

On Redundant Multipath Operating System Support for Wireless Mesh Networks

Yair Amir¹, Claudiu Danilov², Michael A. Kaplan¹, Raluca Musăloiu-Elefteri¹, Nilo Rivera¹

¹ Johns Hopkins University, Baltimore, MD.

² Boeing Phantom Works, Seattle, WA.

{yairamir, kaplan, ralucam, nrivera}@dsn.jhu.edu

claudiu.b.danilov@boeing.com

Abstract—Low-cost wireless routers are changing the way people connect to the Internet. They are also very cheap, albeit quite limited, Linux boxes. These attributes make them ideal candidates for wireless mesh routers.

This paper presents a minimally invasive mechanism for redundant multipath routing in kernel-space to achieve high reliability with high throughput in a mesh network. This service is essential for achieving fast, lossless handoff as mobile devices roam throughout the wireless mesh coverage area. However, redundant multipath is not natively supported by current operating systems, limiting the routing mechanisms that can be used in these networks to user-level implementations, which can greatly degrade performance.

We show an architecture that integrates this mechanism in a wireless mesh system, resulting in a high-throughput 802.11 mesh network with fast handoff over low-cost routers.

I. INTRODUCTION

Low-cost wireless routers are revolutionizing the way people connect to the Internet. The ease of deployment at home or office on one hand, and the good Internet connectivity they provide on the other, have made these wireless routers ubiquitous. These routers are also a revolution from another, less known perspective: They are very cheap, albeit quite limited, Linux boxes (around \$50 a piece). These attributes make them very attractive and convenient for developers to implement their own applications.

Wireless mesh networks extend the connectivity area of mobile devices beyond the limited range of a single access point. They are also very affordable when implemented with off-the-shelf low-cost wireless routers. These networks can be easily deployed inside a building, campus, or on a large geographical area without requiring every access point to be physically connected to the Internet. A lot of research over the last few years has focused on making wireless mesh networks a reality, including MIT Roofnet [10], Microsoft MCL [1], and UCSB MeshNet [15].

Redundant multipath routing (i.e., the ability to simultaneously send the same packet over multiple routes) is an essential service for increasing the reliability of wireless mesh networks. For example, as mobile clients (laptops, PDAs) roam throughout the area covered by the mesh network, their point of attachment (access point) must change to avoid loss of connectivity. Redundant multipath can help achieve uninterrupted connectivity during handoff by: (1) sending packets through multiple access points to the mobile client, to deal with unexpected client movements while the best access

point for the client is chosen, and (2) avoiding loss while route changes take place in the wireless mesh. Related work has looked into these benefits in wireless environments ([19], [13], [4]). Other applications can also benefit from redundant multipath routing (Section V). However, redundant multipath is not natively supported by current operating systems, limiting the routing mechanisms that can be used in these networks to user-level implementations.

In this paper, we propose a new approach to support redundant multipath forwarding in kernel space. When a packet is received for the first time by a node in the mesh network (referred to as mesh entry-point), the routing service encodes the mesh entry-point in the packet. Then, at every mesh node, packets are routed (possibly to multiple next-hops) according to the entry-point embedded in the packet. Our solution does not add any overhead to data packets, utilizing existing space on the IP packet header instead.

The contributions of this work are: (1) A novel mechanism for redundant multipath routing in kernel space. (2) An architecture that integrates kernel-level redundant multipath routing in a wireless mesh system, resulting in the first high-throughput 802.11 mesh with fast handoff over low-cost routers. (3) A publicly available open-source implementation of the kernel modules for Linux.

II. BACKGROUND

The routing process involves computing routes to a destination, usually taking into account the distributed and dynamic nature of the underlying network, and the actual forwarding of the packets. Packet *routing*, which specifies the rules on how packets will be forwarded, is commonly performed in user space, allowing different protocols to be easily deployable and upgradable (OSPF [18], RIP [12]). Packet *forwarding*, on the other hand, typically resides in the kernel to forward packets as fast as possible. Thus, routing relies on the forwarding capabilities provided by the operating system. This approach allows operating systems to be both flexible and efficient.

Previous work has looked into bypassing the kernel from the communication path [20]. However, it requires a sophisticated architecture, in addition to special support from the hardware. A common approach to extend the routing services provided by the kernel (without requiring kernel modifications) is to build user-level routers that forward packets at application level. RON [6] uses this approach to route packets through

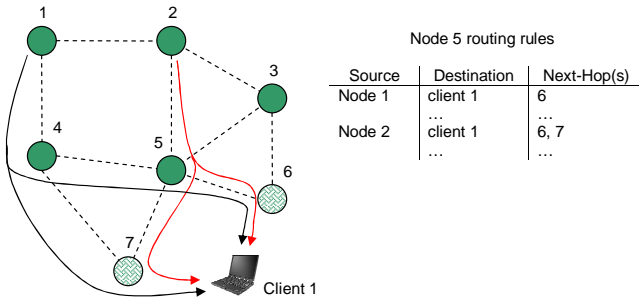


Fig. 1. The routes to a mobile client (multipath routing).

an overlay network to increase the reliability of the end-to-end path. End-System-Multicast [14] and Spines [3] also route through an application router to support overlay multicast without infrastructure support.

Other work has looked into operating system support for wireless ad-hoc routing protocols. Chakeres and Belding showed in [9] an in-kernel design and implementation of the ad-hoc AODV protocol using Netfilter modules, and showed performance improvement compared to user-level ad-hoc protocols. Kawadia et al. [16] proposed a complete architecture to support ad-hoc protocols in-kernel and a generic ad-hoc support library for user-level programs to control different ad-hoc protocols.

SMesh [4] recently showed how overlay multicast can be used in wireless mesh networks to provide fast handoff to unmodified 802.11 clients that connect transparently using DHCP. SMesh allows *multiple* access points to service the client during handoff, as it works in ad-hoc mode. In SMesh, packets sent by the mobile client are diverted from the kernel to the Spines user-level overlay router. Multicast trees are calculated in a way similar to that of MOSPF [17]. SMesh encapsulates client packets and sends them through the overlay network to the access points serving the destination. Once the packets are received by the destination's access points, SMesh strips the overlay headers and forwards the original packet to the mobile client using a raw socket.

We used the SMesh wireless mesh system as a testbed to implement and evaluate the benefits and drawbacks of the approach presented in this paper.

III. OS SUPPORT FOR REDUNDANT MULTIPATH ROUTING

User-level overlay routing allows users to implement any protocol without requiring any special support from the kernel. The SMesh system uses this approach to implement redundant multipath, essential for achieving a seamless handoff. While very convenient, routing the entire traffic through user space becomes problematic for low-cost wireless routers which have limited processing power. It is widely known that forwarding packets through user space results in higher CPU utilization when compared to kernel space. The overhead can be attributed to two primary factors: memory copies and context switches. Each routing node must copy the packets from kernel

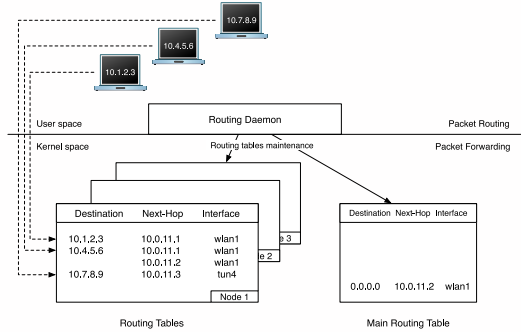


Fig. 2. Architecture.

space to user space in order to determine the next hop. After a routing decision is made, the packet must be returned to kernel space where it is sent on the network. That is, the user-kernel boundary must be crossed a minimum of two times per hop.

We describe next a mechanism that achieves efficient redundant multipath routing in kernel space. The idea behind it is simple: each node maintains multiple kernel routing tables, one for each node in the mesh network, with route entries set according to the multicast trees determined by our routing protocol.

A. Architecture

In wireless mesh networks, it is possible for packets to start flowing in the mesh from different sources. Several Internet gateways may coexist in multi-homed wireless mesh networks [5], any of which may need to forward packets to the client. Also, since clients may communicate with other clients in the mesh network, virtually every access point is capable of being the source of packets in the mesh. To provide optimal redundant multipath routing in these networks, each node must consider the mesh source and the destination of each packet in order to determine the appropriate forwarding rule for that packet. Inside the mesh network, this can be viewed as a multicast routing problem (multi-source multi-destination multicast routing)¹. Figure 1 shows how packets are forwarded to the mobile client from two different sources (node 1 and node 2)². Note that node 5 must forward the packets differently depending on the source of the packet.

Based on the description above, it is clear that a node must decide what are the next hops for a packet based on the mesh *entry-point* as well as the destination address of the packet. However, the entry-point cannot be determined by just looking at the packet destined to the client (the source address in the IP header is not the address of a mesh entry-point, but the actual address of the sender, which can be another mesh

¹Tunneling the unicast packet in an IP-multicast tunnel is not a viable solution as it suffers from lower reliability due to the lack of 802.11 link-layer retransmissions. An alternative is to use IP-multicast tunnels, with an additional *unicast tunnel* on each hop, but this approach incurs additional space in the packet as well as processing overhead on each node.

²Note that these mesh nodes are not the actual sources. Rather, they are the nodes that first received the packet in the mesh network, either from Internet or from clients connected to them.

client or an Internet address). One solution to keep track of the entry-point is to tunnel each packet from the entry-point in the mesh to the mobile client. However, we need to instruct the kernel to remove the tunnel in the last hop, right before sending the packet to the client, which requires new kernel functionality. Otherwise, the mobile client may discard these packets. Another less obvious solution is to encode the mesh entry-point in some of the existing space in the IP header of the packet. Specifically, we can encode the IP address of the entry-point into the identification field from the IP header (also referred to as IPID). This is a 16-bit field used to identify the fragments of the IP datagrams. Together with the offset field, it is used by the IP layer to reassemble the fragmented datagrams. As the packet travels in the network, the intermediate routers must leave the IPID field unchanged. To make a distinction between the original source (entry-point in the mesh) and the rest of the nodes (routers) along the path, we use a bit from *type of service* (TOS) field, specifically, the *cost* bit. This approach benefits from no overhead in the mesh, in terms of packet size, and from the fact that no modification to the packet is necessary except at the entry-point of the mesh.

Even if very convenient, modifying the IPID of the packets may create problems in the case of fragmented traffic. However, current studies show that IP packet fragmentation is not commonly used today, and it amounts to between 1 and 2% [8] of the overall traffic. While advocating for or against the use of fragmentation [11] is outside the scope of this paper, we choose to ignore the mesh entry-point when the packet is fragmented, and forward it through a single path. Other solutions to encode both the fragmentation and the mesh entry-point are possible. However, considering the small amount of the fragmented traffic, we believe this is a practical way for our system to support it.

Our approach to routing packets in the kernel is as follows: We first define a routing table in the kernel for each access point in the mesh, i.e., for each possible mesh entry-point. In each table, we add a route entry for each possible destination, i.e., for each client (Figure 2). This entry may include *several* next-hops, depending on the multicast trees determined by our routing daemon. Then, we instruct the kernel to encode in each packet the mesh entry-point when the packet is first seen in the mesh network. Each node looks for this information at each incoming packet to select what routing table to use. Then the kernel forwards the packet according to the entry that has the client address as the destination. For the example presented in Figure 1, router 5 will use different forwarding tables if the packets come from source 1 or source 2. In addition, to route packets from the clients to the Internet, each node sets in its main routing table the route to the closest Internet gateway.

B. Implementation

There are several challenges in implementing such an architecture with the current services provided by the Linux kernel networking stack. Essentially, we need to be able to forward packets simultaneously on multiple paths to the same

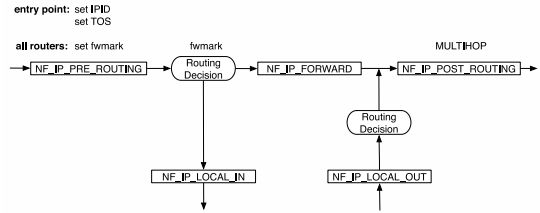


Fig. 3. Implementation in Linux.

destination. Fortunately³, since version 2.2, the Linux kernel supports defining multiple routing tables and permits policy routing (a.k.a. rule based routing), which allows selecting different routing tables based on criteria other than the destination address. In our case, each mesh node maintains a routing table for each entry-point in the network, and includes in each routing table an entry for each multicast group. That is, one routing table corresponds to all multicast trees that have that node as a source.

To alter the IPID field of the IP header, we wrote a simple Netfilter kernel module. The selection itself of which routing table to use, given the IP encoded in IPID, is done with policy routing using *fwmark*, a mark carried by the kernel as the packet travels through the kernel stack. Note that the Netfilter rules required to alter IPID field and to set *fwmark* are added/deleted at run-time since all possible entry-points are not known in advance.

Normally, a routing table specifies a single forwarding action to be taken in a deterministic manner for a given packet. The `CONFIG_IP_ROUTE_MULTIPATH` option in the kernel configuration permits specifying several alternative paths for a destination. If no weight is given, the kernel considers all these paths to be of equal cost and chooses in a non-deterministic way which *one* to use when a packet arrives. Instead, we would like to send the packet to multiple nodes *simultaneously*, if the multicast tree indicates that. Therefore, we wrote a Netfilter target module, called MULTIHOP, which sends a copy of the packet to each next-hop found in the routing rule for a given destination. In order to use this module, one needs to recompile the kernel to export a function required to access the routing table (*fib_lookup*). Other than this, no changes are required in the kernel. The module is available for download from our website.

Figure 3 shows the path of a packet through the Linux kernel and the places where it interacts with our scheme. Immediately after the packet gets in, the entry-point of the mesh must change the IPID field and set the TOS bit. Both the entry-point and any intermediate router set the *fwmark* when processing a packet for routing⁴. We use the `NF_IP_PRE_ROUTING` Netfilter hook to do these modifications. The packet is then passed back to the kernel networking stack, where it goes

³One of our goals is also to do as little changes as possible to the kernel, if any at all.

⁴*fwmark* is an internal mark in the kernel's packet data structure, and does not involve changing the packet as in the case of IPID and TOS.

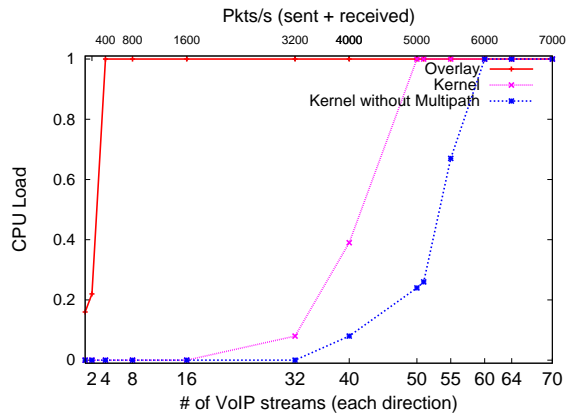


Fig. 4. Average CPU load.

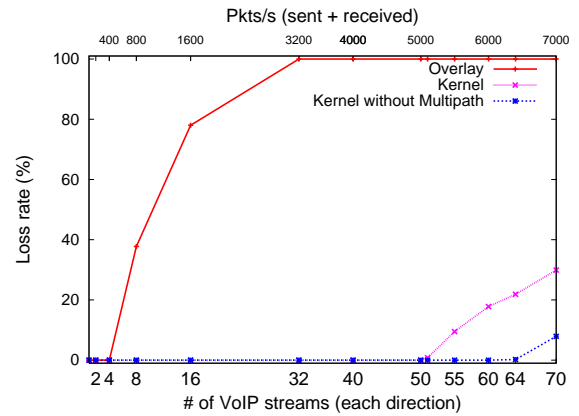


Fig. 5. Loss rate.

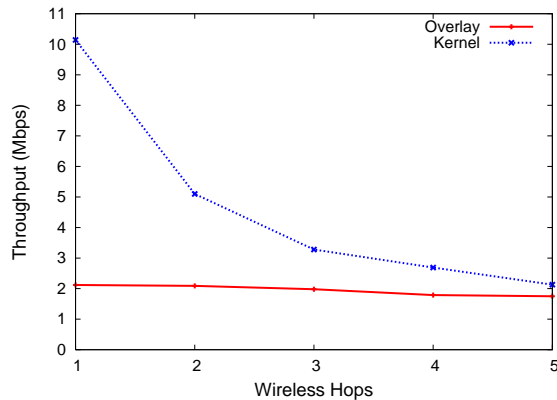


Fig. 6. Average TCP throughput.

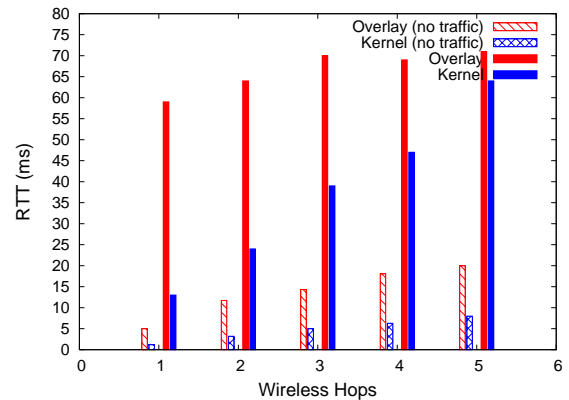


Fig. 7. Average RTT.

into the routing mechanism in which the *fwmark* is used to choose the appropriate routing table. After the routing decision is taken, but before leaving the interface, the packet reaches the *MULTIHOP* module. Additional copies of the packet are created if there is more than one next-hop in the route. The module will simply exit if there is only one next-hop. We register the *MULTIHOP* module at the `NF_IP_POST_ROUTING` Netfilter hook, such that there will be no routing decisions afterward.

C. Other considerations

The maximum number of clients supported by our architecture is limited only by the internal memory of the routers. We maintain one entry in each routing table per client, which requires a total of $32 \times N$ bytes in kernel memory, where N is the number of nodes in the mesh network. As an example, the Linksys WRT54G wireless router has 16 MB of RAM, and assuming only 5 MB is available to be used for the routing, we could theoretically support at least 9,000 mobile clients. In reality, this number is much greater because an entry is added in a routing table only if the router is on the path towards that client. As the size of the mesh network increases, more routing tables need to be maintained; however, as the clients are likely to be spread evenly throughout the network, the number of entries maintained by each router does not significantly grow.

IV. EXPERIMENTAL RESULTS

We evaluate our kernel redundant multipath scheme using low-cost Linksys WRT54G wireless routers running with a third-party OpenWrt firmware [2]. As this firmware was initially built using Linux kernel 2.4, we implemented our kernel modules for this version of Linux. Some of the experiments were performed using only wired connections, showing the CPU limitation in routers' performance, while other were performed in a wireless testbed of 17 nodes deployed among three buildings in our campus. In these experiments, the data rate was set to 24 Mbps⁵, the transmission power to 50 mW, and the 802.11 link-layer retransmission limit to 7.

A. Overlay vs Kernel CPU comparison test

The first experiment demonstrates the CPU limitation and its effect on the packet loss rate as the transmission rate increases. A client computer was connected to a router that was setup as an Internet gateway. Packets were routed to and from the client computer. To avoid losses caused by the wireless link-layer contention, we performed this test connecting the client to the router with a network cable. We sent 160-byte UDP packets, emulating VoIP streams, at rates corresponding to an increasing number of full-duplex VoIP

⁵Higher data rates are possible in denser networks.

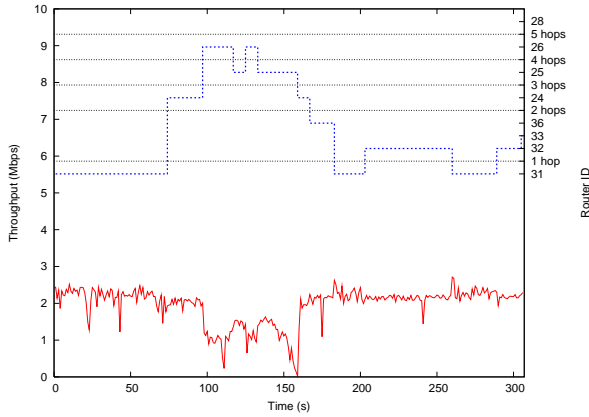


Fig. 8. TCP throughput: Overlay (moving).

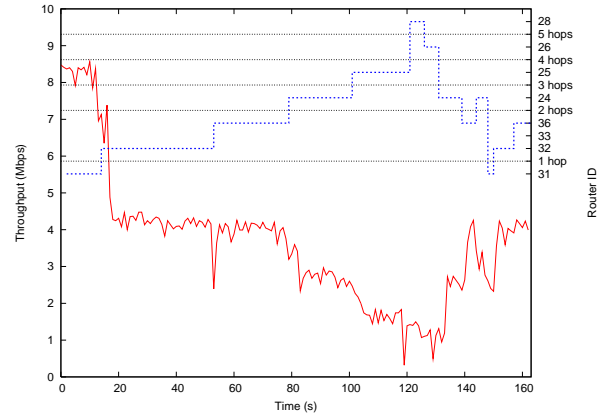


Fig. 9. TCP throughput: Kernel (moving).

streams between the client and the Internet. Each (one-way) stream has a rate of 64 Kbps. We monitored the average CPU load (Figure 4) and the loss rate (Figure 5). The top x -axis shows the corresponding number of packets/sec.

To understand better the overhead of our kernel approach, we included an additional scenario: kernel routing without the overhead of Netfilter rules required by our scheme. We can see that with the overlay implementation, the CPU starts to be saturated at 400 pkts/s (4 full-duplex VOIP streams, or 512 Kbps), in our kernel implementation at 5,000 pkts/s (50 streams or approx 6.4 Mbps) while in kernel “without redundant multipath” implementation at 6,000 pkts/s (60 streams or approx 7.6 Mbps). In each of these three scenarios, after a while, the loss rate starts to be non-zero: less than 8 streams (1 Mbps) for overlay, 51 streams (6.5 Mbps) for kernel and 64 streams (8 Mbps) for kernel without our additions.

B. Overlay vs Kernel TCP throughput and RTT test

This experiment evaluates the *maximum* throughput that can be achieved in a multi-hop wireless network when forwarding through user and kernel space with our modifications. We connected 5 Linksys WRT54G routers in a simple “line” topology, and measured the TCP throughput while sending traffic from Internet to the client. Note that this continues to be a very controlled test. We only use one client, and we do not use background traffic to influence our results, as our goal is to obtain the throughput upper bound. We performed tests with the client placed 1, 2, 3, 4 and 5 hops away from the Internet gateway. The throughput results are presented in Figure 6. We also measured the round trip time (RTT) for both, overlay and kernel routing, with and without background traffic (Figure 7). TCP throughput for 1 hop improved from a CPU-limited amount of 2.1 Mbps to a bandwidth limited amount of about 10 Mbps in our setup. The round-trip latency from overlay is more than 3 times the one from kernel for 1 hop, and even at 4 hops it is much above the kernel implementation.

C. Overlay vs Kernel in the deployed testbed:

1) *TCP throughput*: Figures 8 and 9 present the TCP throughput achieved over time in both overlay and kernel

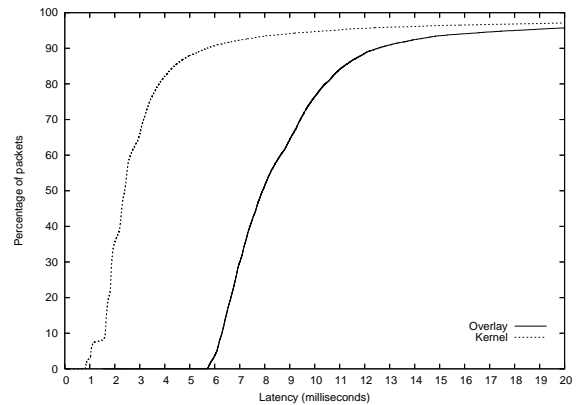


Fig. 10. Packet latency CDF.

modes while moving throughout our 17 node wireless mesh network. Note that, as the two experiments were done sequentially to test the maximum achievable throughput, the sequence of handoffs could not be replicated between the two runs, even though the mobility pattern was similar. The mesh network opportunistically associates the mobile device to access points depending of current conditions. In order to see how far we are from the Internet gateway, we plot with a dotted line the access point that currently services the client. The horizontal lines mark when the number of hops increases by one. To simplify the graphs, we included only the routers that were involved in handling the client. With the overlay routing, the throughput is just above 2 Mbps if the client is 1 or 2 hops away from the Internet gateway (routers 31 and 32, 33 and 36). This is consistent with the throughput reported in the previous test. As the number of hops increases, the throughput drops to 1 Mbps and even lower. In the kernel routing, the throughput was about 8.5 Mbps for 1 hop access points, 4.3 Mbps for 2 hops and it drops to 1 Mbps when the client is 6 hops away from the gateway (router 28).

2) *UDP latency*: We now compare the improvement of packets latencies when routing in overlay and in kernel modes. We use a full-duplex UDP traffic, consisting of 160-

byte packets sent every 20 ms at a rate of 64 Kbps, for 5 minutes. As before, we run this experiment with a client moving throughout the network. We focused on a VoIP-like traffic as a representative application that poses severe latency requirements. As opposed to the RTT presented in the controlled test (Figure 7), we measure the one-way latency of the packets received by the client. Figure 10 shows a CDF of the packets latency. We can see that about 80% of the packets were delivered in under 10.5 ms using the overlay while in the kernel implementation they arrived within 3.7 ms.

V. ADDITIONAL APPLICATIONS

Traditionally, redundancy is equivalent to resiliency, allowing the entire system to function as a whole when some of its components fail. In the case of routing, existing protocols are self-healing such that, upon detection of a route failure, the protocol will search and/or establish a new routing path if one exists. The problem lies in the process of detecting the failure, which usually involves a timeout to expire after a relatively large amount of time. While reducing these timeouts is not desirable in routing protocols, as they may generate dangerous oscillations, we believe that using redundant multipath forwarding in case of uncertainty offers a viable solution for resilient routing.

With our redundant multipath scheme, wireless networks could route each packet with a different level of resilience. Consider the network shown in Figure 11 where the source can reach the destination through one or more paths. These paths may be disjoint, and the number of paths may depend on the current state of the network and the levels of resilience needed (note that these routes cannot be managed by multicast trees). The graph shows the paths for three possible levels of redundancy. Utilizing our approach, different routing tables can have entries that correspond to different levels of redundancy. Then, each packet can carry its desired reliability level in its IPID so that each router is able to forward the packet using the appropriate forwarding table.

Redundant multipath can also be used to send time-sensitive information in hostile environments where nodes may be compromised by an adversary. While existing byzantine routing protocols [7] try to detect bad links and avoid routing through compromised nodes, this detection will take some time. With redundant multipath, the redundancy level of packets can be based on the importance of the data as well as on the current threat level in the working environment.

VI. CONCLUSION

This paper introduced a kernel-level redundant multipath mechanism to increase reliability in mesh networks. The architecture was integrated in a real system consisting of low-cost routers, resulting in a high throughput wireless mesh network with fast-handoff. We also showed that redundant multipath is a useful component to have in Unix-like kernels for other applications. Experimental results show that using the kernel-level redundant multipath service achieves close to optimal performance. An open source implementation of

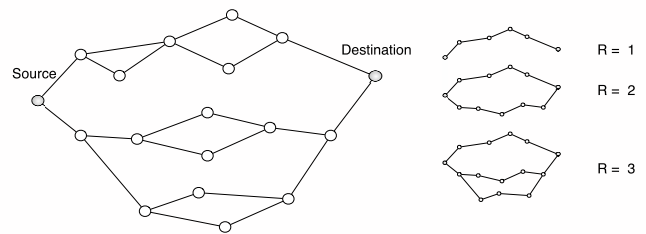


Fig. 11. Generic Redundant Multipath.

the kernel modules described in this paper is available at www.smesh.org.

REFERENCES

- [1] Microsoft research networking research group. <http://research.microsoft.com/mesh>.
- [2] OpenWrt. <http://openwrt.org>.
- [3] Yair Amir and Claudiu Danilov. Reliable communication in overlay networks. In *Proceedings of the IEEE DSN*, pages 511–520, June 2003.
- [4] Yair Amir, Claudiu Danilov, Michael Hilsdale, Raluca Musaloiu-Elefteri, and Nilo Rivera. Fast handoff for seamless wireless mesh networks. In *MobiSys 2006*, pages 83–95, New York, NY, USA, 2006. ACM Press.
- [5] Yair Amir, Claudiu Danilov, Raluca Musaloiu-Elefteri, and Nilo Rivera. An inter-domain routing protocol for multi-homed wireless mesh networks. *International Symposium on a World of Wireless, Mobile and Multimedia Networks (WoWMoM 2007)*, Helsinki, Finland, June 2007.
- [6] D. Andersen, H. Balakrishnan, F. Kaashoek, and R. Morris. Resilient overlay networks. In *Proc. of the 18th Symposium on Operating Systems Principles*, pages 131–145, October 2001.
- [7] Baruch Awerbuch, David Holmer, Cristina Nita-Rotaru, and Herbert Rubens. An on-demand secure routing protocol resilient to byzantine failures. In *WiSE '02: Proceedings of the 1st ACM workshop on Wireless security*, pages 21–30, New York, NY, USA, 2002. ACM Press.
- [8] K.C. Claffy C. Shannon, D. Moore. Beyond folklore: Observations on fragmented traffic. *IEEE/ACM Transactions on Networking*, 2002.
- [9] Ian D. Chakeres and Elizabeth M. Belding-Royer. Aodv routing protocol implementation design. In *ICDCSW '04: Proceedings of the 24th International Conference on Distributed Computing Systems Workshops - W7: EC (ICDCSW'04)*, pages 698–703, Washington, DC, USA, 2004. IEEE Computer Society.
- [10] Benjamin A. Chambers. The grid roofnet: a rooftop ad hoc wireless network. Master's thesis, Massachusetts Institute of Technology, May 2002.
- [11] Girish P. Chandranmenon and George Varghese. Reconsidering fragmentation and reassembly. In *PODC '98: Proceedings of the seventeenth annual ACM symposium on Principles of distributed computing*, pages 21–29. ACM Press, 1998.
- [12] C. Hedrick. Routing Information Protocol. *RFC1058*, June 1988.
- [13] Ahmed A.-G. Helmy, Muhammad Jaseemuddin, and Ganeshha Bhaskara. Multicast-based mobility: A novel architecture for efficient micromobility. *IEEE Journal on Selected Areas in Communications*, 2004.
- [14] Yang hua Chu, Sanjay G. Rao, and Hui Zhang. A Case For End System Multicast. In *Proceedings of ACM SIGMETRICS*, June 2000.
- [15] G. Chandranmenon S. Miller K. Almeroth E. Belding-Royer K. Ramachandran, M. M. Buddhikot. On the design and implementation of infrastructure mesh networks. *WiMesh*, 2005.
- [16] Vikas Kawadia, Yongguang Zhang, and Binita Gupta. System services for ad-hoc routing: Architecture, implementation and experiences. In *MobiSys 2003*, pages 99–112, New York, NY, USA, 2003. ACM Press.
- [17] J. Moy. Multicast extensions to OSPF. RFC 1584, IETF, March 1994.
- [18] J. Moy. OSPF Ver 2. *RFC2328*, April 1998.
- [19] S. Seshan, H. Balakrishnan, and R. Katz. Handoffs in Cellular Wireless Networks: The Daedalus Implementation and Experience. *Kluwer Journal on Wireless Personal Communications*, 1996., 1996.
- [20] T. von Eicken, A. Basu, V. Buch, and W. Vogels. U-net: a user-level network interface for parallel and distributed computing. In *SOSP*, pages 40–53, New York, NY, USA, 1995. ACM Press.