the shared WWAN link selection to the specific application requirements of the communication flows. As an illustration, we demonstrate how the quality of a hierarchically-layered video stream transmitted over lossy channels can be improved by a priority/application-aware traffic assignment.

We introduced our  $MC^2$  system in [40], and in this paper, we present a thorough design and extended evaluation through various scenarios and analysis in simulations and testbed experiments.

The rest of the paper is organized as follows. Section II explores the issues and tradeoffs faced in creating a decentralized inverse multiplexing system. Section III introduces algorithms for the assignment of application flows to heterogeneous WWAN channels, and Section IV describes the specific system architecture we chose to study. Performance evaluation results from an *ns*-based simulation are presented in Section V, and Section VI describes the implementation of a prototype system used to corroborate our findings. Related work is summarized in Section VII, and our conclusions are presented in the final section.

## II. ISSUES, CHALLENGES, AND APPROACHES

### A. Basic Setup

Let's first consider a relatively basic channel aggregation system. Assume that each shared channel contains only the single WAN link between each participating  $MC^2$  member and the inverse multiplexer. Suppose we seek to provide a single, bidirectional, unicast connection between an Internet source and a single  $MC^2$  node. End-system applications are oblivious to the presence of the aggregated channel in the downstream path; all upstream traffic follows the single WAN link associated with the shortest return path. No cross traffic is present on either LAN or WAN links. Each packet flowing downstream is received by a  $MC^2$  member, and immediately forwarded to the destination via the receiving device's LAN interface.

At the receiving end we assume that devices belong to a single, self-organizing  $MC^2$ . Each device is exactly one wireless LAN hop away from any other device.  $MC^2$  membership is dynamic; a newly-arriving device can join an existing community and contribute its (partial or full) WAN channel to the resource pool. Member devices can also leave the community — typically without prior announcement — due to either device failure or movement out-of-range of LAN communications. We assume that a mechanism exists to detect such departures from the resource pool, though packets may be lost until the system can detect and recover from that resource loss.

Even with this remarkably simple system model we are immediately faced with several intriguing questions. What is the performance loss associated with an unannounced departure of a single, actively-used WAN channel? How does this performance loss vary with the traffic type traversing the channel? What is

the minimum time duration that a newly-arriving channel participates in an aggregated channel such that its throughput is increased?

To begin to address this last question, consider the decision of whether to run a TCP connection over a single, persistent link, or an inverse-multiplexed connection comprising that same link *plus* a second link of equal capacity alternating between connected and disconnected states. It is intuitive that the multiplexed connection might *promise* greater average bandwidth capacity, but the fluctuating presence of the second link may result in TCP window size reductions in response to packet losses, such that the two links can have lower throughput than the single persistent link. Selecting links to maximize the throughput of a TCP connection is an even greater challenge when more than two links are available, and each of the links has different communication characteristics; see the appendix of our technical report [39] for a more rigorous development of a method for best channel selection.

The challenge of designing an effective inverse multiplexing system becomes far harder when we recognize that the components are heterogeneous, imperfect, and supporting time-varying workloads. For example, WAN link transmission characteristics (i.e., bandwidth, packet latency, loss) will vary, possibly dramatically as end devices move around. Links from different service providers may be of dissimilar technologies with different costs, complicating link selection. Links of the same type from a single network operator might have dependent or correlated transmission characteristics or outages.

The potentially large latencies introduced by packet forwarding through power- and processing-limited mobile computing devices is also a challenge. Disparities in the forwarding latency on different paths traversing heterogeneous computing devices with time-varying computing workloads can introduce packet misordering in the end-to-end path that can affect certain applications adversely. For example, non-interactive multimedia streaming applications will typically be lightly affected, though larger client buffer capacities might be desired. Although packet reordering might not reduce multimedia application performance noticeably, it can complicate TCP RTT computation and decrease TCP throughput. Packet reordering is not uncommon in today's Internet [4], and in the event that reordering becomes significant, there are approaches that can mitigate performance degradation [5].

## B. Multiplexing Layer

A key issue in our overall system design is the identification of the preferred protocol layer for the multiplexing function. Since IP performs routing and multiplexing, it is natural to consider a network layer multiplexing implementation. An IP-based solution could be implemented exclusively at the communicating end-systems; in this case any packet scheduling, reordering, and reassembly would occur, as usual, only at the source and the destination. Though such a network layer implementation can be achieved in several ways, each requires end-system kernel modification, restricting the availability of channel aggregation to

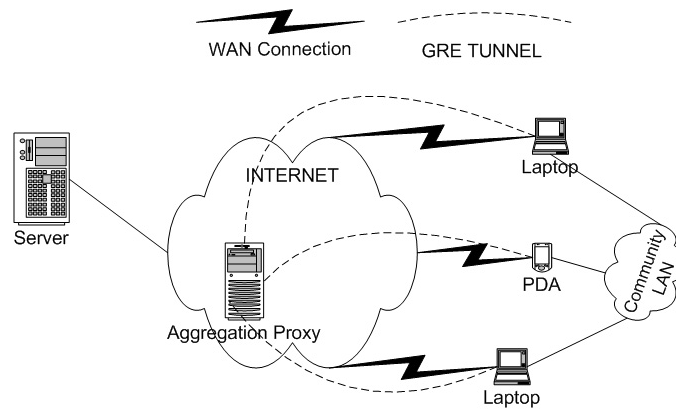


Fig. 1. A bandwidth aggregation service architecture.

data transfers between modified end-systems. An additional disadvantage of network layer striping is that it could restrict the channel assignment policies (i.e., the intelligent mappings of flows to available channels) that we might seek to implement, since the network layer is generally not aware of application characteristics and requirements. Performing multiplexing at the network layer, however, does have the advantage that it would not require any changes to existing applications.

An alternative solution is to perform multiplexing at the transport layer. Once again, end-system protocol stacks would require modifications, though transport-layer channel assignment policies could potentially be made more easily aware of application requirements. The obvious deployment issues associated with either network- or transport-layer multiplexing suggest solutions using application-layer multiplexing. Although such an implementation would incur more packet processing overhead, it requires no kernel modification and is easy to install, maintain and monitor. Application-layer multiplexing also permits controlling packet scheduling on a per-application, per-connection or per-packet priority basis. Like transport-layer and network-layer multiplexing, application-layer multiplexing is also transparent to the applications and does not require modifications to the applications themselves.

### C. Forwarding Mechanism

What forwarding mechanism should an inverse multiplexer use to transmit a packet over a chosen channel? Irrespective of a packet's destination, different packets must traverse different routes. There are several means of achieving this. One approach is to change each packet's destination address to the IP address of the appropriate  $MC^2$  member's WAN interface. When a packet arrives at the  $MC^2$ , its destination address would be reverted back to the original  $MC^2$  member destination address. This would, in a sense, be similar to providing a Network Address Translation (NAT) service, albeit in a distributed manner. But packet modification and processing overhead at the forwarding nodes associated with this approach might be significant.

Another packet forwarding approach could use *loose source routing* to forward a packet through the intermediary interfaces associated with the desired WAN channel to traverse. This would avoid the need to provide a special NAT-like packet forwarding service beyond ordinary IP routing itself. However, loose source routing has multiple, well-known weaknesses (e.g., use of IP options, extra router processing) as well as limited router support, making its use largely unworkable.

A preferred packet forwarding implementation would use *tunnels* between the inverse multiplexer and each  $MC^2$  node. Tunneling has long been used to establish static paths [47], and most OS network stacks today have built-in support for tunnels. In such a system packet forwarding would operate as follows. Unicast packets sent from an Internet-connected source would be routed normally to the inverse multiplexer, where each would then be forwarded, according to the multiplexer’s flow-to-channel assignment policy, to the tunnel corresponding to the appropriate WAN channel. Upon arrival at the  $MC^2$  node, the packet would be decapsulated and forwarded on the wireless LAN to its intended destination. In this simple case, all upstream traffic would be sent over a single WAN link, typically — but not necessarily — the receiver’s own. Figure 1 shows a bandwidth aggregation service architecture using Generic Routing Encapsulation (GRE) [13] tunnels.

#### D. Proxy Placement

Another key question in the design of our system is the appropriate placement of the inverse multiplexer in the end-to-end connection. In principle, this function can be located at almost any point between the WAN link terminations and the connection end-point (e.g., origin server), including the end-point itself. The preferred location depends on many factors including the type of WAN links, who is providing the aggregation service, whether collaborating devices agree to connect to a common multiplexing point, and how generally accessible the multiplexing service must be from a wide range of origin servers.

If all the WAN links from a  $MC^2$  terminate at the same point, a preferred location for the inverse multiplexer is that termination point. It is natural to think of a *proxy* providing this service, and to ease our discussion, we will simply use this term to refer to the location of the inverse multiplexer, regardless of whether a distinct physical component is used to implement the function. If the proxy is located near the WAN link termination points, then it is likely easier and more efficient for a wide range of services to use the proxy to transfer data to the  $MC^2$ . In such a case the aggregation service can be provided as a value-added service by the  $MC^2$  members’ common Internet Service Provider (ISP). The proxy can alternatively be located at the network edge close to the origin server, or even at the origin server itself. While this location avoids the potential restriction of requiring a common WAN link termination point,  $MC^2$  members might have to communicate with different aggregation services to communicate with different servers. As we will

describe later, one can also envision an aggregation service being provided by a third party by placing proxy in the middle of the network.

### *E. Communication Dynamics*

Both the availability and the quality of shared communication channels can be expected to vary with time. Yet these dynamics must be monitored and communicated to the task mapper to provide proper task assignments. Community membership itself will change as hosts join and leave the community, due to either end-system failures (e.g., power loss) or simply moving out-of-range of WLAN communications. WWAN channel quality may change often and unpredictably because of fading, interference, and location-dependent coverage gaps. Delay and delay jitter will change as the heterogeneous, CPU-limited devices forwarding packets between WWAN and WLAN interfaces are subject to time-varying computing workloads.

How do we monitor these communication dynamics? One possible approach is to deploy a *monitoring agent* within each community to oversee the communication environment. The main advantage of deploying a monitoring agent here is the ease of detecting changes to community membership quickly. The drawback, however, is that if the agent resides on a single node it is difficult to monitor the WAN channel performance of other nodes in the community. Moreover, relying on a single (or even a few) monitor(s) can result in both a performance and reliability bottleneck.

This problem can be solved by either replicating the monitoring agent or making every member a monitoring agent, i.e., distributed monitoring. Distributed monitoring works as follows. Each member broadcasts its channel characteristics and associated information (e.g., communication costs and its energy balance) either periodically, upon detection of an event, or when a threshold is exceeded. Each broadcast is timestamped. Upon receiving such a broadcast all the other members update the corresponding entry of their copy of the community communication status database. The task mapper can obtain a copy of the database in either of two ways. First, the task mapper can request a copy of the database from any community member. Requests can be sent to a randomly-selected member, or a member identified by inspection of the most recent database the task mapper has received. For example, an inquiry might be directed to a member with ample advertised available processing power, residual energy, or network bandwidth. An inquiry might be issued periodically, or be driven by an event such as the need to remap channels for a newly-arriving flow. The second way that a task mapper can obtain the database is simply by receiving an update report periodically or when a monitoring agent observes a significant local event (e.g., sudden channel failure).

An alternative approach is to perform monitoring at the location where the task mapping is done. When transport protocols that require end-to-end feedback (e.g., TCP) are used, end-hosts can piggyback their measured channel condition and membership information (e.g., will move out, will turn the device off because of low battery power, etc.) on the acknowledgment/feedback packets. Out-of-band signaling may

be utilized when transport protocols without end-to-end handshaking such as UDP are running. Although this approach has the advantage of easily monitoring communication channel state, there may be delay in obtaining the information from every member. To make matters worse, the feedback messages may not reach the other end, due to wireless transmission errors or unidirectional/asymmetric links. Hence, the right preference could be to use a combination of the two approaches.

### F. Collaboration Incentives

Aggregated bandwidth channels can be realized only when hosts willingly collaborate by sharing their communication channels. Willingness to collaborate is not an issue for a single user with multiple mobile devices (e.g., cell phone, PDA, laptop, etc.) forming a ‘community’ (i.e., personal area network), nor might it be an issue for colleagues or acquaintances. But what are the incentives for collaboration between hosts owned by multiple parties with little or no pre-existing relationship? Clearly, if many community members seek access to the same content (e.g., multicast video) then the members will be well-motivated to take advantage of faster download or streaming. But if each host seeks to receive unique content, will the community members each be willing to sacrifice their computing and communications resources?

Such questions regarding the social use of shared resources have occurred in many settings, and have arisen recently in investigations of forwarding incentives in ad hoc network routing [7, 10, 36, 43, 48]. In ad hoc networks, when the destination node is outside the radio transmission range of the source node, the communication end-points rely on intermediate nodes on the path to forward the packets for them. The source and destination nodes, in turn, will forward data packets when they become the intermediate nodes for other communication pairs.

We have reason to be optimistic that collaborating communities will occur spontaneously in public settings. The benefits of aggregated channels are high, particularly when members seek to access large files or high-quality media objects. In some cases no viable alternative means of access will exist. In addition, the ‘cost’ of the resources being offered (e.g., bandwidth) is relatively low, members can experiment at low risk by offering partial resources, and offered resources may be reclaimed with little difficulty. Moreover, users will tend to share their WWAN bandwidth if the aggregation is performed only on *idle* channels. The recent success of peer-to-peer file sharing leads us to believe that device owners may be willing to share communication and computational resources as readily as they do information, particularly if they directly and immediately benefit from resource sharing.

In the following, we provide a simple analytical model of incentives for sharing bandwidth in our  $MC^2$  network. Let the set of nodes  $n$  in a community be represented by  $\mathcal{I}$ . Every node  $i \in \mathcal{I}$  is connected to all other nodes via a WLAN with bandwidth  $B$  in addition to its WWAN connection whose bandwidth is



represented by  $O_i$ . We denote the bandwidth that node  $j$  offers to node  $i$  by  $O_{ji}$ . We also assume that  $O_i = O_j, \forall i, j \in \mathcal{I}$ .

We devise the following algorithm for the community members for sharing bandwidth. Every node  $i$  in the system keeps a set of normalized weights that we denote by  $\mathcal{W}_i$ , that contains a value for every node in the system. It is computed as follows:

$$\mathcal{W}_i = [w_{i,0}, w_{i,1}, \dots, w_{i,n-1}] \quad (1)$$

$$w_{i,j} = \frac{O_{ji}}{\sum_k O_{ki}}, \forall k \in \mathcal{I} \quad (2)$$

Then node  $i$  will compute how much bandwidth to offer a node  $j$  as:

$$O_{ij} = w_{i,j} \cdot O_i \quad (3)$$

We claim that this algorithm offers incentives to nodes in the system to share their bandwidth as it will directly affect their performance. In what follows, we prove this claim by showing that the system has a Nash Equilibrium.

We define the utility of node  $i$  as:

$$u_i = \frac{\sum_k O_{ki}}{\sum_k O_{ik}} \quad (4)$$

Intuitively speaking, a node  $i$  wants to offer as little as possible of its bandwidth while getting as much as possible from all the WWAN bandwidth of the nodes in the network. Node  $i$  would want to maximize its utility. The obvious solution of setting the denominator to zero (offering no bandwidth to any node) will backfire, since all nodes will set their  $w_{j,i}$  to zero and  $u_i$  will become zero. Thus, we need to maximize the numerator. In fact, we have:

$$\begin{aligned} \max \sum_j O_{ji} &= \sum_j \max(O_{ji}) = \sum_j \max(w_{j,i} \cdot O_j) \\ \max(w_{j,i} \cdot O_j) &= \max\left(\frac{O_{ij} \cdot O_j}{\sum_k O_{kj}}\right) \end{aligned}$$

However, we already know that  $O_j$  is constant, and thus the problem becomes that of maximizing:

$$\max \frac{O_{ij}}{\sum_k O_{kj}} = \max \frac{1}{1 + \rho}$$

where  $\rho = \frac{\sum_{k \neq j} O_{kj}}{O_{ij}}$  should be minimized by maximizing  $O_{ij}$ ; the bandwidth that node  $i$  offers to node  $j$ .

In fact, the best a node  $i$  can do is to offer all of its WWAN bandwidth  $O_i$  and compete for  $\sum_{j \neq i} O_j$ . As a result, it will acquire WWAN bandwidth  $BO_i$  that is bounded as:

$$O_i \leq BO_i \leq (n-1) \cdot O_i$$

In the remainder of this paper we will focus on the proper design of a proxy's channel allocation and packet striping algorithms, and show that such a design can achieve significant performance gains. There are

many other intriguing issues beyond the scope of this paper, including policing malicious  $MC^2$  community members [26], retaining privacy of information transmitted through  $MC^2$  devices, as well as a host of questions associated with more sophisticated topologies involving devices participating in multiple  $MC^2$ s, connections transmitted over one or more inverse-multiplexed hops (in sequence or in parallel), and fairness issues between flows being assigned to channels.

### III. CHANNEL ALLOCATION AND PACKET STRIPING

For each active flow a proxy is responsible for two tasks. First, the proxy must select a set of channels on which to forward packets to the  $MC^2$  destination. Second, the proxy must intelligently stripe arriving packets across those channels. Efficient channel allocation and striping algorithms map or remap the flows to the channels based on both application requirements and the number and the condition of available channels. Hence, the algorithms we examine in this section are both *application-aware* and *channel-adaptive*. As an example, the algorithms we consider would seek to assign a flow from an audio or video source to channels that would maintain that application's stringent delay or delay jitter requirements, while assigning bulk data transfer (e.g., FTP) flows to channels that might incur longer delays but are reliable. Of course, both the number and condition of assignable channels might vary over a given flow's lifetime. Channel allocation and striping algorithms can be categorized along the following orthogonal dimensions:

- **Channel-adaptive:** These algorithms assign packets on different channels according to the channel conditions such as bandwidth, loss, and delay. For example, a Weighted Round Robin (WRR) algorithm stripes packets to channels in proportion to each channel's available bandwidth.
- **Application-aware:** Striping algorithms can also use knowledge or a *profile* of an application flow and its end-system requirements for channel selection and packet striping. Since applications can have different profiles, each application would potentially need a different algorithm. These algorithms promise to provide better performance than application-agnostic algorithms, but they have the burden of obtaining information about a flow's requirements. This information can be obtained explicitly from the traffic source, or may be inferred by examining the flow itself, or some combination of both. For instance, a source might mark its packets (e.g., ToS field in the IP header) or a proxy might infer application type from destination information (e.g., TCP or UDP port numbers) or even the application payload.

A given striping algorithm can be both channel-adaptive and application-aware, as summarized in Table I.<sup>1</sup>

We now (1) define and characterize application requirements, and (2) describe how to monitor and update channel characteristics before presenting illustrative algorithms for intelligently mapping/remapping and scheduling flows to available channels, using video as a motivating example.

<sup>1</sup>(a) We propose layer-priority striping for hierarchically-layered videos in Section III-D. (b) Application-aware algorithms are application-specific and also require channel information.

TABLE I

CATEGORIZATION OF CHANNEL ALLOCATION AND SCHEDULING.

	Application-aware	Application-agnostic
Channel adaptive	Layer Priority Striping <sup>a</sup>	WRR, WFQ
Channel non-adaptive	Not applicable <sup>b</sup>	Random, Round-robin

### A. Application Characteristics

Each application flow can be described by itself (intra-characterization) or against other application flows (inter-characterization). Examples of the former include Multiple Description video Coding (MDC) [2, 46] and the imprecise computation model [22] that is widely used in the real-time computing community. That is, an application flow has multiple representations or versions expressing different degrees of satisfaction (being minimally-to-fully satisfactory). The proxy must allocate and schedule resources to at least guarantee the minimum degree of satisfaction for each given application flow. That is, timely delivery of the base layer or essential part of each application flow must be guaranteed, and the enhancement layer or the optional part receives lower priority. For more general types of applications (including video), an application flow itself is also characterized by its minimum packet interarrival time, burstiness, multiple QoS levels, bandwidth, loss rate, delay, and jitter requirements.

On the other hand, the inter-characterization deals with relative importance among different applications, rendering their priority order. In general, it is more “beneficial” to give more important application flows priority over less important ones in scheduling their data transmission or allocating bandwidth.

### B. Channel Characteristics

The number and condition of channels between the proxy and  $MC^2$  can change with time due to many factors including inter-channel interference, and communication failures due to  $MC^2$  members’ departures, device mobility, or power depletion. While a proxy must be continuously aware of channel conditions, it does not have the benefit of observing packet reception or  $MC^2$  member behavior directly. The proper design of a monitoring system providing such feedback is crucial to achieving superior system performance, and we have investigated it in detail in a separate paper [38].

We studied several approaches to channel monitoring in Section II-E. In this paper we assume that our bandwidth aggregation system has a two-sided channel monitor (i.e., one side on the  $MC^2$  and the other side at the proxy) that is jointly responsible for detecting membership changes, “sensing” channel characteristics (e.g., bandwidth, error rate, latency, security, reliability, cost, etc.) and ensuring that the proxy has reasonably current channel information. Further, to facilitate WAN channel monitoring and overall sys-

tem reliability and responsiveness, the  $MC^2$ -based agents are distributed, operating on many or all of the member devices [38].

The proxy is thus capable of ordering channels in its resource pool according to the application requirements of arriving flows. For example, channels can be sorted according to their delay and reliability characteristics, and then the proxy may choose the  $n$  most reliable channels for transporting the base layer (or essential part) of a video flow while choosing less reliable channels for the enhancement layer.

### C. Allocation/Reallocation of Channels

Each application flow  $f_i$ ;  $1 \leq i \leq k$ , is assumed to have been demultiplexed into an ordered (according to the application characteristics) set of  $n_{f_i} \geq 1$  subflows  $\{sf_j : j = 1, \dots, n_{f_i}\}$ . The traffic of each subflow  $sf_j$  is represented by either a simple token bucket model  $(\rho_j, \sigma_j)$  or a linear bounded arrival process  $(p_j, s_j^{max}, b_j^{max})$  [11], where

$\rho_j$ : average token drain rate,

$\sigma_j$ : bucket size,

$p_j$ : minimum or average time separation between two consecutive packets,

$s_j^{max}$ : maximum packet size (in bytes),

$b_j^{max}$ : maximum burst size (in bytes) for subflow  $j$ .

Let  $C = \{ch_\ell : \ell = 1, \dots, n_c\}$  be an ordered (according to their condition) set of channels available. Note that the size and ordering of this set changes with time and will be updated by the monitor. The problem is now to select one or more channels from  $C$  on which to assign each subflow  $j$ . This selection must also be adapted to reflect the changing number and condition of available channels.

We first treat the simple case of only one application flow between a proxy and a  $MC^2$ , and then the more general case of multiple application flows.

1) *Case I: Single Application Flow:* We want to map a demultiplexed application flow  $f_i = \{sf_j^i : j = 1, \dots, n_{f_i}\}$  to a changing set of channels  $C = \{ch_\ell : \ell = 1, \dots, n_c\}$ . Recall that the subflows of  $f_i$  are ordered according to their importance to the application, while the channels are ordered according to their relevance to the application requirements. For example,  $f_v = \{sf_1^v, sf_2^v\}$  and  $C = \{ch_1, ch_2, ch_3\}$  where  $sf_1^v$  and  $sf_2^v$  represent the base and enhancement layers of a video stream  $f_v$ , respectively, and  $ch_i$ 's are ordered according to their reliability or their signal-to-noise ratio values. In this case  $sf_1^v$  may be transported via  $ch_1$  and  $ch_2$ , and  $sf_2^v$  via  $ch_3$ , assuming that the former requires two channels while the latter requires only one channel.

In general, as many topmost (say,  $k$ ) channels as necessary for transporting  $sf_1^i$  are assigned first to  $sf_1^i$ , and then repeat the same procedure with the remaining channels for  $sf_2^i$ , and so on. If  $sf_1^i$  does not need

the entire bandwidth of channel  $ch_k$ , the remaining bandwidth of this channel is assigned to  $sf_2^i$ , and  $ch_k$  will transmit the packets of  $sf_1^i$  and  $sf_2^i$  using a Weighted Round-Robin (WRR) scheduling algorithm where the weights between the two subflows are determined based on the  $ch_k$ 's bandwidth assigned to  $sf_1^i$  and  $sf_2^i$ . Also, if there is not enough bandwidth available, the least important subflows are not transported at all, realizing a form of imprecise computation [22].

The actual number of channels to be allocated for each subflow are determined by the subflow's requirements of delivery delay or bandwidth. For example, one can compute the bandwidth and delay requirements of both the base and the enhancement layers for layered videos, and derive the effective bandwidth of each channel from its raw bandwidth and loss-rate information.

2) *Case II: Multiple Application Flows:* In the case where there are multiple application flows  $f_i$ ,  $i = 1, \dots, n_f$ , the channels between the proxy and the  $MC^2$  must be shared among these flows according to the relative importance of the flows and the channel condition. We now order applications flows according to their relative importance, and allocate channels to the application flows, exactly in the same way as the channels are allocated to the subflows in the previous subsection (i.e., all subflows of the most important application flow are allocated to channels first, then the subflows of the next important application, and so forth). Multiple application flows of the same importance share some channels using WRR where the weights are assigned according to their bandwidth requirements.

If (weighted) fairness is used instead of importance, or if all channels are of the same quality, one can use a WRR scheduling algorithm to "serve" the different application flows. If multiple flows are multiplexed on a channel, packet transmissions of the multiplexed flows can be scheduled using either WRR or Weighted Fair Queueing (WFQ) to reflect the difference in the flows' bandwidth requirements, or Rate-Monotonic (RM) or deadline scheduling [21] for delay guarantees.

#### D. Example: Assignment of Video Flows

The potential benefit of application-aware channel assignment is best illustrated by considering the case of video traffic. First, a high-quality video flow might be of sufficiently high bandwidth that it could not be transmitted over a single WAN channel. Second, link transmission characteristics can directly affect the perceived quality of the transmission. We present three 'strawman' algorithms, based on simple heuristics, for striping video packets.

- **Layer-Priority Striping (LPS):** This algorithm can be used for video streams that are hierarchically layer-coded [27, 45]. This encoding process generates a base layer  $\ell_0$  containing information required for decoding, and one or more optional enhancement layers ( $\ell_i : i = 1, \dots, n$ ) in a hierarchical structure of cumulative layers. The reconstruction is progressive (i.e., enhancement layer  $\ell_k$  can only be used if all sublayers  $\ell_i : i = 0, \dots, k - 1$  are available). Thus, the layer index  $i$  corresponds to the layer priority.

The LPS algorithm matches the layer-priority to the channel reliability as described in Section III-C.1. For instance, the base layer  $\ell_0$  is assigned to the most reliable channels, where the channel loss rate is used as the metric for reliability. The packets for each layer are striped in WRR fashion onto the allocated channels. If a new channel with higher reliability becomes available, allocation of layers is shifted up to channels with higher reliability. Similarly, if the channel with the highest reliability becomes unavailable, the allocation is shifted down.

- **Frame-Priority Striping (FPS):** This algorithm can be used for MPEG video traffic [18]. The MPEG video stream is separated into three subflows ( $sf_I, sf_P, sf_B$ ) based on frame types. The priority order for the frames in MPEG Group of Pictures (GoP) is I>P>B. Similar to the LPS algorithm, the channels are allocated according to the subflow priority. The I-frame subflow ( $sf_I$ ) is sent over the most reliable channels, and so on.
- **Independent-Path Striping (IPS):** This algorithm is well suited to multiple state video coding [2, 46], where a stream is encoded into multiple *independently* decodeable subflows. Moreover, information from one subflow can be used to correct the errors in another subflow. Hence, it is important for a receiver to successfully receive as many complete subflows or components as possible, and it is desirable to achieve a low correlation of loss across different subflows.

The IPS algorithm tries to achieve path diversity by allocating a separate channel for each description. Since the video can be reconstructed (albeit at lower quality) even if one or more entire subflows are lost, video reception is protected against one or more complete channel failure(s).

We will later show using simulation that even these simple algorithms based on heuristics can improve video quality significantly in realistic settings.

#### IV. ARCHITECTURE

Considering the many systems issues identified in Section II, we chose a channel-aggregation architecture that is both simple and scalable. Figure 1 shows the proposed architecture which permits deployment by various types of network transport and service providers, including content owners, Internet access providers, wireless telecommunication service providers, or content distribution network operators.

The system architecture has three principal components: a dedicated appliance providing channel-aggregation proxy services, standard LAN-based announcement and discovery protocols, and standard protocol tunnels. The dedicated aggregation proxy performs inverse multiplexing at the application layer.

Generic Routing Encapsulation (GRE) [13] tunnels are used to create channels between the proxy and participating  $MC^2$  members, and support packet forwarding. This approach requires no modification to  $MC^2$  members, as most contemporary OSs (e.g., Linux, FreeBSD, Windows) have built-in support for GRE

tunnels. Each packet received by a  $MC^2$  member over a GRE tunnel is automatically decapsulated and forwarded via the wireless LAN to the destination device. Since the destination is oblivious to which  $MC^2$  node forwarded the data packets, no additional data reassembly functionality is required at the receiver.

To participate in  $MC^2$  formation and channel aggregation, a standard announcement and discovery protocol is required on end-devices. The choice of a standard protocol enables end-devices to participate in other types of resource or service discovery and access. Though the specifics of these protocols are beyond the scope of this paper, Jini, Universal Plug and Play (UPnP), and the Service Location Protocol (SLP) [15] may all be suitable candidates. Besides the basic join-leave protocol, we have designed a distributed monitoring architecture. The proposed architecture is decentralized; every community member participates in monitoring. Each member has a monitoring agent which joins a well-known multicast group  $G_m$  for exchanging community status information. Each monitoring agent broadcasts a *local report*  $R_l$  addressed to  $G_m$  on the LAN. Each local report contains information about the current state of the member and its offered WAN link(s).

Upon receiving a local report from member  $m_i$ , each member updates the information about member  $m_i$  in its locally-maintained community status database. In steady-state each member has up-to-date information about all community members. Each member issues a single packet containing the local report once every local reporting interval  $I_l$ . The collective information about the community members is sent to the inverse multiplexing proxy in *proxy reports*  $R_p$ . The community reports its current state to the proxy once every proxy reporting interval  $I_p$ . Instead of electing a particular member to send proxy reports, every member shares the responsibility. Each member sets a suppression timer for duration of  $I_p + \delta$ , where  $\delta$  is a uniform random variable on  $[0, S_d]$ . Upon expiration of its suppression timer, member  $m_i$  sends  $R_p$  to the proxy via its local WAN link, and also multicasts the same report to the community on group  $G_m$ . Upon receipt of the multicast proxy report the members other than  $m_i$  cancel their suppression timers and report transmissions. The details of the monitoring architecture and simulation results are discussed in detailed in [38].

The performance gains that our channel aggregation can realize will be explored through simulation and implementation in Sections V and VI. These benefits come at the expense of some computing and communication overhead. Note, for example, that it will not be possible in general to have a proxy on the shortest path between a source and a destination. Clearly, both an application-layer proxy as well as tunneled channels incur packet-processing overhead.

## V. PERFORMANCE EVALUATION: SIMULATION

We evaluated the proposed bandwidth aggregation system using the *ns-2* [31] simulator. Figure 2 shows the network topology we used for simulating an entire end-to-end system. The number of  $MC^2$  members was varied from 2 to 14, and the  $MC^2$  members were interconnected via an 11 Mb/s wireless LAN. In our

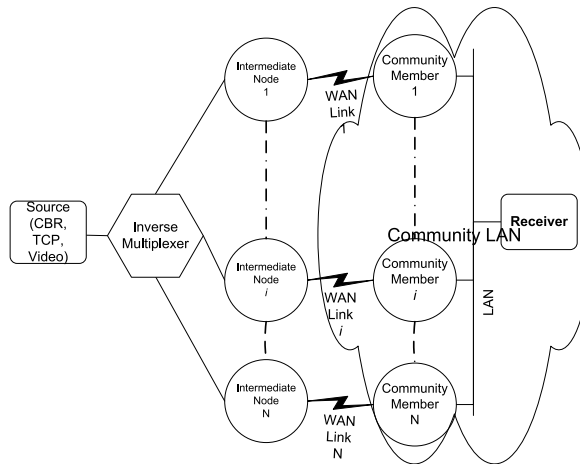


Fig. 2. Simulation topology.

experiments with homogeneous WAN links, the link bandwidth was set at 115.2 kb/s, roughly consistent with currently-available 2.5G cellular services. With the exception of the single dedicated receiver, each  $MC^2$  member was equipped with both a WAN and a LAN interface. The receiver could communicate upstream only using one of the other members as a gateway. We consider a variety of scenarios with varying link characteristics such as bandwidth, loss, and membership dynamics. We first evaluate the benefits of bandwidth aggregation for different applications: we use (1) bulk file transfer over TCP and measure TCP throughput, and (2) CBR traffic over UDP and measure packet loss rate and delay jitter. We then study how much performance improvements application-aware striping can make using layered video as an example application. We also evaluate the importance of efficient channel monitoring system by studying the impact of feedback latency on application-aware striping. For experiments with TCP and UDP traffic we implemented three application-agnostic striping algorithms: random, round-robin (RR), and weighted round-robin (WRR).<sup>2</sup> We implemented the LPS algorithm described in Section III-D for application-aware, channel-adaptive striping algorithms.

### A. TCP Throughput

We first evaluate the effect of the addition or deletion of a WAN link in an aggregated channel on TCP throughput. Let's consider the simple case of a fixed membership  $MC^2$ . We measured TCP throughput by transferring a 1 MB file from a data source to a  $MC^2$  receiver using  $2 \sim 14$  identically-configured links aggregated into the shared pool. To provide a baseline for measured TCP throughput, we also performed the experiment with a single channel (i.e., no aggregation).

Figure 3 plots the measured TCP throughput as the  $MC^2$  size changes. The average throughput achieved with a single link was 103.2 kb/s. As expected, the TCP throughput increases nearly linearly as the number

<sup>2</sup>To be precise, since packets are not fragmented in the proxy we have implemented the Surplus Round Robin approximation of bit-WRR.



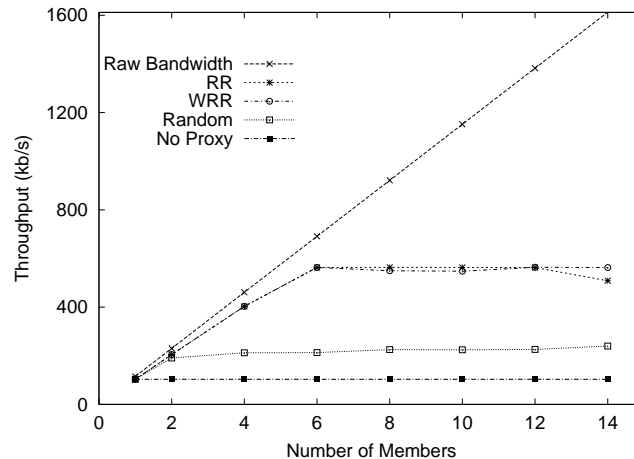


Fig. 3. TCP throughput as a function of  $MC^2$  size.

of links grows under both RR and WRR policies until saturation occurs with six links. This saturation occurs due to the limit imposed by the receiver's maximum window. As the number of available channels increases, the bandwidth-delay product increases, but TCP cannot utilize all the available bandwidth because of the small receiver window. The TCP throughput continues to increase linearly if the receiver-advertised window is increased to accommodate a larger bandwidth-delay product. The random policy does not perform as well as (W)RR because it causes undesired side effects, such as packet reordering and unstable RTT calculation, thus reducing the TCP throughput.

We next explore TCP performance for the highly-dynamic case where the channels were frequently added or removed from the pool. It is difficult to predict the likely rates of joins and leaves in a  $MC^2$ , as the behavior will likely change dramatically with the actual setting (e.g., a bus or a conference room). Hence, we conducted a variety of experiments to study join and leave dynamics, repeating the file-transfer scenario described earlier. In this set of experiments there was no significant difference in the achieved throughput for RR and WRR striping. Hence it is difficult to distinguish between the two in the figures presented here.

1) *Case I: three persistent links, one transient link:* In this scenario, three links always remain active in the pool. The fourth link periodically joins the pool for *up-time* and leaves for *down-time*. The sum of *up-time* and *down-time* was kept constant at 20 seconds. That is, an *up-time* of 20 seconds is same as striping continually over four links (i.e., 100% duty cycle) and a *down-time* of 20 seconds is the same as continually striping over only three links (i.e., 0% duty cycle). Figure 4 shows that as the duty cycle increases, the TCP throughput increases for RR and WRR schemes, whereas the random striping cannot effectively utilize the available bandwidth of the transient link.

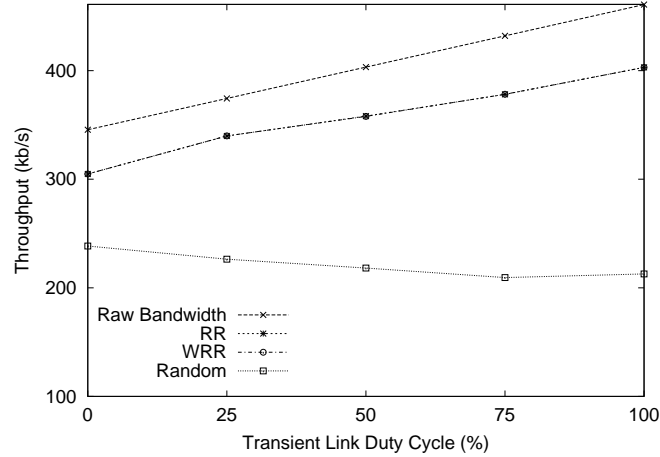


Fig. 4. TCP throughput with three persistent links and one transient link.

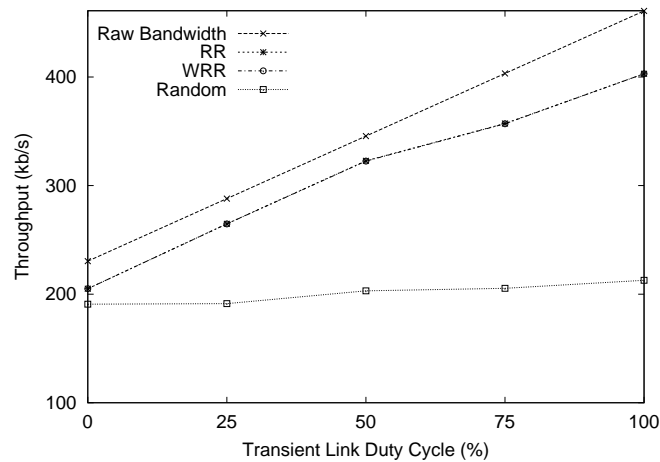


Fig. 5. TCP throughput with two persistent links and two transient links.

2) *Case II: two persistent links, two transient links*: This scenario is identical to the previous one, except that there are two links remaining active and two links being simultaneously added and removed from the pool. Figure 5 shows that as the duty cycle increases, the average TCP throughput increases for RR and WRR. Even though two of the links in the pool are rather short-lived, channel-adaptive striping is able to utilize their capacity to improve the transfer rate.

3) *Case III: one persistent link, one transient link*: In this scenario only one link is persistent and one link is periodically added and removed from the pool. We varied the length of the *up-time* interval from one second to five seconds. The duty cycle was kept constant at 50% by using the same value for *down-time* and *up-time* intervals. Figure 6 shows the TCP throughput as the interval is varied. Although the duty cycle is

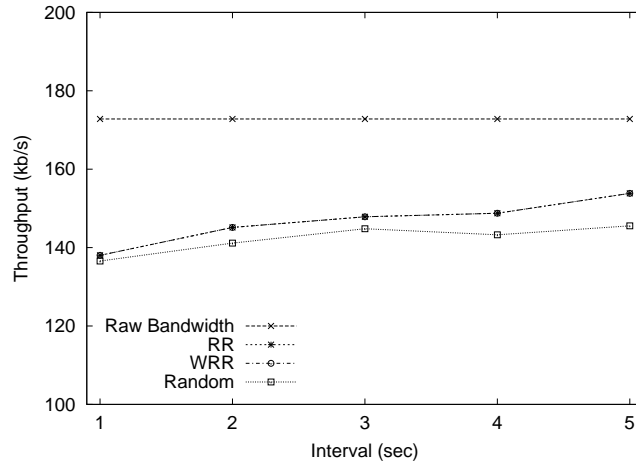


Fig. 6. TCP throughput as a function of *up-time* interval.

constant, the TCP throughput slightly increases with the length of *up-time* interval. Thus, we observe that TCP throughput varies with not only the frequency of change in the number of links, but also with the length of the change intervals.

We also measured the TCP throughput by transferring a 1 MB file over an aggregated channel consisting of four links with unequal bandwidths of 128 kb/s, 64 kb/s, 32 kb/s, and 16 kb/s. The throughput achieved by Random, RR, and WRR striping was measured at 41.2 kb/s, 44 kb/s, and 55.6 kb/s, respectively. It is interesting to note that — even for WRR — the throughput for the aggregated channel is less than the highest bandwidth of a single link. Since the proxy does not fragment packets and, instead, uses an approximation of bit-WRR, there is frequent packet misordering if the link bandwidths vary greatly. The effect of link bandwidth disparity in TCP throughput is explored in [33]. Several techniques such as weighted packet fragmentation [42], and multiple parallel TCP connections [17, 19] have been proposed. These techniques are orthogonal to our proposed architecture and can be very easily adopted to address this problem.

### B. CBR Media Traffic over UDP

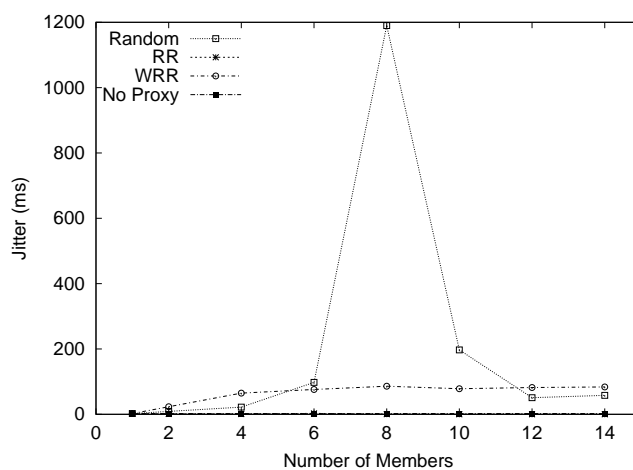
Many media applications generate CBR traffic carried over UDP. We studied the loss and jitter observed for an  $8 \times 115 \text{ kb/s} = 920 \text{ kb/s}$  CBR stream from a video source to a  $MC^2$  destination. The RTP delay jitter as described in RFC 1889 [37] was measured at the receiver. The topology used for this set of experiments was the same as the one for the TCP throughput experiments.

Table II shows the packet loss rate as a function of the  $MC^2$  size. Without channel aggregation we observe 87.6% loss as the CBR stream rate was eight times the bandwidth of a single link. As more links are pooled, the loss rate decreases. Figure 7 shows that except for random striping, the jitter values remain

TABLE II

CBR LOSS RATE (%) AS A FUNCTION OF  $MC^2$  SIZE.

# of members	Random	RR	WRR	No proxy
2	75.15	75.15	75.15	87.57
4	50.31	50.3	50.32	87.57
6	25.48	25.45	25.5	87.57
8	1.14	0.61	0.59	87.57
10 or more	0	0	0	87.57

Fig. 7. CBR jitter as a function of  $MC^2$  size.

largely unaffected by channel aggregation. With random striping, the jitter increases as the number of  $MC^2$  members increases to eight. As in the case of any switching fabric, the delay jitter increases as the offered load approaches saturation. The maximum jitter value of 425 ms was observed with eight members, i.e., when the offered CBR traffic load is equal to the sustainable throughput of the pool.

We also studied the performance of different striping algorithms for UDP streaming over four heterogeneous links of 128 kb/s, 64 kb/s, 32 kb/s, and 16 kb/s, respectively. Table III shows the loss rates when a CBR stream of 256 kb/s is sent over the aggregated channel. Random and RR algorithms do not adapt to channel bandwidth and allocate an equal number of packets to each channel. Hence, the lower bandwidth links drop larger amounts of traffic, resulting in higher total loss rates. In contrast, WRR achieves a low overall loss rate by assigning packets proportionally to the bandwidths of various links and distributing the loss uniformly over different links. A small penalty is paid through a very slight increase in jitter under the WRR algorithm as shown in Table IV, but this small increase in jitter can be easily absorbed in the receiver

TABLE III

CBR LOSS RATE (%) OVER FOUR HETEROGENEOUS LINKS.

	Random	RR	WRR
Link 1 (128 kb/s)	0	0	14.1
Link 2 (64 kb/s)	6.93	7.9	13.75
Link 3 (32 kb/s)	53.67	53.95	13.06
Link 4 (16 kb/s)	77.15	76.97	11.54
Total	34.18	34.4	13.25

TABLE IV

CBR JITTER (MS) OVER FOUR HETEROGENEOUS LINKS.

Protocols	Random	RR	WRR
Jitter	7.24	2.45	13.25

buffer.

We also evaluated how CBR streaming over UDP is affected by the dynamics of  $MC^2$  membership. Under the same join and leave dynamics as for the TCP throughput experiments, the loss rate decreased with the increase in duty cycle.

### C. Application-Aware Striping

We now present the results from the application-aware striping experiments. We experimented with the application-aware, channel-adaptive *LPS* algorithm introduced in Section III-D. The scenarios were so chosen as to elucidate the key benefits of application-aware mechanisms in comparison with application-agnostic schemes.

1) *Availability of Extra Channels*: Let's consider a scenario where the proxy has ten channels available for striping data. All the channels are identical except for having different error rates that vary from 1 to 10%. The error rate  $e_i$  for channel  $ch_i$  was set at  $i\%$ . The traffic source generated CBR traffic at 30 kb/s and the bandwidth of each channel was 20 kb/s. Thus, at least two channels are required for the transfer. Table V shows the average loss rates for the different striping algorithms. If the proxy is unaware of the application profile/requirements, then it will use all the available channels indiscriminately. Hence, the observed loss rate is higher for the application-agnostic striping algorithms. But the proxy using an application-aware algorithm achieves better performance by striping data over only the two channels with minimum loss. Hence, even minimal information, such as the bandwidth requirements of the application, can make a significant

TABLE V

LOSS RATE (%) WITH EXTRA AVAILABLE CHANNELS.

Protocols	Random	RR	WRR	LPS
Loss rate	5.56	5.58	5.68	1.41

TABLE VI

LOSS RATE (%) FOR LAYERED VIDEO WITH STATIC CHANNELS.

	Application-agnostic			Application-aware
	Random	RR	WRR	LPS
Layer $\ell_0$	5.07	9.97	6.05	1
Layer $\ell_1$	5.28	1.02	4.89	4.96
Layer $\ell_2$	5.53	4.81	5.16	9.72

improvement in the system performance.

2) *Priority-awareness*: As we discussed in Section III, different packets in an application flow can have higher priority than others, such as base layer or I-frame packets. We now present the results for striping a hierarchically-layered video stream with a base layer  $\ell_0$  and two enhancement layers  $\ell_1$  and  $\ell_2$ . Each layer was modeled as a 15 kb/s CBR stream. The topology consists of three  $MC^2$  members, each with a 20 kb/s WAN link. The error rate on the channels was 1, 5 and 10%, respectively. Table VI shows the percentage loss rate suffered by each layer. As expected, the random striping indiscriminately distributes the loss over all the layers. Since all the layers are constant bit-rate with equal bandwidth and the number of channels is same as the number of layers, the RR algorithm stripes all the packets from one layer to one channel. Instead of the loss being spread over all the layers equally, the layer sent over the most unreliable link suffers the most loss. The loss rate for the base layer is significantly less with the LPS algorithm. LPS uses priority-awareness to assign the base layer to the most reliable link, and the highest enhancement layer to the link with the highest error rate.

The striping algorithms utilize application-awareness to intelligently drop lower-priority subflows when an insufficient amount of resource is available. To demonstrate this benefit of application, we simulated a scenario with two  $MC^2$  members connected to the Internet via 20 kb/s WAN link. The error rate of the channels was 1 and 5%, respectively. Note that the offered traffic rate exceeds the aggregated channel bandwidth. Table VII shows the loss experienced by different layers while streaming the same video traffic as described above. Since the two available channels cannot handle the offered load of all the three video layers, the LPS algorithm drops the layer  $\ell_2$  entirely, improving the loss suffered by the base layer  $\ell_0$  and the

TABLE VII

LOSS RATE (%) FOR LAYERED VIDEO IN LIMITED STATIC CHANNELS.

	Application-agnostic			Application-aware
	Random	RR	WRR	LPS
Layer $\ell_0$	18.99	23.57	18.76	0.96
Layer $\ell_1$	19.64	12.44	20.53	5.15
Layer $\ell_2$	19.89	22.4	19.25	100

TABLE VIII

LOSS RATE (%) FOR LAYERED VIDEO IN DYNAMIC CHANNELS.

	Application-agnostic			Application-aware
	Random	RR	WRR	LPS
Layer $\ell_0$	3.87	4.09	4.04	0.91
Layer $\ell_1$	3.99	3.93	4.18	1.08
Layer $\ell_2$	4	4.24	3.97	10.11

enhancement layer  $\ell_1$ . The application-agnostic mechanisms end up spreading the loss over all the layers.

3) *Dynamic Channel Adaptation*: What happens if in the above scenarios the link error rates change dynamically? Let us assume that each link has an error rate of 1% for 100 seconds and then 10% for 50 seconds, repeating this cycle several times during the lifetime of the flow. The changes in error rates are distributed such that at any instant two links have error rate of 1% and one link has error rate of 10%. Thus, the total of the error rates of all the links is the same throughout the experiment. Table VIII shows the measured loss rates for this experiment. Once again, in the case of application-agnostic schemes, lack of application knowledge leads to uniform loss rates for all the layers of the flow. In contrast, LPS is able to protect the base layer from loss, and instead increase the loss rate of enhancement layers.

We also simulated the limited channel scenario described earlier with varying channel error rates. At any instant one link experiences 1% error rate and the other 10%. Table IX shows the measured loss for each layer. In this case too, with random, RR and WRR striping, all the layers suffer similar loss rates. As before, LPS entirely drops the enhancement layer  $\ell_2$  due to limited channel availability, to shield layers  $\ell_0$  and  $\ell_1$  from loss. Also, it remaps the base layer to the more reliable channel as the channel error rates change. Hence, the loss suffered by the base layer is lower.

TABLE IX

LOSS RATE (%) FOR LAYERED VIDEO IN LIMITED DYNAMIC CHANNELS.

	Application-agnostic			Application-aware
	Random	RR	WRR	LPS
Layer $l_0$	18.88	22.9	18.9	1.01
Layer $l_1$	19.73	12.78	20.75	4.97
Layer $l_2$	19.96	22.76	18.88	100

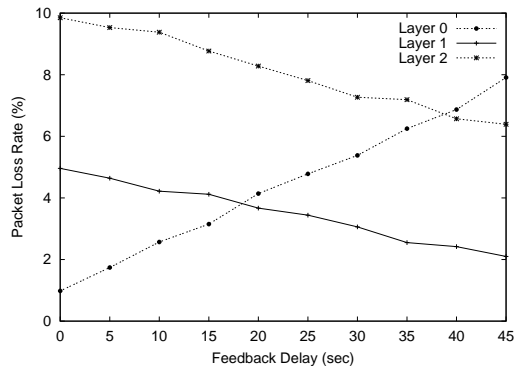


Fig. 8. The effect of reporting latency on aggregation performance.

#### D. Channel Monitoring Feedback Latency

To illustrate the importance of low latency in reporting WAN channel status to the aggregation proxy in improving the performance of an aggregated channel, we simulated an aggregation system with three community members. Each member offered a WAN channel with 20 kb/s bandwidth. Each channel has a time-varying packet loss rate (unknown to the proxy) that cycles as follows: a loss rate of 1% for 50 seconds, followed by a loss rate of 5% for 50 seconds, and then a loss rate of 10% for 50 seconds. The cycle is repeated multiple times during the lifetime of the session. The changes in loss rates across the three links are synchronized such that at any instant there is exactly one channel that has error rate of 1%, one channel with 5% and one channel with 10%. Thus, the total of the error rates of all the channels is the same throughout the experiment.

Similar to the earlier experiments, the simulated layered video consists of base layer ( $l_0$ ) and two enhancement layers ( $l_1$  and  $l_2$ ). Each layer is modeled as a 20 kb/s CBR stream. Using the channel loss rate as the reliability metric, the aggregation proxy uses the *Layer Priority Striping* (LPS) algorithm to map each layer onto one of the three channels, with higher layers assigned to increasingly less reliable channels. Figure 8 shows the packet loss rate of each layer when the reporting latency (i.e., feedback delay) is varied. The feedback delay is defined as the time difference between the instant when the channel error rate changes and the time when the aggregation proxy remaps the flows onto the channels based on the newly-available



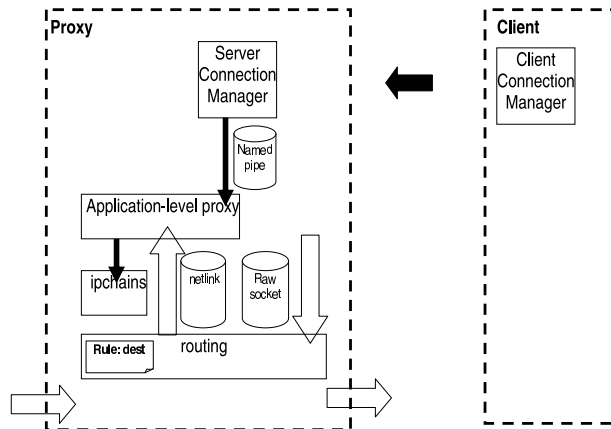


Fig. 9. Linux-based implementation of an aggregation proxy.

information. As expected, the feedback delay decreases aggregated channel performance; the base layer is not transmitted over the most reliable channel during the feedback delay period following each loss rate transition event. In fact, when the feedback latency is larger than 18 seconds, the loss rate of the base layer exceeds that of enhancement layer  $l_1$ .

## VI. IMPLEMENTATION, EXPERIMENTS AND RESULTS

We now present a detailed description of the channel aggregation testbed we built and the experiments we performed. The principal goals of the testbed were to validate our proposed architecture, corroborate our simulation results, and explore deployment issues that might not readily emerge from our simulations.

### A. Testbed Implementation

Figure 9 shows a block diagram of the prototype channel aggregation system we constructed, with dark arrows representing control messages and light arrows representing data traffic. Each  $MC^2$  member runs a compact *Client Connection Manager* (CCM) application. The CCM participates in the announcement and discovery of  $MC^2$  members (and their associated WAN links). Though we anticipate that standard announcement and discovery protocols would be used in an actual system, resource discovery was done manually in our testbed. This gave us precise control over  $MC^2$  membership, facilitated automated testing, and allowed us to modify resource availability on very short time scales.

The CCM communicates the addition or deletion of links to the *Server Connection Manager* (SCM) which resides on the proxy and maintains the channel resource pool. The CCM also monitors link transmission characteristics such as bandwidth and delay that is provided to the striping proxy. This information can be

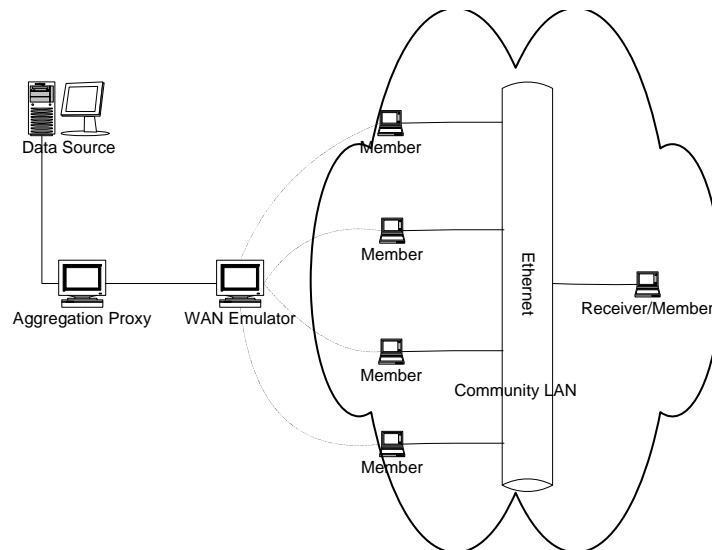


Fig. 10. Experimentation testbed topology.

used as input parameters to the striping algorithm. The CCM can also request the striping algorithm to be used for an aggregated channel. The SCM and CCM together also coordinate setup and teardown of the GRE tunnels [13] between the proxy and the  $MC^2$  members.

We implemented a Linux-based inverse multiplexing proxy. The proxy intercepts each packet destined for a  $MC^2$  and forwards it to the GRE tunnels corresponding to each active channel. Packet interception at the proxy is handled by *Netfilter* [29], a packet filtering subsystem in Linux that is primarily used for building firewalls and NATs. For each channel aggregate, the proxy sets up *Netfilter*'s forwarding rules to intercept appropriate data traffic and passes it to the proxy's user-layer forwarding engine. The forwarding engine currently implements both random and round-robin data striping policies. Use of the IP address of a  $MC^2$  member's WAN interface to set up the tunnel ensures that each packet sent over the GRE tunnel traverses the desired WAN channel.

Data reassembly at the receiving side is automatic and straightforward. Packet forwarding is enabled at each  $MC^2$  node sharing a WAN link. When a packet is received by a node over a GRE tunnel, it is decapsulated and passed to the node's routing engine. Since the destination address of the decapsulated packet corresponds to the receiver's LAN address, the packet is forwarded to the LAN.

Figure 10 shows the topology of the testbed we used for emulating an entire end-to-end system. The membership of the  $MC^2$  in our experiments varied from two to five notebook computers running Linux (2.2.16 kernel), each with built-in support for GRE tunnels. We selected relatively low-performance systems with an eye toward ultimately supporting even lower performing handheld PDAs. Forwarding was enabled on each  $MC^2$  node. Our proxy was implemented on a Linux-based desktop PC.

The  $MC^2$  members were connected to each other via a 10 Mb/s Ethernet. WAN links were emulated by

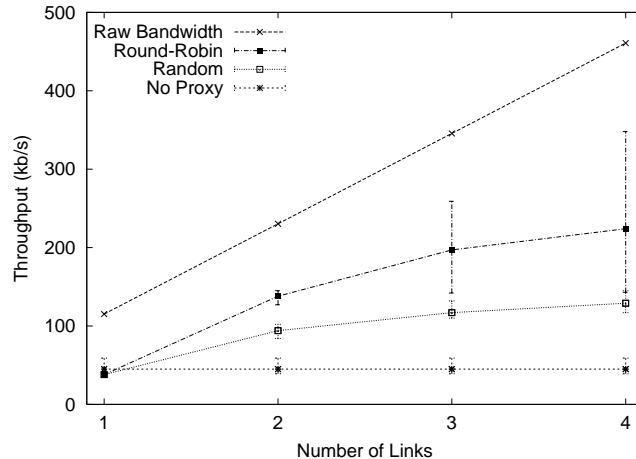


Fig. 11. Effect of  $MC^2$  size on TCP throughput.

connecting a wired serial null modem running PPP to the *NISTnet* [28] network emulator whose transmission link characteristics we could control. As in simulations presented in Section V, the transmission speed of each serial link was set at 115.2 kb/s. Each  $MC^2$  member, with the exception of the dedicated data receiver, had both an emulated WAN interface and an Ethernet interface. The data receiver could only communicate upstream to the Internet using one of the other members as a gateway.

Traffic generation, collection and measurement was performed using NetIQ's *Chariot* network measurement tool version 4.2 [30]. *Chariot end-points* running on the data source and receiver generated various packet flows, emulating reliable data transfers, streams, etc.

## B. Experimental Results

1) *TCP Throughput*: To validate our simulation results in practice, we measured TCP throughput by transferring a 1MB file from a data source to a  $MC^2$  receiver using two to four identically-configured, aggregated links. To provide a baseline for measured TCP throughput we also performed the experiment with a single channel (i.e., no aggregation) both with and without the proxy in the data path. Performance was measured using both round-robin and random striping policies. Figure 11 plots the measured TCP throughput as the number of links in the aggregate pool changes, with error bars showing the minimum and maximum measured throughput among the 50 trials.

The average TCP throughput achieved with no proxy was 45 kb/s. The TCP throughput with a single link and the proxy in the data path is 38 kb/s, not significantly lower than the throughput achieved without a proxy, indicating that the proxy does not introduce a long delay. The TCP throughput measured in the testbed was lower than the simulation results due to PPP overhead and the presence of background traffic.

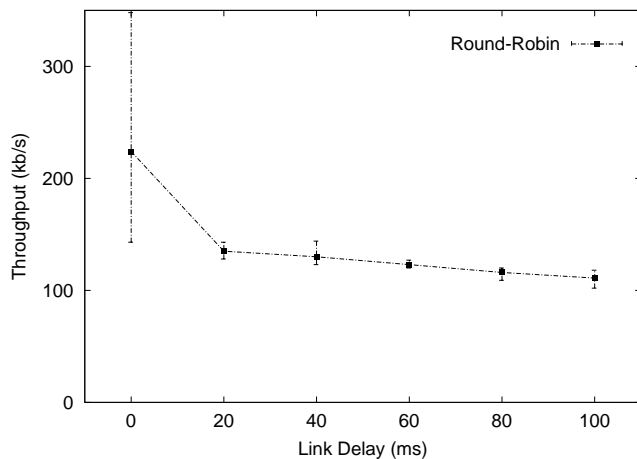


Fig. 12. Effect of link latency variation on TCP throughput.

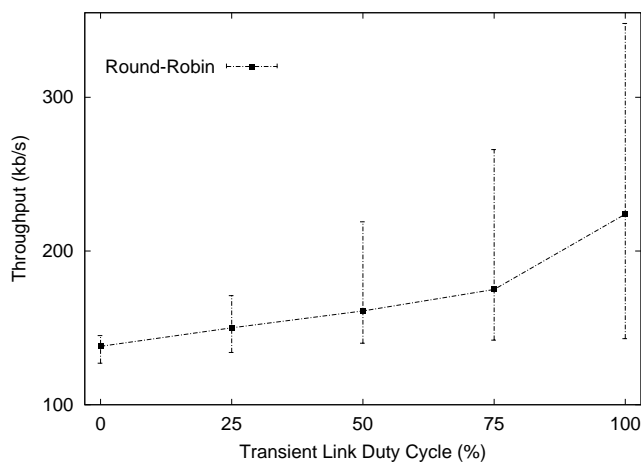


Fig. 13. TCP throughput (kb/s) with two transient links.

However, the trends with respect to the number of  $MC^2$  members were similar in both the cases.

To study the effect of varying transmission link parameters of different WAN channels on TCP throughput, we used the *NISTnet* emulator to add extra transmission delay to one of the WAN channels. Figure 12 shows the change in TCP throughput as the extra delay of one of the four links is varied from 0 to 100 ms. As expected, increasing the link delay decreases throughput. There are two reasons for this. First, the increased delay can cause additional packet misordering, introducing reassembly delays. Second, the extra delay results in a larger computed value for RTT, directly decreasing throughput. With unequal link latencies, round-robin is not the optimal scheduling algorithm. WFQ techniques such as that proposed in [42] can reduce packet misordering and hence improve TCP throughput.

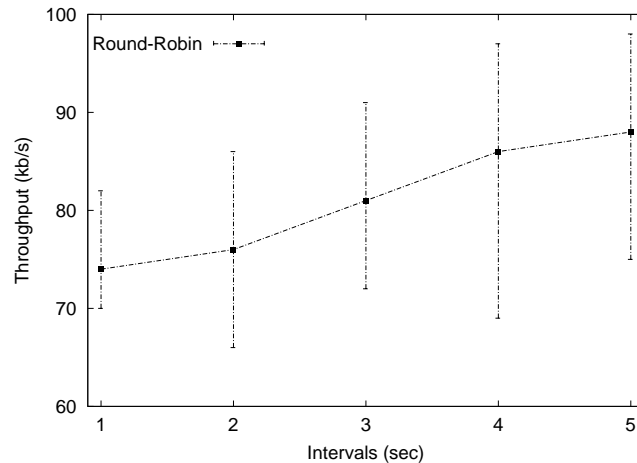


Fig. 14. Effect of *up-time* interval length on TCP throughput.

We now measure the TCP throughput in a highly dynamic  $MC^2$  where channels are added and removed from the resource pool. The topology is the same as we used in Section V-A. Figure 13 shows the TCP throughput as the duty cycle was changed for two transient links among four total links. We observe that as the duty cycle increases, the average TCP throughput increases in most cases. Although some of the links in the pool are rather short lived, link aggregation is nonetheless able to use their capacity to improve the transfer rate.

Figure 14 shows the TCP throughput when the length of *up-time* interval is changed for one link while the other link is persistent. The result verifies our earlier observation from simulation (Figure 6) that the interval duration as well as the frequency of change in the number of active channels affect TCP throughput.

2) *Streaming Media via UDP*: We next conducted experiments to study how high bandwidth streaming is enabled with channel aggregation. In these experiments a server streams a stored media file to the receiver at one of various bit rates (64, 128, 175, and 256 kb/s). Chariot generates a traffic pattern intended to resemble the video transmission of Cisco’s IP-TV. RTP [37] is used as the stream transport protocol. Each experiment was repeated 25 times, measuring the loss rate and RTP delay jitter observed by the receiver.

Without channel aggregation the receiver can only receive a stream with negligible loss at the 64 kb/s rate. Higher bit-rate streams suffered more than 70% loss, and due to this high loss rate, the tests were prematurely terminated by Chariot. Note that the limited available bandwidth precludes use of retransmission for loss recovery. Techniques such as Forward Error Correction (FEC) cannot be used in this setting, especially for low-bandwidth links, as it further increases the bandwidth required. Such a high loss rate can severely degrade the perceived stream reception quality, making it unwatchable. Yet striping over just two links reduced the loss rate dramatically for the 128 kb/s stream; every 128 kb/s stream test completed with a loss

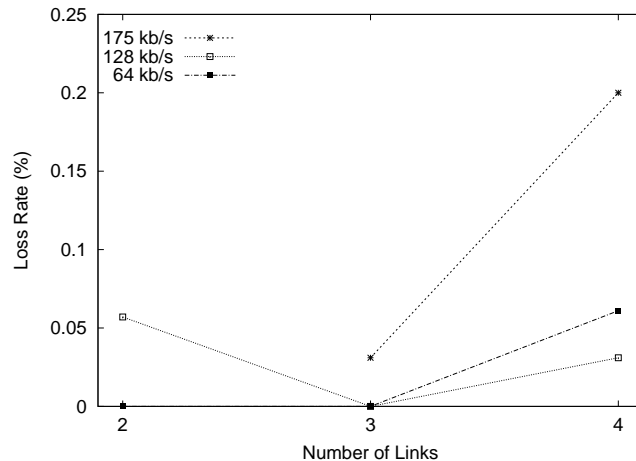


Fig. 15. Effect of  $MC^2$  size on RTP stream loss rate.

rate of less than 0.1%. The 175 kb/s streaming experiment with striping over two links was also terminated before completion due to high loss rate. Striping over four links was capable of sustaining a 256 kb/s stream without significant packet loss.

Figure 15 shows the streaming data loss rates. Observe that the data loss rate does not exceed 0.2% when dynamic link striping is performed. This result confirms that bandwidth aggregation enables high bandwidth multimedia streams to be delivered to  $MC^2$ s, which would otherwise be impossible.

Figure 16 shows RTP jitter values. Note that the system generates relatively little jitter. In most cases, the jitter is less than 10 ms with the maximum jitter occasionally exceeding 20 ms. Such small amounts of jitter can be easily absorbed by the receiver buffer in multimedia applications and will have negligible effect on the viewing experience of the video receiver.

## VII. RELATED WORK

A combination of a high demand for communication bandwidth and high tariffs on WAN links has long made inverse multiplexing a popular technique [14]. In the early 1990s the Bandwidth on Demand Interoperability Group (BONDING) created a standard inverse multiplexing mechanism to achieve a virtual high capacity WAN link using  $n \times 56$  (or 64) kb/s links [6]. Network equipment providers supported inverse multiplexing for various link-layer technologies such as frame relay, ISDN, and SMDS. The same technique was later applied and standardized within the context of ATM networks in the Inverse Multiplexing for ATM (IMA) specification [3]. Each of these cases assumed highly reliable, homogeneous links with constant link characteristics such as capacity, delay, and error rates. Moreover, each WAN connection being bundled together originated from, and terminated at, the same endpoints.

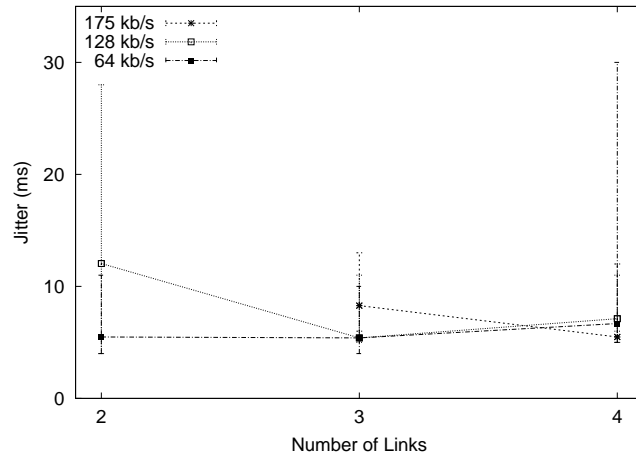


Fig. 16. Effect of  $MC^2$  size on RTP stream jitter.

Various striping algorithms have been proposed and implemented to reduce packet reordering, jitter, and load imbalance. RR scheduling is primarily used for striping data over homogeneous links, while variants of queuing algorithms are used in case of heterogeneous links. It has been shown that maximum throughput is achieved by striping data over each channel in proportion to the channel's bandwidth-delay product [1, 33].

More recent research has explored adaptive inverse multiplexing for CDPD wireless networks [42]. In this scheme the packets are split into fragments of size proportional to the observed throughput of component links. Here the goal is to create variable fragments sizes such that each fragment can be transmitted in roughly the same amount of time. The fragment size of each link is dynamically adjusted in proportion to the measured throughput. The fragmented packets are then tunneled over multiple links using Multilink PPP [41]. In this case the endpoints of the WAN connections forming the virtual link are the same.

The bandwidth of mobile users with multiple interfaces is aggregated at the transport layer in pTCP (parallel TCP) [16, 17]. pTCP is a wrapper that interacts with a modified TCP called TCP-virtual (TCP-v). A TCP-v connection is established for each interface, and pTCP manages send buffers across the TCP-v pipes. The striping is performed by pTCP and is based on congestion window size of each TCP-v connection. When congestion occurs on a certain pipe, pTCP performs data reallocation to another pipe with large congestion window. One possible problem of this approach is that the congestion window size may not accurately reflect the bandwidth-delay product.

A scheduling algorithm for aggregating bandwidth for realtime applications is detailed in [9]. The authors propose a Earliest Delivery Path First (EDPF) scheduling algorithm that is channel and application aware and minimizes the cost of striping traffic over multiple wireless channels of a device. This algorithm can also be adopted in our  $MC^2$  architecture.

Coordinating communications from *multiple* mobile computing devices has become a new focus of interest. Network connection sharing has been proposed in [32]. This architecture permits use of a single, idle WAN connection among collaborating mobile devices but it does not address aggregation of multiple links into a high capacity bundle.

Our goal of cooperation and resource aggregation among collaborating devices is similar to the vision of the mobile grouped devices (MOPED) architecture [8, 20]. The goal of MOPED project is to enable group mobility such that a user's set of personal devices appear as a single mobile entity connected to the Internet. The MOPED routing architecture builds a *multipath* layer to encapsulate packets between the home agent and MOPED devices. Unlike our approach of using GRE tunnels, the home agent and the mobile devices in MOPED must implement a new lightweight encapsulation protocol called *multipath routing encapsulation* (MRCAP). MOPED architecture provides a higher-capacity and better-quality connection to the Internet by adapting the Mobile IP home agent to support aggregation of multiple links at network and transport layers. It uses transport-layer inverse multiplexing for multi-homed devices [23]. Aggregating bandwidth at the transport layer requires different protocols for different applications. MOPED presents two new transport protocols, namely: (1) R-MTP (Reliable Multiplexing Transport Protocol) [25] for data, and (2) M-MTP (Multimedia Multiplexing Transport protocol) [24] for multimedia. Additional transport protocols might be needed as the application requirements change. Modifications to both client and server kernels are also required. Our application-level approach does not require any kernel changes and allows support for different application profiles.

The commuter Mobile Access Router (MAR) project [34] also leverages wireless WAN connection diversity to provide high speed Internet access to mobile users. Instead of using the WAN connections of the users, it relies on pre-provisioning the MAR with different WAN connections, limiting the aggregation to the already existing links.

## VIII. CONCLUSION

We have designed, implemented and evaluated a deployable bandwidth aggregation system providing high-speed Internet access to a collaborating community of wireless end-systems. We have demonstrated that the system not only improves access service quality, but enables otherwise unachievable services such as the delivery of high-bandwidth streaming media. Further, we have shown that network and application-aware allocation and assignment policies do indeed improve system performance.

Though not described in this paper, we performed various experiments with bandwidth-adaptive multimedia applications over aggregated connections. Ideally, such applications would measure available bandwidth and smoothly increase or decrease audio or video quality to optimize perceived reception quality. We typically observed an application decreasing its streaming rate to a predefined fraction of its maximum rate;



often this rate was well below the available bandwidth of the aggregated connection. The application would subsequently maintain that low rate, remaining non-responsive to any increase in available bandwidth, no matter how large it is. Since the widely-used applications we tested were proprietary, we were unable to modify their adaptation algorithms.

To aggregate bandwidth we have relied upon the conventional technique of inverse multiplexing. But rarely, if ever, has inverse multiplexing been applied in such a dynamic and decentralized setting and made to work. As a result of operating in this challenging environment, we have identified a significant number of technical issues that appear to be fertile areas for future research. Some of these include WAN cost sharing, accounting, information privacy, and security. We have also relied on the assumption that an application's networking requirements are well known, though such a characterization remains a formidable and long-standing problem. Not surprisingly, we have found that aggregating relatively homogeneous communication links is often easier and more successful than working with heterogeneous links. Opportunity lies in the development of simple 'rules-of-thumb' that help identify which links are sufficiently homogeneous that aggregation is likely to perform well.

Opportunities for community channel sharing is not limited to cellular links from spontaneously formed ad hoc groups of mobile devices but also in context of aggregating broadband access links (e.g., DSL) among neighboring residences. With the widespread adoption of 802.11 networks in home, neighboring residences can form collaborative communities to share and aggregate their broadband links. And there appears to be no other means of satisfying the growing demand for access bandwidth as quickly and as cheaply.

## REFERENCES

- [1] H. Adishesu, G. Parulkar, and G. Varghese, "A reliable and scalable striping protocol," in *Proceedings of ACM SIGCOMM*, Stanford, CA, Aug. 1996, pp. 131–141.
- [2] J. Apostolopoulos, "Error-resilient video compression via multiple state streams," in *Proceedings of VLBV*, Kyoto, Japan, Oct. 1999, pp. 168–171.
- [3] ATM Forum, "Inverse multiplexing for ATM specification, version 1.0," July 1997.
- [4] J. C. R. Bennett, C. Partridge, and N. Shectman, "Packer reordering is not pathological network behavior," *IEEE/ACM Trans. Networking*, vol. 7, no. 6, pp. 789–798, Dec. 1999.
- [5] E. Blanton and M. Allman, "On making TCP more robust to packet reordering," *ACM SIGCOMM Computer Communication Review*, vol. 32, no. 1, pp. 20–30, Jan. 2002.
- [6] BONDING Consortium, "Interoperability requirements for  $n \times 56/64$  kb/s calls, version 1.0," Sept. 1992.
- [7] S. Buchegger and J.-Y. L. Boudec, "Performance analysis of the CONFIDANT protocol (cooperation of nodes: Fairness in dynamic ad-hoc networks)," in *Proc. of ACM MobiHoc*, June 2002, pp. 226–236.
- [8] C. Carter and R. Kravets, "User device cooperating to support resource aggregation," in *Proceedings of IEEE WMSCA*, Callicoon, NY, June 2002, pp. 59–69.
- [9] K. Chebrolu and R. Rao, "Communication using multiple wireless interfaces," in *Proc. of IEEE WCNC*, Orlando, FL, Mar. 2002, pp. 327–331.
- [10] J. Crowcroft, R. Gibbens, F. Kelly, and S. Ostring, "Modelling incentives for collaboration in mobile ad hoc networks," in *Proc. of WiOpt*, Sophia-Antipolis, France, Mar. 2003.

- [11] R. L. Cruz, "A calculus for network delay, Part I: Network elements in isolation," *IEEE Trans. Inform. Theory*, vol. 37, no. 1, pp. 114–131, Jan. 1991.
- [12] J. Duncanson, "Inverse multiplexing," *IEEE Commun. Mag.*, vol. 32, no. 4, pp. 32–41, Apr. 1994.
- [13] D. Farinacci, T. Li, S. Hanks, D. Meyer, and P. Traina, "Generic routing encapsulation GRE," IETF, RFC 2784, Mar. 2000.
- [14] P. H. Fredette, "The past, present, and future of inverse multiplexing," *IEEE Commun. Mag.*, vol. 32, no. 4, pp. 42–46, Apr. 1994.
- [15] E. Guttman, C. Perkins, J. Veizades, and M. Day, "Service location protocol, version 2," IETF, RFC 2608, June 1999.
- [16] H.-Y. Hsieh and R. Sivakumar, "pTCP: An end-to-end transport layer protocol for striped connections," in *Proceedings of IEEE ICNP*, Paris, France, Nov. 2002, pp. 24–33.
- [17] —, "A transport layer approach for achieving aggregate bandwidth on multi-homed mobile hosts," in *Proceedings of ACM MobiCom*, Atlanta, GA, Sept. 2002, pp. 83–94.
- [18] ISO/IEC moving picture experts group (MPEG). [Online]. Available: <http://mpeg.telecomitalia.com>
- [19] K.-H. Kim and K. G. Shin, "Improving TCP performance over wireless networks with collaborative multi-homed mobile hosts," in *Proc. of USENIX/ACM MobiSys*, Seattle, WA, June 2005, pp. 107–120.
- [20] R. Kravets, C. Carter, and L. Magalhaes, "A cooperative approach to user mobility," *ACM SIGCOMM Computer Communication Review*, vol. 31, no. 5, pp. 57–69, Oct. 2001.
- [21] C. M. Krishna and K. G. Shin, *Real-Time Systems*. McGraw Hill, 1997.
- [22] J. W. S. Liu, K.-J. Lin, W. K. Shih, R. Bettati, and J. Chung, "Imprecise computations," *Proc. IEEE*, vol. 82, no. 1, pp. 1–12, Jan. 1991.
- [23] L. Magalhaes and R. Kravets, "End-to-end inverse multiplexing for mobile hosts," *Journal of the Brazilian Computer Society*.
- [24] —, "MMTP: Multimedia multiplexing transport protocol," in *Proceedings of ACM SIGCOMM-LA*, San Jose, Costa Rica, Apr. 2001, pp. 220–243.
- [25] —, "Transport level mechanisms for bandwidth aggregation on mobile hosts," in *Proceedings of IEEE ICNP*, Riverside, CA, Nov. 2001, pp. 165–171.
- [26] S. Marti, T. J. Giuli, K. Lai, and M. G. Baker, "Mitigating routing misbehavior in mobile ad hoc networks," in *Proc. of ACM MobiCom*, Boston, MA, Aug. 2000, pp. 255–265.
- [27] S. McCanne and M. Vetterli, "Joint source/channel coding for multicast packet video," in *Proceedings of IEEE ICIP*, Washington, DC, Oct. 1995, pp. 25–28.
- [28] National institute of standards and technology's NIST net. [Online]. Available: <http://snad.ncsl.nist.gov/itg/nistnet>
- [29] Netfilter. [Online]. Available: <http://www.netfilter.org>
- [30] NetIQ corporation's Chariot. [Online]. Available: <http://www.netiq.com/products/chr/default.asp>
- [31] ns-2, the network simulator. [Online]. Available: <http://www.isi.edu/nsnam/ns>
- [32] M. Papadopoulou and H. Schulzrinne, "Connection sharing in an ad hoc wireless network among collaborative hosts," in *Proceedings of NOSSDAV*, Florham Park, NJ, June 1999, pp. 169–185.
- [33] D. S. Phatak and T. Goff, "A novel mechanism for data streaming across multiple IP links for improving throughput and reliability in mobile environments," in *Proceedings of IEEE INFOCOM*, New York, NY, June 2002, pp. 773–781.
- [34] P. Rodriguez, R. Chakravorty, J. Chesterfield, I. Pratt, and S. Banerjee, "MAR: A commuter router infrastructure for the mobile internet," in *Proc. of ACM/USENIX MobiSys*, Boston, MA, June 2004.
- [35] K. Salem and H. Garcia-Molina, "Disk striping," in *Proceeding of IEEE ICDE*, Los Angeles, CA, Feb. 1986, pp. 336–342.
- [36] N. B. Salem, L. Buttyan, J.-P. Hubaux, and M. Jakobsson, "A charging and rewarding scheme for packet forwarding in multi-hop cellular networks," in *Proc. of ACM MobiHoc*, Annapolis, MD, June 2003.
- [37] H. Schulzrinne, S. Casner, R. Frederick, and V. Jacobson, "RTP: A transport protocol for real-time applications," IETF, RFC 1889, Jan. 1996.
- [38] P. Sharma, J. Brassil, S.-J. Lee, and K. G. Shin, "Improving aggregated channel performance through decentralized channel monitoring," *Computer Networks*, to appear.
- [39] P. Sharma, S.-J. Lee, J. Brassil, and K. G. Shin, "Handheld routers: Intelligent bandwidth aggregation for mobile collaborative communities," HP Labs, Technical Report HPL-2003-37R1, May 2003.

- [40] ———, “Handheld routers: Intelligent bandwidth aggregation for mobile collaborative communities,” in *Proc. of IEEE BROADNETS*, San Jose, CA, Oct. 2004, pp. 537–547.
- [41] K. Sklower, B. Lloyd, G. McGregor, D. Carr, and T. Coradetti, “The PPP multilink protocol (MP),” IETF, RFC 1990, Aug. 1996.
- [42] A. C. Snoeren, “Adaptive inverse multiplexing for wide area wireless networks,” in *Proceedings of IEEE GLOBECOM*, Rio de Janeiro, Brazil, Dec. 1999, pp. 1665–1672.
- [43] V. Srinivasan, P. Nuggehalli, C.-F. Chiasserini, and R. R. Rao, “An analytical approach to the study of cooperation in wireless ad hoc networks,” *IEEE Trans. Wireless Commun.*, vol. 4, no. 2, pp. 722–733, Mar. 2005.
- [44] C. B. S. Traw and J. M. Smith, “Striping within the network subsystem,” *IEEE Network*, vol. 9, no. 4, pp. 22–32, July/Aug. 1995.
- [45] M. Vishwanath and P. Chou, “An efficient algorithm for hierarchical compression of video,” in *Proceedings of IEEE ICIP*, Austin, TX, Nov. 1994, pp. 275–279.
- [46] Y. Wang, M. Orchard, V. Vaishampayan, and A. R. Reibman, “Multiple description coding using pairwise correlating transforms,” *IEEE Trans. Image Processing*, vol. 10, no. 3, pp. 351–366, Mar. 2001.
- [47] R. Woodburn and D. Mills, “A scheme for an internet encapsulation protocol: Version 1,” IETF, RFC 1241, July 1991.
- [48] S. Zhong, J. Chen, and Y. R. Yang, “Sprite: A simple, cheat-proof, credit-based system for mobile ad-hoc networks,” in *Proc. of IEEE INFOCOM*, San Francisco, CA, Apr. 2003.