

# INSIGNIA: An IP-Based Quality of Service Framework for Mobile ad Hoc Networks

Seoung-Bum Lee, Gahng-Seop Ahn, Xiaowei Zhang, and Andrew T. Campbell

*Center for Telecommunication Research, Columbia University, New York, New York 10027*

Received January 15, 1999; revised August 15, 1999; accepted December 15, 1999

---

We present the design, implementation, and evaluation of INSIGNIA, an IP-based quality of service framework that supports adaptive services in mobile ad hoc networks. The framework is based on an in-band signaling and soft-state resource management approach that is well suited to supporting mobility and end-to-end quality of service in highly dynamic environments where the network topology, node connectivity, and end-to-end quality of service are time varying. Architecturally INSIGNIA is designed to support fast reservation, restoration, and end-to-end adaptation based on the inherent flexibility and robustness and scalability found in IP networks. We evaluate the framework, paying particular attention to the performance of the in-band signaling system, which helps counter time-varying network dynamics in support of the delivery of adaptive services. Our results show the benefit of our framework under diverse mobility, traffic, and channel conditions. © 2000 Academic Press

*Key Words:* quality of service; mobile ad hoc networks; adaptive services; in-band signaling; soft-state resource management.

---

## 1. INTRODUCTION

Mobile ad hoc networks are autonomous distributed systems that comprise a number of mobile nodes connected by wireless links forming arbitrary time-varying wireless network topologies. Mobile nodes function as hosts and routers. As hosts, they represent source and destination nodes in the network while as routers, they represent intermediate nodes between a source and destination, providing store-and-forward services to neighboring nodes. Nodes that constitute the wireless network infrastructure are free to move randomly and organize themselves in arbitrary fashions. Therefore the wireless topology that interconnects mobile hosts/routers can change rapidly in unpredictable ways or remain relatively static over long periods of time. These bandwidth-constrained multi-hop networks typically support best effort voice and data communications where the achieved “goodput” is often lower than the maximum radio transmission rate after encountering the effects of multiple access, fading, noise, and interference, etc. In addition to being bandwidth constrained,

mobile ad hoc networks are power constrained because network nodes rely on battery power for energy. Providing suitable quality of service (QoS) support for the delivery of real-time audio, video and data in mobile ad hoc networks presents a number of significant technical challenges.

Mobile ad hoc networks may be large, which makes the problem of network control difficult. The end-to-end communication abstraction between two communicating mobile hosts can be viewed as a complex “end-to-end channel” that may change route over time. There may be a number of possible routes between two communicating hosts over which data can flow, and each path may have different available capacity that may or may not meet the quality of service requirements of the desired service. Even if the selected path between a source–destination pair meets the user’s needs at the session set-up time, the capacity and error characteristics observed along the path are likely to be time varying due to the multiple dynamics that operate in the network.

The fading effects resulting from host mobility cannot always be masked by the link layer and typically result in discernible effects on the application’s perceptible quality (e.g., assured delivery of audio/video may degrade rapidly). This affects the capacity of a given path through the network, where links tend to degrade slowly at first and then rapidly drop out. This results in topological dynamics that operate on slower time scales than channel fades and other such discontinuities. Reacting to these network capacity dynamics over the appropriate time scale requires fast, lightweight, and responsive protocol operations. Flows must be established, maintained, and removed in mobile ad hoc networks over the course of a user-to-user session. Typically, “connections” (i.e., the establishment of “state” information at nodes along the path) need to be maintained and automatically renegotiated in response to the network topology dynamics and link quality changes. Since resources are scarce in these networks, any protocol signaling overhead needed to maintain connections limits the utilization of the network. Therefore, bandwidth required to support signaling systems must be kept to a minimum. This places emphasis on minimizing the signaling required to establish, maintain, restore, and tear down network states associated with user sessions. In addition, due to the disconnected nature of maintaining state in mobile ad hoc networks, explicit tear-down mechanisms (e.g., disconnect signaling) are impractical. This is due to the fact that it is infeasible to explicitly remove network state (established during session setup) in portions of the network that are out of radio contact of a signaling controller due to topology changes.

There is a need for new mobile ad hoc architectures, services, and protocols to be developed in response to these challenges. New control systems need to be highly adaptive and responsive to changes in the available resources along the path between two communicating mobile hosts. Future protocols need to be capable of differentiating between the different service requirements of user sessions (e.g., continuous media flows, microflows, RPC, etc.). Packets associated with a flow traversing intermediate nodes (as illustrated in Fig. 1) between a source and destination may, for example, require special processing to meet end-to-end bandwidth and delay constraints. When building quality of service support into mobile ad hoc networks the design of fast routing algorithms that can efficiently track

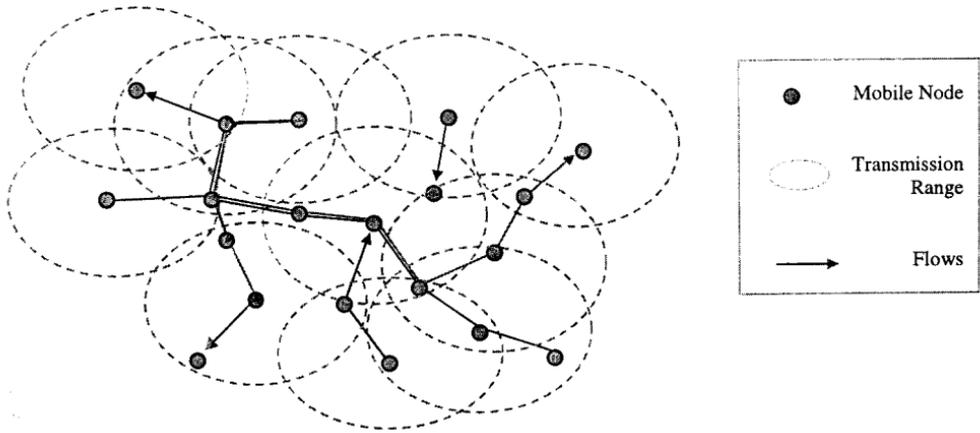


FIG. 1. Mobile ad hoc networking.

network topology-changes is important. Mobile ad hoc network routing protocols need to work in unison with efficient signaling, control, and management mechanisms to achieve end-to-end service quality. These mechanisms should consume minimal bandwidth in operation and react promptly to changes in the network state (viewed in terms of changes in the network topology) and flow state (viewed in terms of changes in the observed end-to-end quality of service).

In this paper, we present the design, implementation, and evaluation of the *INSIGNIA QOS Framework* that supports the delivery of adaptive services in mobile ad hoc networks. A key component of our QOS framework is the *INSIGNIA* signaling system, an in-band signaling system that supports fast reservation, restoration, and adaptation algorithms that are specifically designed to deliver adaptive service. The signaling system is designed to be lightweight and highly responsive to changes in network topology, node connectivity, and end-to-end quality of service conditions. The structure of the paper is as follows. We discuss our framework in the context of the related work and present the main design considerations that have influenced our thinking in Sections 2 and 3, respectively. Section 4 presents an overview of the *INSIGNIA QOS* framework. The detailed design of the *INSIGNIA* signaling system is given in Section 5. We evaluate our QOS framework in Section 6, paying particular attention to the performance of the signaling system under a variety of network conditions. Our simulation results show the benefit of the *INSIGNIA QOS* framework under diverse mobility, traffic, and channel conditions in support of fast reservation, restoration, and adaptation. Finally, we present our conclusion in Section 7.

## 2. RELATED WORK

Research and development of mobile ad hoc networking technology is proceeding in both academia and industry under military and commercial sponsorship. Current military research projects such as the Army Research Office Focused Research Initiatives, the Army Research Laboratory Federated Laboratory, and the DARPA

Global Mobile Information Systems (GloMo) program [23] are producing new technologies.

There has been little research in the area of supporting quality of service in mobile ad hoc networks, however. What work exists tends to be based on distributed scheduling algorithms [32] that address rescheduling when the network topology changes and QOS-based medium access controllers [34]. Typically, these schemes are based on a single link layer network technology and not on an interconnection of different wireless technologies at the IP layer. In addition, the work does not address suitable support for adaptive QOS paradigms that are required to deliver adaptive services in mobile ad hoc networks. In this paper, we propose an IP-based QOS framework, adaptive services, and support protocols incorporating a soft-state [3] resource management system. This system is based on in-band signaling techniques supporting reservation across multiple link layer radio technologies that map to specific link layer access technologies for distributed packet scheduling. Our contribution addresses a suitable IP level control architecture for delivering adaptive services in mobile ad hoc networks. We do not, however, propose any new distributed scheduling techniques. Rather, we leverage the existing body of work found in the literature as a basis for the provision of QOS support over radios.

In [9, 13], multi-hop, multi-cluster packet radio network architectures are proposed. The provisioning of quality of service is discussed based on “dynamic virtual circuit” communications derived from wireline network control and signaling found in ATM networks. This approach relies on a “circuit” model that requires explicit connection management and the establishment of hard state in the network prior to communication. We believe there is a need to investigate alternative network models that are more responsive to the dynamics found in ad hoc networks other than the hard-state virtual circuits. Typically, virtual circuits are established across mobile ad hoc networks using explicit “out-of-band” signaling to set up reservations for the duration of the call/session holding time. We believe that flows/sessions should be established and maintained using a faster, more responsive system based on soft-state and in-band signaling paradigms. We believe that virtual circuits lack the intrinsic flexibility needed to adapt to the dynamics found in mobile ad hoc networks and that the notion of “soft-state connections” driven by in-band techniques is more suitable. There is a need to develop new QOS architectures that can provide fast reservation, responsive restoration, and seamless adaptation to mobile ad hoc network dynamics based on the inherent flexibility, robustness, and scalability found in IP networks.

Delivering end-to-end service quality in mobile ad hoc networks is intrinsically linked to the performance of the routing protocol because new routes or alternative routes between source-destination pairs need to be periodically computed during ongoing sessions. The IETF Mobile Ad Hoc Networks (MANET) Working Group [2] recently began to standardize internetwork layer technologies (i.e., routing and support protocols). As such, it is presently focused on standardizing network-layer routing protocols suitable for supporting best effort packet delivery in IP-based networks. Within this context there have been a number of proposals for efficient routing that dynamically track changes in mobile ad hoc network topology, including the Temporally Ordered Routing Algorithm (TORA) [1], Dynamic Source

Routing [7], Zone Routing Protocol [5], and Ad Hoc On Demand Distance Vector Routing Protocol [6]. The performance of a QOS framework will rely on the speed at which routing protocols can compute new routes (if no alternative route is currently cached) after topology changes have occurred. The delay in computing new routes will have an impact on the QOS delivered to ongoing sessions. For a comparison of mobile ad hoc routing protocols see [21].

### 3. DESIGN CONSIDERATIONS

#### 3.1. Adaptive Services

The most suitable service paradigm for mobile ad hoc networks is adaptive in nature. We observe that adaptive voice and video applications operating in mobile cellular networks are capable of responding to packet loss, delay jitter, changes in available bandwidth, and handoff while maintaining some level of service quality [28]. While adaptive multimedia applications can respond to network dynamics they typically require some minimum bandwidth assurance below which they are rendered useless.

The INSIGNIA QOS framework is designed to support adaptive services as a primary goal. In this context, adaptive services provide minimum bandwidth assurances to real-time voice and video flows and data allowing for enhanced levels (i.e., maximum bandwidth) of service to be delivered when resources become available. A flow represents a sequence of packets sent from a single source to one or more destinations representing a single media type (e.g., voice, video, etc.). Flows require admission control, resource reservation, and maintenance at all intermediate routers between a source and destination to provide end-to-end quality of service support. Typically, continuous media flows are long lived in comparison to microflows, which represent short-lived flows (e.g., web style client/server interactions) that comprise a limited train of data packets. We use the terms “session,” “flow,” “continuous media flow,” and “microflow” interchangeably in this paper. The INSIGNIA QOS framework is designed to transparently support the requirements of continuous media flows and microflows. Adaptive services support applications that require *base QOS* (i.e., minimum bandwidth) and *enhanced QOS* (i.e., maximum bandwidth) assurances, respectively. The semantics of the adaptive service provides preference to packets associated with the base QOS over enhanced QOS.

Adaptation is an application-specific process. Some applications may be incapable of adapting while others may adapt discretely (e.g., scalable profiles of MPEG2) or continuously (e.g., dynamic rate-shaped applications [28]). The time scale over which applications can adapt is also application specific. For example, greedy data applications (e.g., image downloads) may want to take advantage of any change in available bandwidth at any time. In contrast, adaptive continuous media applications (e.g., audio and video) may prefer to follow trends (via some low pass filtering scheme) in available bandwidth based on slower adaptation time scales, preferring some level of “stable” service delivery rather than responding to every instantaneous change in bandwidth availability. Adaptive applications therefore should

manage the adaptation process and dictate the time scales and semantics of their adaptation process. Given this observation, our QOS framework is designed to adapt user sessions to the available level of service without explicit signaling between source–destination pairs. In this case the network and application adapt to different dynamics. The network adapts (via restoration algorithms) to changes in topology and measured channel conditions while trying to deliver base and enhanced QOS. Applications adapt to the observed end-to-end QOS fluctuations within the prescribed max-min limits based on application specific adaptation time scales. This observation drives a number of architectural design decisions.

### *3.2. Separation of Routing, Signaling, and Forwarding*

There has been a growing amount of work in the area of QOS routing for fixed networks. Here the routing protocols interact with resource management to establish paths through the network that meet end-to-end QOS requirements (i.e., delay, bandwidth, possibly multi-metrics demands). In this case there is a certain level of integration of resource management and routing. One could apply such an approach to MANET routing protocols given that the time scales over which new routes are computed are much faster than traditionally found in the case of routing in fixed infrastructures. While we believe this a promising approach (see the CEDAR [35] proposal) we note that the time scales over which session setup and routing (i.e., computing new routes) operate are distinct and functionally independent tasks. Therefore, we believe that signaling, resource management, and routing should be modeled independently in the network architecture.

We consider that MANET routing protocols should not be burdened with the integration of QOS functionality that may be tailored toward specific QOS models. Rather, we argue that it is better to maintain a clean separation between routing, signaling, and forwarding. These architectural components are rather different from one another in the algorithms they implement and in the time scales over which they operate. Our approach is to develop a QOS framework that can “pluggin” a wide variety of routing protocols. In this case, resource reservation and signaling will be capable of interacting with any number of routing protocols to provide end-to-end QOS support. Different MANET routing protocols clearly perform differently [21] in response to topology changes while the QOS framework attempts to maintain end-to-end service quality.

### *3.3. In-band Signaling*

In-band signaling systems are capable of operating close to packet transmission speeds and are therefore well suited to responding to the fast time scale dynamics found in mobile ad hoc environments, as illustrated in Fig. 1. The term “in-band signaling” refers to the fact that the control information is carried along with data. In contrast, out-of-band signaling systems (e.g., Internets RSVP, ATMs UNI, etc.) are incapable of responding to such fast time-scale dynamics because out-of-band signaling systems require maintenance of source route information and respond to

topology changes by directly signaling “affected mobiles” to allocate/free resources. In some cases, this is impossible due to lack of connectivity between “affected routers” and the signaling entity that attempts to deallocate resources over the old path.

The term “out-of-band signaling” refers to fact that the control information is typically carried in separate control packets and on channels that may be distinct from the data path. Based on an in-band approach, the INSIGNIA signaling system can restore the flowstate (i.e., a reservation) in response to topology changes within the interval of two consecutive IP packets under ideal conditions. INSIGNIA performance relies on the speed at which the routing protocol can recompute new routes if no alternative route is cached after topology changes. Out-of-band signaling systems, for example, would need to maintain source route information and respond to topology changes by directly signaling intermediate routers on an old path to allocate/free radio resources. In many case, this is impossible to do if the affected router is out of radio contact from the signaling entity that attempts to deallocate resources over the old path.

### *3.4. Soft-State Resource Management*

Maintaining the QOS of adaptive flows in mobile ad hoc networks is one of the most challenging aspects of the INSIGNIA QOS framework. In wireline networks that support quality of service and state management, the route and the reservation between source–destination pairs remain fixed for the duration of a session. This style of hard-state connection-oriented communications (e.g., virtual circuit) guarantees quality of service for the duration of the session holding time. However, these techniques are not flexible enough in mobile ad hoc networks, where the path and reservation need to dynamically respond to topology changes in a timely manner.

We believe that a soft-state approach to state management at intermediate routing nodes is suitable for the management of reservations in mobile ad hoc networks. Such an approach models the transient nature of network reservations, which have to be responsive to fast time-scale wireless dynamics, moderate time-scale mobility changes, and longer time-scale session “holding times.” Based on the work by Clark [3], soft state relies on the fact that a source sends data packets along an existing path. If a data packet arrives at a mobile router and no reservation exists then admission control and resource reservations attempt to establish soft state. Subsequent reception of data packets (associated with a reservation) at that router are used to refresh the existing soft-state reservation. This is called a “soft connection” when considered on an end-to-end basis and in relation to the virtual circuit hard-state model. When an intermediate node receives a data packet that has an existing reservation it reconfirms the reservation over the next interval. Therefore the holding time of a soft connection is based on a soft-state timer interval and not on session duration holding time. If a new packet is not received within the soft-state timer interval then resources are released and flow states removed in a fully decentralized manner.

We believe that the development of new QOS frameworks based on the notion of in-band signaling and soft-state resource management and constructed with separation of routing, signaling, and forwarding functions will provide a responsive, scalable, and flexible solution for delivering adaptive services in mobile ad hoc networks.

#### 4. THE INSIGNIA QOS FRAMEWORK

The INSIGNIA QOS framework allows packet audio, video, and real-time data applications to specify their maximum and minimum bandwidth needs and plays a central role in resource allocation, restoration control, and session adaptation between communicating mobile hosts. Based on the availability of end-to-end bandwidth, QOS mechanisms attempt to provide assurances in support of adaptive services. To support adaptive service, the INSIGNIA QOS framework establishes and maintains reservations for continuous media flows and microflows. To support these communication services the INSIGNIA QOS framework comprises the following architectural components as illustrated in Fig. 2:

- *In-band signaling* establishes, restores, adapts, and tears down adaptive services between source-destination pairs. Flow restoration algorithms respond to dynamic route changes, and adaptation algorithms respond to changes in available bandwidth. Based on an in-band signaling approach that explicitly carries control information in the IP packet header, flows/sessions can be rapidly established, restored, adapted, and released in response to wireless impairments and topology changes.

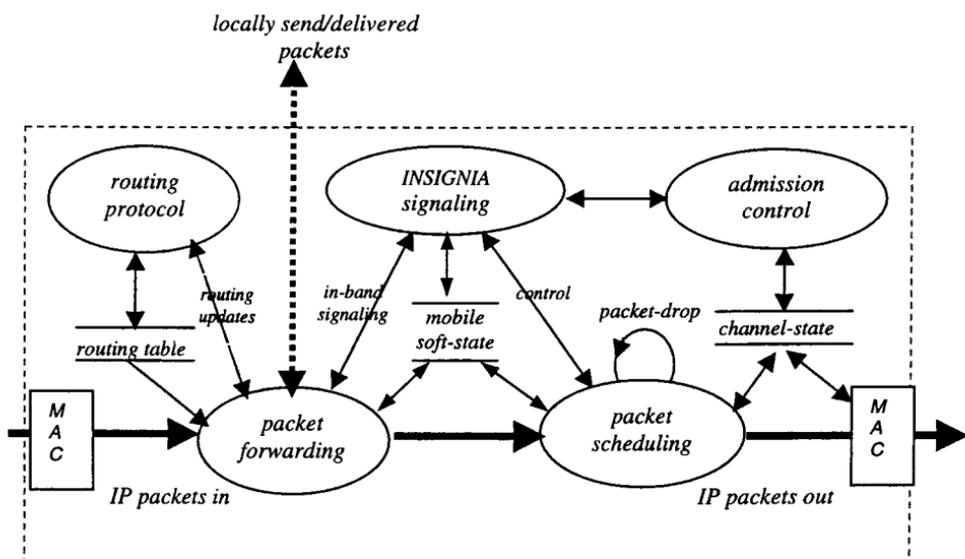


FIG. 2. INSIGNIA QOS framework.

- *Admission control* is responsible for allocating bandwidth to flows based on the maximum/minimum bandwidth (i.e., base and enhanced QOS) requested. Once resources have been allocated they are periodically refreshed by a soft-state mechanism through the reception of data packets. Admission control testing is based on the measured channel capacity/utilization and requested bandwidth. To keep the signaling protocol simple and lightweight, new reservation requests do not impact existing reservations.

- *Packet forwarding* classifies incoming packets and forwards them to the appropriate module (viz. routing, signaling, local applications, packet scheduling modules). Signaling messages are processed by INSIGNIA signaling, and data packets are delivered locally (as illustrated by the dashed line in Fig. 2) or forwarded to the packet scheduling module (as illustrated by the bold line in Fig. 2) for transmission on to the next hop.

- *Routing* dynamically tracks changes in ad hoc network topology, making the routing table visible to the node's packet forwarding engine. The QOS framework assumes the availability of a generic set of MANET routing protocols [2] that can be plugged into the architecture. The QOS framework assumes that the routing protocol provides new routes, either proactively or on demand, in the case of topology changes.

- *Packet scheduling* responds to location-dependent channel conditions when scheduling packets in wireless networks [22]. A wide variety of scheduling disciplines can be used to realize the packet scheduling module and the service model. Currently, we have implemented a weighted round robin [22, 25] service discipline based on an implementation [29] of deficit round robin that has been extended to provide compensation in the case of location-dependent channel conditions between mobile nodes.

- *Medium access control (MAC)* provides quality of service-driven access to the shared wireless media for adaptive and best effort services. The INSIGNIA QOS framework is designed to be transparent to any underlying media access control protocols and is positioned to operate over multiple link layer technologies at the IP layer. However, the performance of the framework is strongly coupled to the provisioning of QOS support provided by specific medium access controllers.

## 5. THE INSIGNIA SIGNALING SYSTEM

The INSIGNIA signaling system plays an important role in establishing, adapting, restoring, and terminating end-to-end reservations. In what follows, we describe the INSIGNIA in-band signaling approach. The signaling system is designed to be lightweight in terms of the amount of bandwidth consumed for network control and to be capable of reacting to fast network dynamics such as rapid host mobility, wireless link degradation, and intermittent session connectivity. We discuss the protocol commands and protocol mechanisms.

### 5.1. Protocol Commands

Protocol commands are encoded using the IP option field and include service mode, payload type, bandwidth indicator, and bandwidth request field, as illustrated in Fig. 3. By adopting an INSIGNIA IP option in each IP packet header the complexity of supporting packet encapsulation inside the network is avoided. These protocol commands support the signaling algorithms discussed in Section 5.2 including flow reservation, restoration, and adaptation mechanisms. The protocol commands drive the state operations of the protocol. Figure 4 presents a simplified view of the finite state machines for a source host, intermediate router, and destination host. These three state machines capture the major event/actions and resulting state transitions. We use these state machines to illustrate the dynamics of the INSIGNIA signaling system.

**5.1.1. Service mode.** When a source node wants to establish a fast reservation to a destination node it sets the *reservation (RES) mode* bit in the INSIGNIA IP option service mode of a data packet and sends the packet toward the destination. On reception of a RES packet intermediate routing nodes execute admission control to accept or deny the request. When a node accepts a request, resources are committed and subsequent packets are scheduled accordingly. In contrast, if the reservation is denied, packets are treated as *best effort (BE) mode* packets.

In the case where a RES packet is received and no resources have been allocated, the admission controller attempts to make a new reservation. This condition commonly occurs when flows are rerouted during the lifetime of an ongoing session due to host mobility. When the destination receives a RES packet it sends a QOS report to the source node indicating that an end-to-end reservation has been established and transitions its internal state from best effort to reservation state, as illustrated in Fig. 4c.

The service mode indicates the level of service assurance requested in support of adaptive services. The interpretation of the service mode, which indicates a RES or BE packet, is dependent on the payload type and bandwidth indicator discussed in Sections 5.1.3 and 5.1.4, respectively. A packet with the service mode set to RES and bandwidth indicator set to *MAX* or *MIN* is attempting to set up a max-reserved or min-reserved service, respectively. The bandwidth requirements of the flow are carried in the bandwidth request field, as illustrated in Fig. 3. A RES packet may be degraded to BE service in the case of rerouting or insufficient resource availability along the new/existing route. Note that a BE packet requires no resource reservation to be made.

The IP option also carries an indication of the payload type, which identifies whether the packet is a base QOS (*BQ*) or enhanced QOS (*EQ*) packet, as discussed

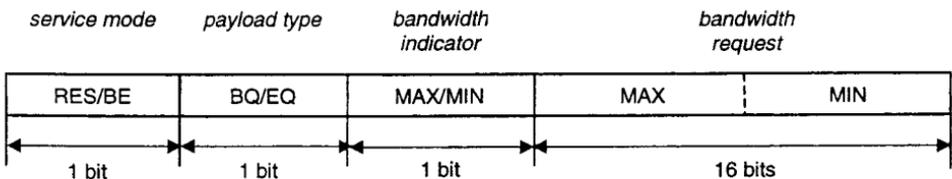


FIG. 3. INSIGNIA IP option.

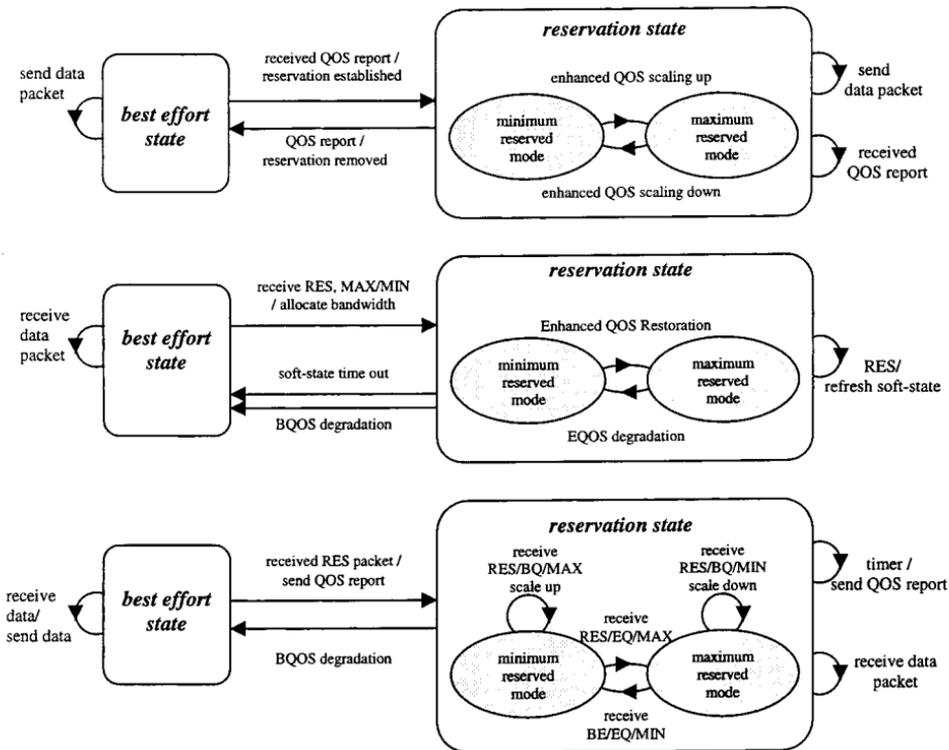


FIG. 4. State machine (a) a source mobile host, (b) an intermediate mobile node, and (c) a destination mobile host.

in Section 5.1.3. Using the “packet state” (service mode/payload type/bandwidth indicator) one can determine which component of the flow is degraded. Reception of a BE/EQ/MIN packet or RES/BQ/MIN indicates that the enhanced QOS packets have been degraded to best effort service. By monitoring the packet state the destination node can issue scaling/drop commands to the source based on the destination state machine illustrated in Fig. 4c.

As shown in Fig. 4 the source, intermediate, and destination state machines support two reservation substates:

- *max-reserved mode* provides reservation for a flow’s base and enhanced QOS packets. This type of service requires successful end-to-end reservation to meet a flow’s maximum bandwidth needs (e.g., RES/EQ/MAX).
- *min-reserved mode* provides reservation for the base QOS and best effort delivery for the enhanced QOS component (if it exists). This service mode typically occurs when max-reserved flows experience degradation in the network. For example, max-reserved flows may encounter mobile nodes that lack the resources to support both the base and enhanced QOS, resulting in the degradation of enhanced QOS packets to best effort delivery (e.g., BE/EQ/MIN).

5.1.2. *Bandwidth request.* The bandwidth request allows a source to specify its maximum (*MAX*) and minimum (*MIN*) bandwidth requirements for adaptive services. This assumes that the source has selected the RES service mode. A source

may also simply specify a minimum or a maximum bandwidth requirement. For adaptive services the base QOS (min-reserved service) is supported by the minimum bandwidth, whereas the maximum bandwidth supports the delivery of the base and enhanced QOS (max-reserved service) between source–destination pairs. Flows are represented as having minimum and maximum bandwidth requirements. This characterization is commonly used for multi-resolution traffic (e.g., MPEG audio and video), adaptive real-time data that has discrete max-min requirements, and differential services that support prioritization of aggregated data in the Internet.

**5.1.3. Payload type.** The payload field indicates the type of packet being transported. INSIGNIA supports two types of payload called base QOS and enhanced QOS, which are reserved via distributed end-to-end admission control and resource reservation. The semantics of the adaptive services are related to the payload type and available resources (e.g., enhanced QOS requires that maximum bandwidth requirements can be met along the path between a source–destination pair). The semantics of the base and enhanced QOS are application specific. They can represent a simple prioritization scheme between packets, differential services, or self-contained packet streams associated with multi-resolution flows. The adaptation process may force adaptive flows to degrade when insufficient resources are available to support the maximum bandwidth along the existing path or during restoration when the new path has insufficient resources. For example, if there is only sufficient bandwidth to meet the minimum bandwidth requirement needs of the base QOS, enhanced QOS packets are degraded to best effort packets at bottleneck nodes by simply flipping the service mode for EQ packets from RES to BE. When a downstream node detects degraded packets, it releases any resources that may have previously been allocated to support the transport of enhanced QOS packets. The adaptation process (discussed in Section 5.2.5) is also capable of scaling flows up by taking advantage of any additional bandwidth availability that may be encountered along a new/existing path. In this case, a flow could be “scaled-up” from min-reserved to max-reserved mode delivery, as indicated in Figs. 4a and 4c.

**5.1.4. Bandwidth indicator.** A bandwidth indicator plays an important role during reservation setup and adaptation. During reservation establishment the bandwidth indicator reflects the resource availability at intermediate nodes along the path between source–destination pairs. Reception of a setup request packet with the bandwidth indicator bit set to MAX indicates that all nodes encountered enroute have sufficient resources to support the maximum bandwidth requested (i.e., max-reserved mode). In contrast, a bandwidth indicator set to MIN implies that at least one of the intermediate nodes between the source and destination is a bottleneck node and insufficient bandwidth is available to meet the maximum bandwidth requirement; that is, only min-reserved mode delivery can be supported. In this case, adaptation algorithms at the destination can trigger the signaling protocol to release any overallocated resources between the source and bottleneck node by issuing a “drop” command to the source node (see Section 5.2.5 on adaptation). A bandwidth indicator set to MIN does, however, indicate that the mobile ad hoc network can support the minimum requested bandwidth (i.e., min-reserved mode). The bandwidth indicator is also utilized during the adaptation of

ongoing sessions in this manner. The adaptation mechanism resident at the destination host continuously monitors the bandwidth indicator to determine whether any additional bandwidth is available to support better service quality.

## 5.2. Protocol Operations

In what follows, we provide an overview of the main protocol mechanisms and state machines for the source, intermediate router, and destination nodes as illustrated in Fig. 4. The key signaling components include fast reservation, QOS reporting, soft-state resource management, restoration, and flow adaptation.

**5.2.1. Fast reservation.** To establish adaptive flows, source nodes initiate reservations by setting the appropriate field of the IP option in data messages before forwarding “reservation request” [5] packets toward destination nodes. A reservation request packet is characterized as having the service mode set to RES, payload set to BQ/EQ, bandwidth indicator set to MAX/MIN, and valid bandwidth requirements. Reservation packets traverse intermediate nodes executing admission control modules, allocating resources, and establishing low state at all intermediate nodes between source–destination pairs, as illustrated in Fig. 5. A source node continues to send reservation packets until the destination node completes the reservation setup phase by informing the source node of the status of the flow establishment phase using QOS reporting, as shown in Fig. 5.

The establishment of an adaptive flow is illustrated in Fig. 5. A source node ( $M_s$ ) requests maximum resource allocation and node  $M_1$  performs admission control upon reception of the reservation packet. Resources are allocated if available, and the reservation packet is forwarded to the next node  $M_2$ . This process is repeated on a hop-by-hop basis until the reservation packet reaches the destination mobile  $M_D$ . The destination node determines the resource allocation status by checking the packet state (i.e., service mode, payload type, and bandwidth indicator). The QOS reporting mechanism is used to inform the source node of the reservation status en route. As far as the destination node is concerned the reservation phase is complete on reception of the first RES packet. From the example shown in Fig. 5, we see that

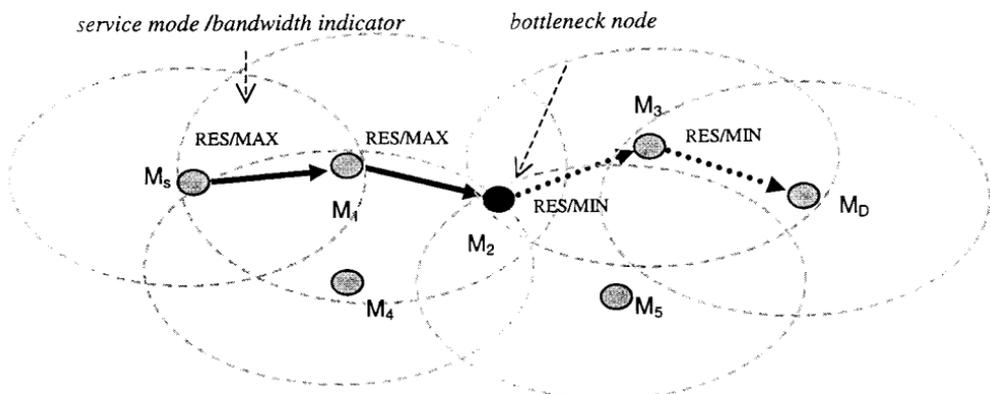


FIG. 5. Fast reservation.

only the minimum bandwidth is supported between  $M_2$  and  $M_3$  and subsequent nodes receiving the request packet avoid allocating resources for the maximum.

When a reservation is received at the destination node, the signaling module checks the reservation establishment status. The status is determined by inspecting the IP option field service mode, which should be set to RES. If the bandwidth indication is set to MAX, this implies that all nodes between a source-destination pair have successfully allocated resources to meet the base and enhanced bandwidth needs in support of the max-reserved mode. On the other hand, if the bandwidth indication is set to MIN this indicates that only the base QOS can be currently supported (i.e., min-reserved mode). In this case, all reservation packets with a payload of EQ received at a destination will have their service level flipped from RES to BE by the bottleneck node. As a result “partial reservations” will exist between the source and bottleneck node (e.g., between  $M_S$  and  $M_2$  in Fig. 5). In the case of partial reservations, resources remain reserved between the source and the bottleneck node until explicitly released. Release of partial reserved resources can be initiated by the source based on feed back during the reservation phase or as part of the adaptation process where the destination can issue “scale-down/drop” commands to a source node. This will have the effect of clearing any partial reservation (e.g., between  $M_S$  and  $M_2$  in Fig. 5). An application may choose not to deallocate a partial reservation, hedging that bandwidth will become available at the bottleneck node allowing for a full end-to-end reservation to be made in due course.

Note that if a reservation has been established for the maximum reserved state and a RES/BQ/MIN packet is persistently received in this substate then the state machine determines that the enhanced QOS packets have been degraded and transitions to minimum reserved state in anticipation of scaling back up. This behavior is illustrated in Fig. 4c. Degradation of this sort can occur at intermediate nodes due to insufficient resources to support a new reservation, or an ongoing flow can be degraded due to rerouting or insufficient resource availability on the new/existing path. The state information maintained at the destination can decode which of these conditions occurred.

*5.2.2. QOS reporting.* QOS reporting is used to inform source nodes of the ongoing status of flows. Destination nodes actively monitor ongoing flows, inspecting status information (e.g., bandwidth indication) and measuring the delivered QOS (e.g., packet loss, delay, throughput, etc.). QOS reports are also sent to source nodes for completing the reservation phase and on a periodic basis for managing end-to-end adaptations. QOS reports do not have to travel on the reverse path toward a source. Typically they will take an alternate route through the ad hoc network, as illustrated in Fig. 6. Although the QOS reports are basically generated periodically according to the applications’ sensitivity to the service quality, QOS reports are sent immediately when required (i.e., typically actions related to adaptation).

In the case where only the BQ packets can be supported, as is the case with the min-reserved mode, the signaling system at the source “flips” the service mode of the BQ packets from RES to BE sending all “degraded” packets sent as best effort.

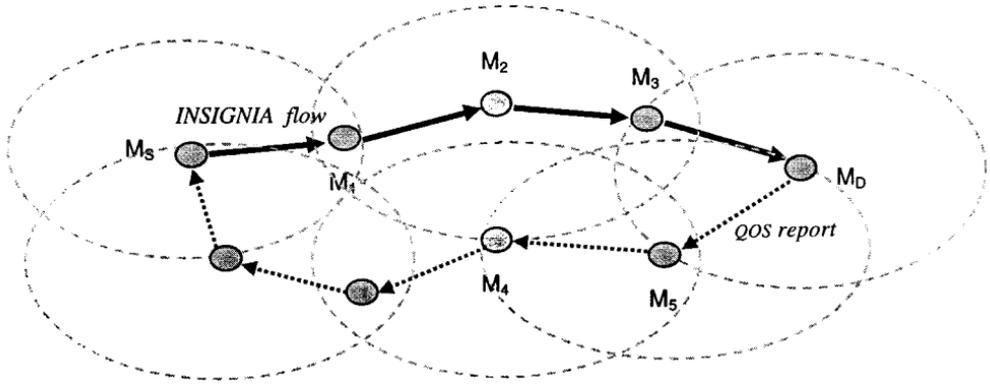


FIG. 6. QoS reporting.

Any partial reservations that may exist between source and destination nodes are automatically timed out after “flipping” the state variable in the EQ packets. Since there is a lack of EQ packets with the RES bit set at intermediate routers any associated resources are released (e.g., between  $M_S$  and  $M_2$  in Fig. 5), allowing other competing flows to contend for these resources. In a similar fashion QoS reports are also used as part of the ongoing adaptation process that responds to mobility and resource changes in the network. The adaptation process is discussed in Section 5.2.5.

**5.2.3. Soft-state resource management.** Reservations made at intermediate routing nodes between source and destination pairs are driven by soft-state management, as indicated in Fig. 4b. A soft-state approach is well suited for management of resources in dynamic environments, where the path and reservation associated with a flow may change rapidly. The transmission of data packets is strongly coupled to the maintenance of flow states (i.e., reservations). In other words, as the route changes in the network, new reservations will be automatically restored by the restoration mechanism. A major benefit of the soft state is that resources allocated during flow establishment are automatically removed when the path changes. For example, the

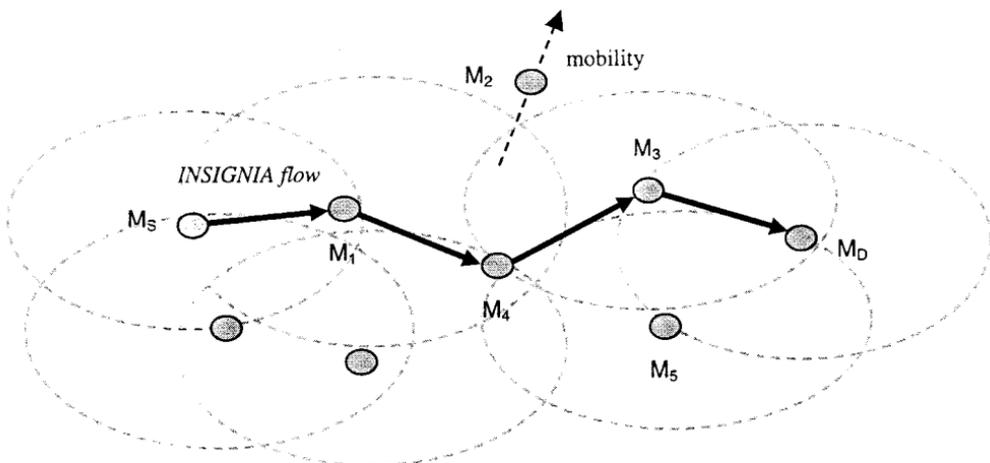


FIG. 7. Rerouting and restoration.

mobility of node  $M_2$  in Fig. 7 will cause flows to be rerouted via intermediate routers  $M_1-M_4-M_3$ . Due to the absence of reserved mode data packets at node  $M_2$  the node will automatically release resources associated with the flow without any interaction from any explicit controller.

Once admission control has accepted a request for a new flow, soft-state management will start the soft-state timer associated with the new or rerouted flow. The soft-state timer is continually refreshed as long as packets associated with a flow are periodically received at intermediate routers. In contrast, if packets are not received (e.g., due to rerouting) then the soft state is not refreshed but times out with the result of deallocating any resources. Since data packets are used to maintain the state at intermediate nodes we couple the data rate of flows to the soft-state timer value. In Section 6.4, we evaluate the performance of a fixed and dynamic scheme for determining the soft-state timer value. The fixed scheme simply sets a value for all flows regardless of the data rate of individual flows (e.g., RSVP recommends 30 s), and the dynamic scheme tracks the changing data rate of individual flows and sets the soft-state timer accordingly.

*5.2.4. Restoration.* Flows are often rerouted within the lifetime of ongoing sessions due to host mobility. The goal of flow restoration is to reestablish reservation as quickly and efficiently as possible. Rerouting active flows involves the routing protocol (to determine a new route), admission control, and resources reservation for nodes that belong to a new path. Restoration procedures also call for the removal of old flow state at nodes along the old path. In an ideal case, the restoration of flows can be accomplished within the duration of a few consecutive packets given that an alternative route is cached. We call this type of restoration “immediate restoration.” If no alternative route is cached, the performance of the restoration algorithm is coupled to the speed at which the routing protocols can discover a new path.

As illustrated in Fig. 7, network dynamics trigger rerouting and service degradation. In this example, mobile host  $M_2$  moves out of radio contact and connectivity is lost in Fig. 7. The forwarding router node,  $M_1$ , interacts with the routing protocol and forwards packets along a new route. The signaling system at intermediate router  $M_4$  receives packets and inspects its flow soft-state table. If a reservation does not exist for newly arriving packets then the signaling module invokes admission control and attempts to allocate resources for the flow. Note that when a rerouted packet arrives at node  $M_3$  the forwarding engine detects that a reservation exists and treats the packet as any other packet with a reservation. In other words, the packets are routed back to the existing path, where a reservation is still present. Such scenarios are frequently observed in our experimental systems, discussed in Section 6, with the result of minimizing any service disruption due to rerouting. Soft-state timers ensure that the flow state is still intact at  $M_3$  and that state along the old path (i.e., mobile host  $M_2$ ) is removed in an efficient manner.

When an adaptive flow is rerouted to a node where resources are unavailable, the flow is degraded to best effort service. Subsequently, downstream nodes receiving these degraded packets do not attempt to allocate resources or refresh the reservation state associated with the flow. In this instance the state associated with a flow

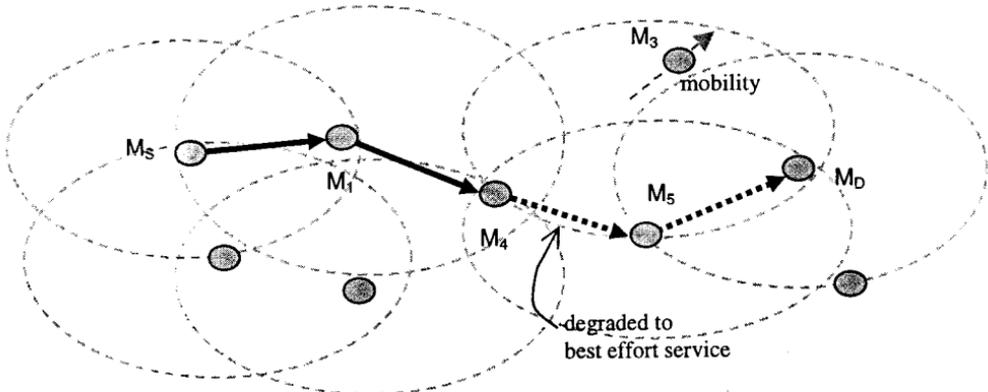


FIG. 8. Rerouting and degradation.

is timed out and resources are deallocated. A reservation may be restored if the resources free up at a bottleneck node (e.g., mobile node  $M_4$  in Fig. 8) or further rerouting may allow the restoration to complete. We call this type of restoration “degraded restoration.” A flow may remain degraded for the duration of the session and never be restored; this is described as “permanent degradation.” The enhanced QoS component of an adaptive flow may be degraded to best effort service (i.e., min-reserved mode) during the flow restoration process if the nodes along the new path can only support the minimum bandwidth requirement. If the degradation of enhanced QoS packets persist, it may cause service disruption and trigger the destination mobile node to invoke its adaptation procedure to “scale down” or “drop” packets rather than live with degraded quality. Adaptation mechanisms located at destination nodes are capable of responding to changes in network resource availability through scale down, scale up, and drop actions in response to network conditions.

During the restoration process, the INSIGNIA framework does not favor rerouted flows over existing flows (e.g., by forcing existing flows to scale down to their minimum requirements to allow rerouted or new flows to be admitted). In this sense, INSIGNIA avoids the introduction of additional service fluctuations to existing flows in support of the restoration of rerouted flows. As a result of this policy, admission control simply rejects/scales down any rerouted flows when insufficient resources are available along a new path.

Three types of restoration are supported by the INSIGNIA QoS framework:

- An *immediate restoration* occurs when a rerouted flow immediately recovers its original reservation; that is, a max-reserved mode flow is immediately restored as a max-reserved mode flow and a min-reserved mode flow as a min-reserved mode flow.
- A *degraded restoration* occurs when a rerouted flow is degraded for a period ( $T$ ) before it recovers its original reservation. Two forms of degraded restoration can occur: (i) a max-reserved mode flow operates at min-reserved mode and/or best effort mode and eventually recovers its original max-reserved mode service after some interval; (ii) a min-reserved mode flow operates at best effort mode and eventually recovers its original min-reserved mode service after some interval.

- A *permanent degradation* occurs when the rerouted flow never recovers its original reservation.

Figure 8 illustrates the topology changes that occur after rerouting based on the initial topology shown in Fig. 7. After rerouting link  $M_4$ – $M_5$  can only support best effort services. This type of restoration represents either a degraded restoration or a permanent degradation. In this scenario the destination node clears the partial reservation between mobile nodes  $M_5$ – $M_4$  by issuing a drop adaptation command to the source. The process of restoration can be immediate or delayed. Adaptation is application specific where the application can choose to respond to the network conditions and the delivered QOS.

**5.2.5. Adaptation.** The INSIGNIA QOS framework actively monitors network dynamics and adapts flows in response to observed changes based on user-supplied adaptation policy. Flow reception quality is monitored at the destination node and based on application-specific adaptation policy, actions are taken to adapt flows under certain observed conditions. Action taken is conditional on what is programmed into the adaptation policy by the user. For example, an adaptation policy could be to maintain the service level under degraded conditions or scale down adaptive flows to their base QOS in response to degraded conditions; other policy aspects could be to always scale up adaptive flows whenever resources are available. The application is free to program its own adaptation policy, which is executed by INSIGNIA through interaction of the destination and source nodes.

INSIGNIA provides a set of adaptation levels that can be selected. Typically, an adaptive flow operates with both its base and enhanced components being transported with resource reservation. Scaling flows down depends on the adaptation policy selected. A flow can be scaled down to its base QOS delivering enhanced QOS packets in a best effort mode, hence releasing any partial reservation that may exist. On the other hand, the destination can issue a drop command to the source to drop enhanced QOS packets (i.e., the source stops transmitting enhanced QOS packets). Further levels of scaling can force the base and enhanced QOS packets to be fully transported in best effort mode. In both cases, the time scale over which the adaptation actions occur is dependent on the application itself. These scaling actions could be instantaneous or based on a low pass filter operation [19].

During restoration of the flow state, admission control and resource reservation are invoked. This can lead to changes in a flow's observed quality at the destination mode in terms of having to scale down flows in response to observed resource bottlenecks along the new path or of having to scale up flows when additional resources are made available along the new path.

The INSIGNIA signaling system supports three adaptation commands that are sent from the destination host to the source using QOS reports:

- A *scale-down command* requests a source node to send its enhanced QOS packets as best effort or its enhanced QOS and base QOS as best effort.
- A *drop command* requests a source node to drop its enhanced QOS packets or enhanced and base QOS packets (where the term “drop” means the source node stop transmitting these packets).

- A *scale-up command* requests a source node to initiate a reservation for its base and/or enhanced service quality.

The scale down, drop, and scale up actions are driven by adaptation policy implemented at the destination, as illustrated in Fig. 9. Note that preference is given to base over enhanced QOS components in the event reserved packets have to be degraded to the best effort mode at bottleneck nodes, as illustrated in the figure. The scale down command is issued when the degradation of enhanced QOS packets persists. This action forces source nodes to send the enhanced QOS packets as best effort packets, thereby effectively removing any partial reservations that may exist, as illustrated in Fig. 9. A drop command is issued only when a destination node determines that degraded packets render insufficient quality. It is up to the applications to decide whether the reception of degraded packets is acceptable and take appropriate action. An adaptation policy handler at the destination is free to issue scale down commands or, in the case of persistent degradation (possibly including best effort delivery of both the base and enhanced QOS components), to terminate the session.

Mobility results in the release of resources along old paths and session dynamics result in additional resources becoming available along existing paths when sessions terminate. These released resources help other source-destination pairs support higher levels of quality for their sessions assuming they share a common path with the released resources. In such a case, the signaling system sets the bandwidth indication in the packet's INSIGNIA IP option field to indicate to adaptation handlers (located at the receiver) that sufficient resources may be available to support the delivery of base and enhanced QOS. The signaling system uses the bandwidth indication field to inform the destination host of the availability of new network resources should they become available along an existing path. Bottleneck nodes set the bandwidth indicator to MIN when enhanced QOS packets are scaled back in response to degraded conditions. Since each packet carries the max-min bandwidth requirements of each flow, bottleneck nodes can update a packet's bandwidth indicator in the event that resources become available to meet enhanced QOS needs. If all nodes along a path have resources to support enhanced QOS, then the bandwidth indicator received at the destination will indicate MAX in the bandwidth indicator field. This does not imply that a reservation has been made or that a reservation could be made with 100% assurance. Rather, it indicates to the source node that a reservation may be possible and that at the time the bandwidth

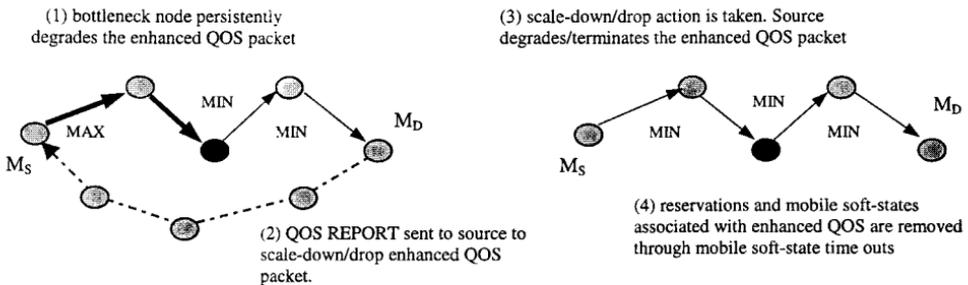


FIG. 9. Flow adaptation.

indicator bit was set resources were available. To initiate the reservation for the enhanced QOS, adaptation handlers send scale-up commands to their respective source nodes. In this sense the bandwidth indicator represents a good resource hint that additional service quality is possible. All messaging between source–destination pairs in support of scaling or dropping flow components is achieved using QOS reports.

## 6. EVALUATION

In what follows, we present the evaluation of the INSIGNIA QOS framework through simulations, with emphasis on the performance of the signaling system. The goal of the simulations is to evaluate the suitability of the INSIGNIA to support adaptive flows in a mobile ad hoc network under various traffic, mobility, and channel conditions. In particular, we are interested in evaluating systemwide restoration and adaptation dynamics and the impact of soft-state mechanisms and mobility on end-to-end sessions.

### 6.1. Simulation Environment

The INSIGNIA simulator consists of 19 ad hoc nodes, as illustrated in Fig. 10. Each mobile node has a transmission range of 50 m and shares a 2 Mbps air interface between neighboring mobile nodes within that transmission range. Time-varying wireless connectivity between nodes is modeled using 42 links. The mobility model is based on link failure and recovery characteristics defined in [11]; that is, connectivity is randomly removed and recovered with an arbitrary exponential distribution. Typically, mobile ad hoc networks do not have full connectivity between all mobile nodes at any given time due to the mobility behavior of mobile nodes and time-varying wireless link characteristics. With this in mind, maximum network connectivity is set at 85%, so that 15% of the mobile nodes within their transmission range remain disconnected.

We discuss the implementation of our INSIGNIA QOS framework where the generic MANET routing protocol used is based on an implementation of the TORA [1].

The QOS architectural components implemented in our simulator include the following:

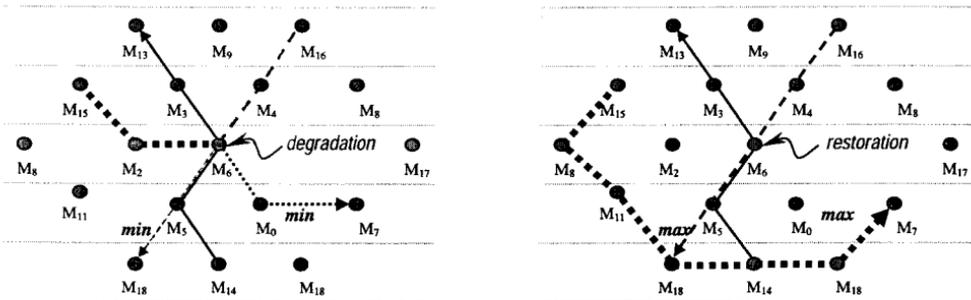


FIG. 10. (a) Degradation due to lack of resources. (b) Restorations through rerouting of a flow.

- The TORA [1] provided by the Naval Research Lab is used as a generic MANET routing protocol. The INSIGNIA framework is designed to “plug in” any MANET routing protocol.
- A packet scheduler is based on a deficit round robin implementation [29].
- An admission controller is simply based on peak allocation of bandwidth.

For simulation purposes 10 adaptive flows with different bandwidth requirements ranging from 75 to 500 kbps are operational throughout the simulation. An arbitrary number of best effort flows are randomly generated to introduce different loading conditions distributed randomly throughout the network (i.e., in different parts of the networks) during the simulation. We also chose an arbitrary traffic pattern/load with an average packet size of 2 Kbytes. Identical traffic/loads are used for all scenarios under investigation. The base QOS component of adaptive flows corresponds to 50–70% of an adaptive flow’s bandwidth needs, whereas enhanced QOS corresponds to 30–50%. For example, an adaptive flow of 300 kbps operating between nodes  $M_{14}$ – $M_{13}$  (as illustrated in Fig. 10) has 150 kbps for both its base and enhanced QOS such that minimum and maximum requirement is set to 150 and 300 kbps, respectively.

The mobility model used throughout the simulations supports three different rates of mobility. Moderate mobility represents slow vehicular mobility ranging from 9–18 km/h. Mobility conditions slower than moderate mobility are defined as slow mobility (i.e., speed less than 9 km/h) while rates faster than moderate mobility models are categorized as fast mobility (i.e., speed exceeding 18 km/h). We inherit the mobility model that was used in the TORA simulation [27]. In this simulation, we adopted a simple model for a mobility pattern [27] that abstracts the mobility and wireless link characteristics into link failure and link recovery characteristics. A shortcoming of this approach is that mobile nodes have a fixed set of neighboring mobile nodes, limiting the set of possible neighbors to communicate with. Therefore, the relative—and not absolute—mobility of the nodes is modeled. For the purpose of evaluating our framework, we measure per-session and aggregate network conditions for a number of experiments that analyze restoration, adaptation, soft-state management, and host/router mobility. We observe throughput, delays, out-of-order sequence packets, packet loss, percentage of delivered degraded packets for different mobility rates, and systemwide configuration (e.g., changing soft-state timers). We are particularly interested in the percentage of reserved and degraded packets delivered to all the receivers. This metric represents the ability of our framework to deliver assurance in mobile ad hoc networks. We also observe the number of rerouting, degradation, restoration, and adaptation events that took place during the course of each experiment as a measure of the dynamics of the system under evaluation.

## 6.2. Restoration Analysis

In the following experiment we investigate the impact of rerouting and restoration on adaptive flows. Since rerouting of flows requires admission control, resource allocation, state creation, and removal of old state we track the rerouting and

restoration events and any degradation that takes place. Typically, adaptive flows experience continuous rerouting during their session holding time. This is certainly the case for flows that represent continuous audio and video but not necessarily the case for microflows. These flows may be rerouted over new paths that have insufficient resources to maintain the required QOS. A key challenge for restoration is the speed at which flows can be restored. This is dependent on the speed at which new routes can be computed by the routing protocol if no alternative routes are cached and the speed at which the signaling system can restore reservations. The speed at which old reservations are removed is a direct function of the soft-state timer. The mobility rate impacts the number of restorations observed in the system and therefore the QOS delivered by the INSIGNIA QOS framework. As the rate of mobility increases (e.g., from moderate to fast), restoration algorithms need to be scalable and highly responsive to such dynamics in order to maintain end-to-end QOS.

In Section 5.2.4 we identified three types of restoration supported by the INSIGNIA model: immediate restoration, degraded restoration, and permanent degradation. Figures 10a and 10b illustrate the number of restorations and degradations that are associated with three randomly selected adaptive flows in our simulation. Due to the lack of resources at mobile node  $M_6$ , only flow  $M_{14}-M_{13}$  (i.e., the flow that traverses nodes  $M_{14}-M_{13}$ ) is transported in max-reserved mode, while flows  $M_{16}-M_{18}$  and  $M_{15}-M_7$  are transported in min-reserved mode. As a consequence, only the base QOS packets of flows  $M_{16}-M_{18}$  and  $M_{15}-M_7$  are delivered as reserved mode packets, while enhanced QOS packets are transported as degraded best effort packets. As illustrated in Fig. 10b, flow  $M_{16}-M_{18}$  transported in min-reserved mode regains its max-reserved service through the rerouting of flow  $M_{15}-M_7$ . Rerouting of flow  $M_{15}-M_7$  causes resources (i.e., 200 kbps) to be released by soft-state management. Consequently, this action allows mobile router  $M_6$  to restore the reservation requirement for the enhanced QOS of flow  $M_{16}-M_{18}$ , which requires 80 kbps. The rerouting of flow  $M_{15}-M_7$  finds sufficient resource availability on the new path (i.e.,  $M_{15}-M_8-M_{11}-M_{18}-M_{14}-M_{10}-M_7$ ), restoring its enhanced QOS.

Figures 11a and 11b illustrate immediate and degraded restorations observed under various mobility conditions. As indicated in the figures an increase in network dynamics increases the number of observed immediate and degraded restorations. The network experiences a total of 38 (61%) immediate restorations and 24 (39%) degraded restorations in the course of the simulation for a mobility rate of 3.6 km/h, as illustrated in Fig. 11a. As the mobility condition increases, the ratio between immediate restoration and degraded restoration changes. More immediate restorations are observed in comparison to degraded restorations for slow and moderate mobility conditions, as illustrated in Fig. 11b. However, when mobility conditions exceed 45 km/h, degraded restoration becomes dominant, as illustrated in Fig. 11b. The connectivity between mobile nodes becomes problematic as the mobility of nodes increases, causing the network topology to rapidly change. Consequently, the number of available routes between source and destination nodes diminishes and the contention for network resources increases. This phenomenon introduces service fluctuations and degradation. Figure 11 illustrates the different types of restoration

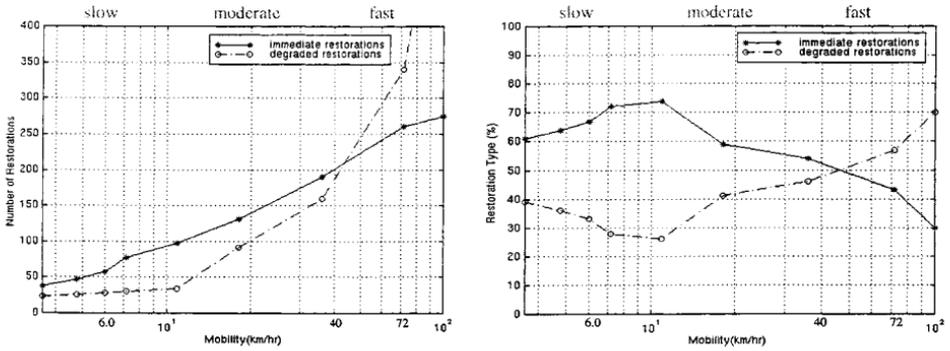


FIG. 11. (a) Number of restorations. (b) Percentage of restorations.

discussed in Section 5.2.4. Adaptive flows experience frequent rerouting with increased mobility, causing a rise in the number of observed degraded restorations.

The INSIGNIA framework adopts a simple admission control test that does not favor rerouted flow over existing flows. A rerouted flow is denied restoration along a new route when insufficient resources are available to meet its minimum bandwidth requirement. This approach minimizes any service disruptions to existing flows, preventing a wave of service fluctuation to propagate throughout the network. When a mobile host loses its connectivity to neighboring nodes due to mobility, reservations along the old path are automatically removed. In the case of degraded restoration or permanent degradation, flows are degraded to min-reserved mode or best effort mode because of the lack of resources to restore the flows during rerouting. We observed that max-reserved adaptive flows are more likely to be degraded to best effort service than are min-reserved mode adaptive service. This is mainly due to the admission control policy adopted and semantics of base QoS and enhanced QoS components of flows where the base QoS of a typical adaptive flow consists of 50–70% of the overall bandwidth needs. The admission controller will attempt to support the base and enhanced bandwidth needs of flows. This leads to a situation where most mobile nodes mainly support max-reserved mode flows and a few min-reserved mode flows to fill the remaining unallocated bandwidth. This leads to the blocking of max-reserved flows, and due to this behavior the vast majority of degraded flows are max-reserved to best effort. Therefore, degraded restorations of best effort to min-reserved (meaning that the min-reserved flow is degraded to best effort before being restored to min-reserved) only occur when the rerouted adaptive flows encounter resources to support only min-reserved service. We observed that degraded restoration for best effort to max-reserved (meaning that the max-reserved flow is degraded to min-reserved and/or best effort before being restored to max-reserved) is the most dominant degraded restoration type observed, as shown in Fig. 12. This is because rerouted flows are more likely to be accepted or denied rather than degraded to min-reserved flows under slow and moderate mobility conditions. However, we observe that when mobility exceeds 72 km/h, best effort to min-reserved degraded restoration becomes the dominant type, as shown in Fig. 12. In the case of high mobility, only a limited number of routes exist to route flows, causing service degradation. Rapid fluctuations in the monitored QoS cause the adaptation processes at the destination to

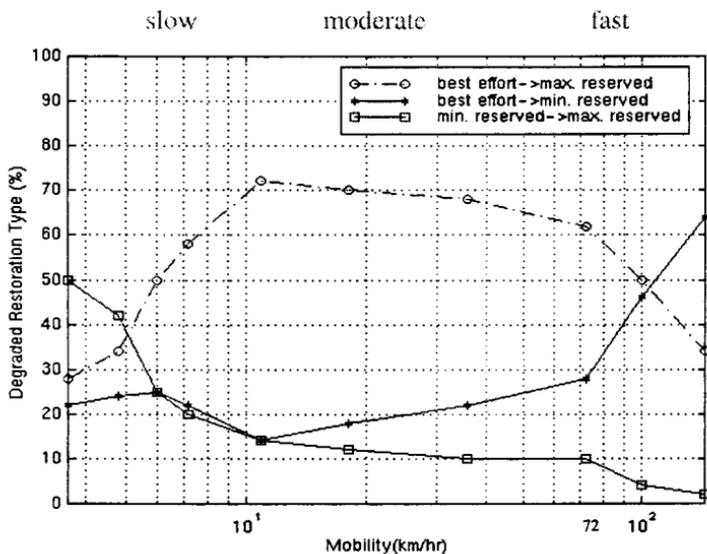


FIG. 12. Degraded restorations types.

request that the degraded flows be scaled down to their min-reserved mode. In this instance, the best effort to min-reserved restoration becomes the dominant type, as shown in Fig. 12.

Increased mobility forces mobile hosts to adapt flows to their min-reserved modes preventing adaptive flows from scaling back up due to the fast time scale dynamics and rerouting observed. When the mobility exceeded 72 km/h, all adaptive flows are scaled down to their min-reserved service 90 s into the trace. Only two scale up adaptations actions were observed during the complete trace. The number of best effort to max-reserved and min-reserved to max-reserved degraded restoration types decrease as mobility is increased beyond 72 km/h, as shown in Fig. 12. The best effort to min-reserved degraded restoration continues to increase, implying that most of the flows scale down to their minimum requirements and operate at the min-reserved mode.

Figure 13 shows the restoration times across the complete mobility range. The base QOS restoration time corresponds to the time taken to regain the min-reserved service for a flow that has been temporarily degraded to a best effort mode service. The enhanced QOS restoration time corresponds to the time taken for the max-reserved service to restore from the best effort service or from min-reserved service. We observe that the average required restoration time for immediate restoration is relatively constant at 0.2–0.9 s under all mobility conditions. We observe that immediate restoration only requires an interval of two consecutive packets to restore the reservation. However, mobility conditions impact the average degraded restoration time, unlike the immediate restorations, as shown in Fig. 13.

### 6.3. Adaptation Analysis

The adaptation process operates on an end-to-end basis and is driven by the observed service quality and adaptation policy of the destination node. This is in

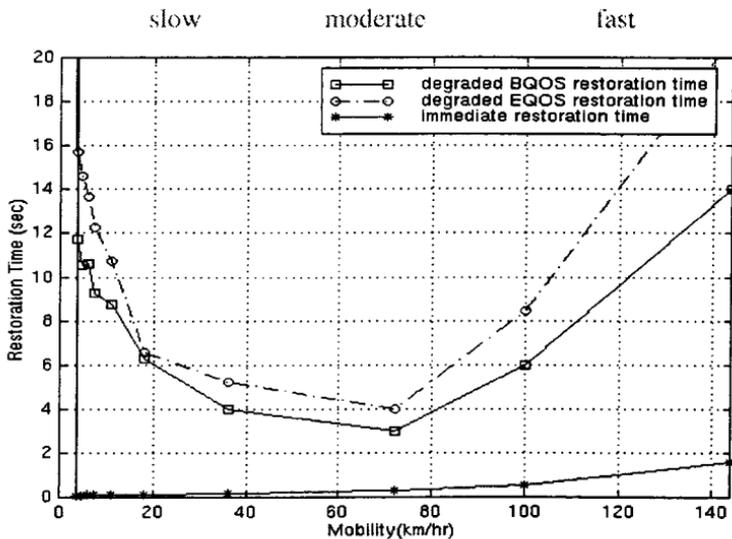


FIG. 13. Time spent for immediate restorations and degraded restorations.

contrast to restoration, which operates on the rerouting time scale. Typically, adaptation operates over longer time scales associated with end-to-end applications and their adaptation strategies. Monitoring modules residing at destination nodes actively measure the delivered service quality. As discussed in Section 5.2.5, destination nodes can issue adaptation commands to source nodes using QOS reports to scale down, drop, and scale up flows. For example, when the degradation of enhanced QOS packets persists beyond an acceptable period, the destination can issue a scale down adaptation command to the source node, removing any partial reservations that may exist between the source host and the bottleneck host. The INSIGNIA system is also capable of scaling up flows (e.g., from a min-reserved to a max-reserved service). The bandwidth indicator plays a central role in the adaptation process, as discussed in Section 5.2.5.

To observe the dynamics associated with the adaptation process, two adaptive flows are arbitrarily chosen and their associated throughputs measured (at their destination nodes) over the course of the simulation. The simulation results reflect moderate mobility conditions of 11 km/h. Moderate mobility conditions were chosen because slow mobility lacks network dynamics and fast mobility rarely experiences end system-initiated adaptation due to the rapid fluctuations in resource availability.

The impact of the adaptation process, degradation, and restoration on flows  $M_{15}-M_7$  and  $M_{16}-M_{18}$  from the previous example is shown in Fig. 14. As shown in the trace, flow  $M_{16}-M_{18}$  is affected by network dynamics at 17 s into the trace. The mobility of the network forces flows to be rerouted and, due to lack of resources along the new path, causes flow  $M_{16}-M_{18}$  to degrade to the min-reserved service, as indicated by (1) in Fig. 14. The degradation of flow  $M_{16}-M_{18}$  enhanced QOS packets is restored at (2) in Fig. 14. Degradation of the base QOS at point (3) is observed at 160 s and it is preceded by degradation of enhanced QOS packets at 145 s into the trace. Due to persistent service disruption the destination node ( $M_{18}$ ) triggers the source node ( $M_{16}$ ) to scale down the flow at 151 s into the trace. The decision

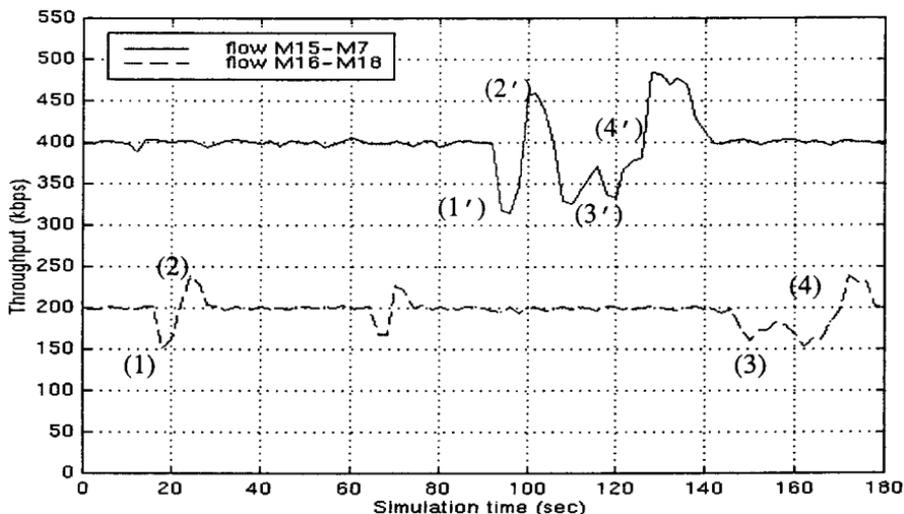


FIG. 14. Trace of adaptive INSIGNIA flows.

to scale down the flow is controlled by an adaptation handler. The source responds by transmitting the enhanced QOS packets as best effort packets. The reservations associated with the enhanced QOS packets is deallocated by soft-state management operating at intermediate routing nodes along the path, allowing other adaptive flows to scale up. Scaling up can be observed at  $t=172$  s into the trace when the destination node ( $M_{18}$ ) detects consistent resource availability through monitoring the bandwidth indicator. Flow  $M_{16}$ - $M_{18}$  restores its max-reserved mode service while flow  $M_{15}$ - $M_7$  first experiences degradation, scaling down and then scaling up. The degradation of the flow  $M_{15}$ - $M_7$  enhanced QOS packets degraded at  $t=92$  s is restored (2') to max-reserved mode service at  $t=98$  s into the trace. However, further network dynamics force the degradation of the enhanced QOS packet at  $t=100$  s into the simulation.

Adaptation policy is application specific in the sense that some flows prefer to instantly scale up when resources become available while others prefer not to follow not instantaneous changes but trends in resource availability. The scaling policy can be based on simple algorithms, for example, a simple state machine that scales flows down or up based on a certain number of degraded packets or packets indicating that additional resources are available, respectively. More sophisticated algorithms could follow statistical observations about network dynamics using low pass filters.

The rate of mobility has a large impact on the observed adaptation dynamics. Fewer instances of adaptation are observed given the same adaptation policy for slow mobility over moderate mobility. For mobility of 3.6 km/h we observe two scale up actions and one scale down action, whereas at 18 km/h we observe seven scale up and four scale down actions. As mobility increases beyond the moderate rate we observe more fluctuation in delivered service quality where scaling down flows to a min-reserved service becomes common. As the mobility speed increases to fast we observe few scaling up actions due to the fast dynamics of the network. Few destinations observe stable enough conditions to issue a scaling up command to their peer source nodes. For example, at 72 km/h we observe that only two scaling up

actions are recorded, with all adaptive flows being forced to scale down to their min-reserved mode during the course of the simulation.

#### 6.4. Soft-State Analysis

Soft-state resource management is used to maintain reservations. The duration of the soft-state timer has a major impact on the utilization of the network. Figure 15 shows the impact of soft-state times on network performance in terms of the number of reserved mode packets delivered. Reception of a reserved mode packet (with the service mode set to RES, as discussed in Section 5.1.1) at the destination indicates that the packet is delivered with max-reserved or min-reserved assurance. Reception of a packet degraded implies that the packet has been delivered without such guarantees. Therefore the percentages of reserved and degraded packets received by destination nodes as a whole indicate the degree of service assurance that an INSIGNIA network can support for different values of soft-state timers.

In what follows, we discuss the impact of soft-state timers on network performance. We set the soft-state timer value in the range of 0.01 to 30 s and observe the corresponding system performance. For each experiment we set the same timer value at each node. As shown in Fig. 15, the mobile soft-state timer value has an impact on the overall network performance. The ability to support adaptive services decreases as the soft-state timer value increases. The percentage of delivered reserved packets decreases as mobile soft-state timer increases. The percentage of degraded packets increases as the soft-state timer value increases, as shown in Fig. 15. Worst case performance is observed when the soft-state timer value is set to 30 s. In contrast, the best performance is observed when soft-state timer is set to 2 s, as shown in Fig. 15. We observed that 69% of the packets are delivered as reserved packets and 31% as best effort packets when the soft-state timer is set at 30 s. Support for QOS substantially improves with 88% of reserved packets being delivered to the receivers with a soft-state timer value of 2 s. Large timeout values tend to lead to under utilization of the network because resources are “locked up” with resources remain allocated long after flows have been rerouted. New flows are

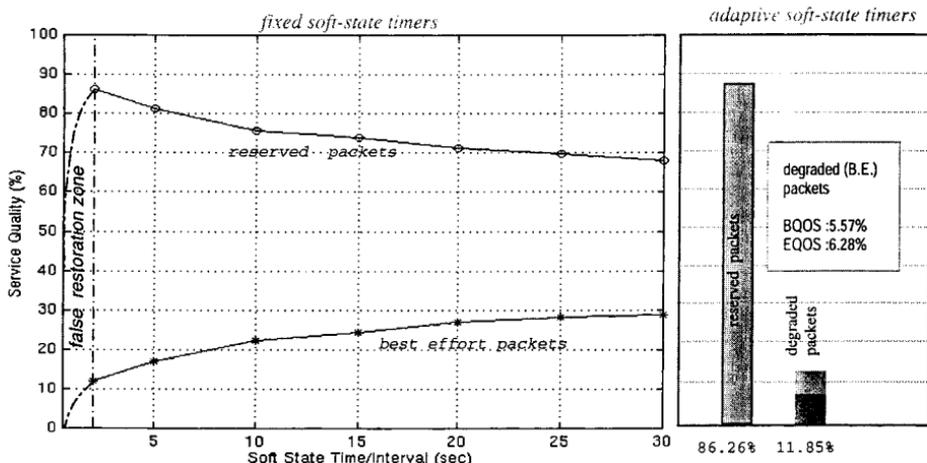


FIG. 15. Soft-state timers and network performance.

unable to use these dormant resources, resulting in the overall degradation of the network due to “resource lockup.”

As the value of the soft-state timer gets smaller fewer resource lock ups are observed and utilization increases. However, when the timer is set to a value smaller than 2 s the network experiences what we describe as “false restoration.” This occurs when a reservation is prematurely removed because of a small soft-state timer value. However, this is a false state because the session holding time is still active and the source node keeps sending packets. In this case, the reservation is removed because of a timeout and then immediately reinstated when the next reserved packet arrives. False restorations occur when the timeout value is smaller than the interarrival time between two consecutive packets associated with a flow. With a soft-state timer of 0.04 s, for example, all the adaptive flows experienced numerous false restorations. Mobile routers often deallocate and reallocate resources without the involvement of any network dynamics due to mobility. In the worst case, every packet can experience a false restoration. Such events not only increase the processing costs of state creation and removal, and resource allocation and deallocation, but also falsely reflect the resource utilization and availability of the system. When the network experiences numerous false restorations, rerouted flows often find nodes with few resources allocated on the new path. This phenomenon causes flows to always gain max-reserved mode resources with mobile nodes accepting the request for resources well beyond their actual capacity. This results in reserved packets experiencing indefinite delays at intermediate nodes even though resource assurances are provided by admission controller, resulting in wide scale packet losses and service degradation. Figure 15 shows a “false restoration region” where there is little distinction between reserved and best effort operational modes and where reservations are typically always granted. Adaptation and restoration algorithms can fail under false restoration conditions due to the perception of unlimited resource availability. Setting a suitable soft-state timer value is therefore essential to preventing both false restoration and resource lockup in our framework.

Each data packet associated with a reserved flow is used to refresh soft-state reservations. We observe that different adaptive flows have different data rates, and thus a fixed timeout value is too limiting. For example, one value may be fine for some set of flows but cause false restorations or resource lockup for others. Clearly there needs to be a methodology for determining the value of the soft-state timer. The issues of false restoration and resource lockup can only be resolved by adjusting the timeout value based on the observed flow dynamics. The timeout should be based on the effective data rate of each flow. More specifically, the soft-state timer should be based on the measured packet interarrival rate of adaptive flows. The signaling system measures packet interarrivals and jitter at each mobile node for each flow, adjusting the soft-state timeout accordingly. In the experimental system we implemented an *adaptive soft-state timer* that is initially set to 4 s, representing an initial safety factor. This allows mobile nodes to set their soft-state timers according to their effective data rate, allowing the timeout to adjust to network dynamics and the variation in the interarrival rates of individual flows traversing nodes. The implementation of an adaptive soft-state timeout effectively removes resource lockups and false restorations, as shown in Fig. 15. We observe that when an adaptive

soft-state timer scheme is used 86% of flows are delivered as reserved packets and 11% as degraded packets. Adaptive soft-state timers greatly reduce resource lockup and false restoration conditions, allowing the network to support better service assurances through the delivery of more reserved packets and fewer degraded packets at destination nodes.

### 6.5. Mobility Analysis

To evaluate the impact of mobility on the INSIGNIA QOS framework, we conduct a set of experiments operating under identical traffic patterns/load conditions and various mobility conditions ranging from 0 to 72 km/h. Figure 16 illustrates the impact of mobility on the delivered service quality. When there is no host mobility, results closely approximate a fixed network infrastructure where admitted flows receive stable QOS assurances. One anomaly is observed, however. Six adaptive flows failed to be granted reservations due to a lack of network resources at intermediate nodes. As a consequence only 49% of the packets are delivered as reserved packets and 51% as best effort packets. This anomaly is a product of the routing protocol, which provides a non-QOS routing solution. Adaptive flows are routed to bottleneck nodes, resulting in the failure of admission control due to the lack of resources. This problem could be resolved by designing a signaling system that takes alternative routes when admission control fails along a selected path.

With the introduction of mobility into the network, the performance improves (i.e., more reserved packets are delivered) as illustrated in Fig. 16. Mobility-induced rerouting allows request packets to traverse alternative paths, increasing the probability of finding a route with sufficient resource availability to admitted flows as reserved mode packets. Figure 16 shows that INSIGNIA supports relatively constant QOS under slow and moderate mobility conditions between 3.6 and 18 km/h. The optimal performance is observed when the average network mobility is approximately 11 km/h. This results in the delivery of 86% of reserved packets. The in-band

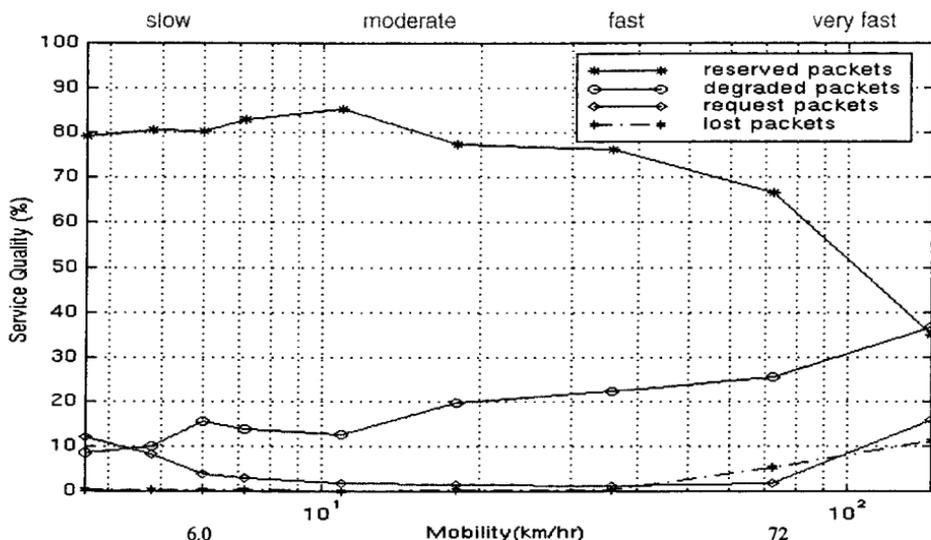


FIG. 16. Mobility and network performance.

nature of INSIGNIA allows the system to cope with fast network dynamics in a responsive manner. In an ideal case, INSIGNIA requires only a single packet reception to set up and restore (i.e., immediate restorations) reservations for the new and rerouted flows, respectively. INSIGNIA supports the delivery of 66% reserved packets even when mobiles are moving at 72 km/h, as shown in Fig. 16. This is a very encouraging result.

Note that the service provided in a mobile ad hoc network has a memoryless property such that adaptive flows require new admission tests along the new path when rerouting occurs. This implies that an increase in mobility may cause fluctuations in the perceived service quality. At 72 km/h all flows are scaled down to min-reserved packets after 90 s into the simulation due to the fluctuations in delivered quality. At this speed only two flows are capable of regaining their max-reserved service. When mobility conditions exceed 72 km/h, support for QOS breaks down rapidly as indicated in Fig. 16. The mobility characteristics overload the system and service assurance for adaptive flows diminishes. In fact, when mobility exceeds 90 km/h, we observe that flows  $M_{12}-M_{11}$ ,  $M_3-M_7$ , and  $M_5-M_{12}$  are transported as best effort packets for more than 70 s because they failed to accomplish their end-to-end flow setup due to persistent loss of RES packets and QOS reports. This phenomenon corresponds to the abrupt loss of reserved packets and degraded packets.

An increase in out-of-sequence packet is also observed at higher speeds, possibly causing service disruption at the receiver. Figure 17 shows the number of out-of-sequence packets under various mobility conditions. The number of out-of-sequence packets generally increases as mobility increases. The number of delivered out-of-sequence packets is influenced by different propagation delay characteristics of reserved and best effort packets associated with the same end-to-end flow. Figure 17 also shows the number of lost packets observed under different mobility conditions. Packets that are delayed for more than 15 s are discarded at intermediate nodes and considered lost. Figure 18 shows the delay characteristics of packets under various mobility conditions. When mobility increases, the connectivity between

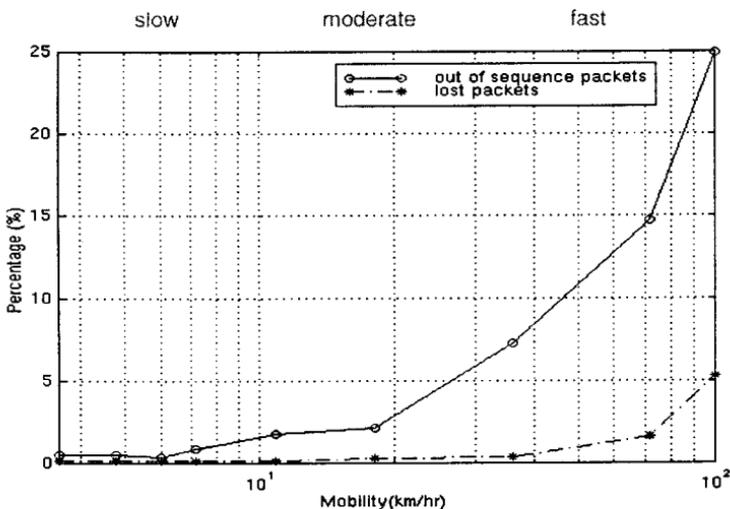


FIG. 17. Impact of mobility on out of order delivery and packet loss.

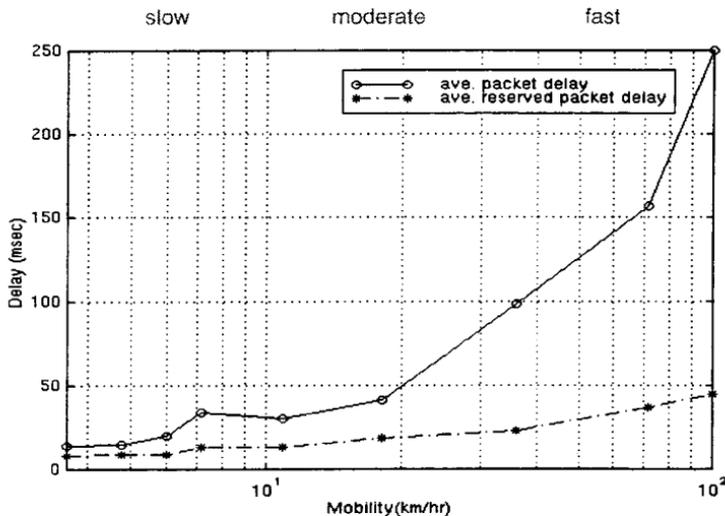


FIG. 18. Packet delays.

nodes becomes problematic. Such network dynamics trigger frequent routing updates and decreased connectivity. Thus, the number of available routes between nodes decreases as mobility increases. Degraded packets queue up at intermediate nodes experiencing long delays. However, the reserved packets are less sensitive to these delays, as indicated in Fig. 18, with all reserved packets being delivered within a period of 40 ms.

## 7. CONCLUSION

In this paper, we have presented the design, implementation, and evaluation of the INSIGNIA QOS framework that supports the delivery of adaptive services in mobile ad hoc networks. A key contribution of our framework is the INSIGNIA signaling system, an in-band signaling system that supports fast reservation, restoration, and adaptation algorithms. The signaling system is designed to be lightweight and highly responsive to changes in network topology, node connectivity, and end-to-end quality of service conditions. We have evaluated our QOS framework paying particular attention to the performance of the signaling system.

The approach discussed in this paper looks promising in terms of the performance results presented. Our simulation results show the benefit of our framework under diverse mobility, traffic, and channel conditions. The use of in-band signaling and soft-state resource management proved to be very efficient, robust, and scalable. Our results highlighted a number of anomalies that emerged during the evaluation phase. However, the use of adaptive soft-state timers seemed to resolve many of these issues (e.g., false restorations and resource lockups).

Currently we are implementing a number of MANET routing protocols using the NS2 simulator to investigate how well INSIGNIA performs in a heterogeneous

routing environment where reservation, restoration, and adaptation are required across multiple MANET routing domains. We are also building an experimental INSIGNIA test bed at Columbia University and intend to investigate how well our framework will operate in support of real-time applications. Results from this phase of our research will be the subject of a future publication.

## ACKNOWLEDGMENTS

This work is supported in part by the Army Research Office under Award DAAD19-99-1-0287 and with support from COMET Group industrial sponsors. The authors thank Javier Gomez Castellanos and Raymond Rui-Feng Liao for their comments on this work. We are particularly indebted to Vicent Park (Naval Research Lab.) and Dr. M. Scott Corson (Institute for System Research, University of Maryland) for providing us with the source code for TORA.

## REFERENCES

1. V. Park and S. Corson, Temporally ordered routing algorithm (TORA) version 1 functional specification, draft-ietf-manet-tora-spec-00.txt, work in progress, November 1997.
2. J. Macker and M. S. Corson, Mobile ad hoc networking (MANET): Routing protocol performance issues and evaluation considerations, draft-ietf-manet-issues-01.txt, work in progress, April 1998.
3. D. D. Clark and D. L. Tennenhouse, Architectural consideration for a new generation of protocols, in "Proceedings, ACM SIGCOMM'90," August 1990.
4. M. Gerla and J. T.-C. Tsai, Multicluster, mobile. Multimedia Radio Network, *Wireless Networks* 1 (3) (1995).
5. Z. Haas and M. Pearlman, The zone routing protocol (ZRP) for ad hoc networks, draft-ietf-manet-zone-zrp-00.txt work in progress.
6. C. Perkins, Ad hoc on demand distance vector (AODV) routing, draft-ietf-manet-aodv-01.txt, work in progress.
7. D. B. Johnson and D. A. Maltz, Dynamic source routing in ad hoc wireless network, in "Mobile Computing," Chap. 5, pp. 153–181.
8. M. S. Corson, Issues in supporting quality of service in mobile ad hoc networks, in "Proceedings, IFIP Fifth International Workshop on Quality of Service (IWQOS '97)," Columbia University, 1997.
9. C. R. Lin and M. Gerla, A distributed architecture for multimedia in a multihop dynamic packet radio network, in "Proceedings, IEEE Globecom'95," pp. 1468–1472, November 1995.
10. V. Park and M. S. Corson, A highly adaptive distributed routing algorithm for mobile wireless networks, in "Proceedings, IEEE INFOCOM '97," April 1997.
11. V. Park and M. S. Corson, A performance comparison on the temporally-ordered-routing algorithm and ideal link-state routing, in "Proceedings, IEEE Symposium on Computers and Communication '98," Athens, Greece, June 1998.
12. W. Almesberger, T. Ferrari, and J. Le Boudec, SRP: A scalable resource reservation protocol for the internet, available at <http://lrcwww.epfl.ch/srp/>.
13. R. Ramanathan and M. Streenstrup, Hierarchically-organized, multi-hop mobile wireless networks for quality-of-service support, available at <ftp://ftp.bbn.com/pub/ramanath/mmwn-paper.ps>.
14. C. R. Lin and M. Gerla, Asynchronous multimedia multihop wireless networks, in "Proceedings, IEEE INFOCOM'97," April 1997.
15. R. Braden, L. Zhang, S. Berson, S. Herzog, and S. Jamin, Resource reservation protocol (RSVP), RFC 2205, September 1997.

16. P. Sharma, D. Estrin, S. Floyd, and V. Jacobson, Scalable timers for soft-state protocols, in "Proceedings, IEEE INFOCOM'97," April 1997.
17. P. Ferguson, Simple differential services: IP TOS and precedence, delay indication and drop preferences, draft-ferguson-delay-drop-00.txt work in progress.
18. M. S. Corson and V. Park, An internet MANET encapsulation protocol (IMEP) specification, internet draft, draft-ietf-amnet-imep-spec-01.txt, work in progress, November 1997.
19. R. R.-F. Liao and A. T. Campbell, On programmable universal mobile channels in a cellular internet, in "Proceedings, 4th ACM/IEEE International Conference on Mobile Computing and Networking (MOBICOM'98)," Dallas, October 1998.
20. M. S. Corson and A. T. Campbell, Toward supporting quality of service in mobile ad hoc networks, in "Proceedings, First Conference on Open Architecture and Network Programming," San Francisco, April 3-4, 1998.
21. J. Broch, D. A. Maltz, D. B. Johnson, Y.-C. Hu, and J. Jetcheva, A performance comparison of multi-hop wireless ad hoc network routing protocols, in "Proceedings, 4th Annual ACM/IEEE International Conference on Mobile Computing and Networking," ACM, Dallas, TX, October 1998.
22. S. Lu, V. Bharghavan, and R. Srikant, Fair scheduling in wireless packet networks, in "Proceedings, ACM SIGCOMM'97," San Francisco, 1997.
23. Global mobile information systems program, available at <http://www.darpa.mil/ito/research/glomo/index.html>.
24. H. Schulzrinne *et al.*, RTP: A transport protocol for real-time applications, RFC, 1989.
25. M. Shreedhar and G. Varghese, Efficient fair queuing using deficit round robin, in "Proceedings, ACM SIGCOMM'95," Berkeley, 1995.
26. A. Balachandran, A. T. Campbell, and M. E. Kounavis, Active filters: Delivering scaled media to mobile devices, in "Proceedings, NOSSDAV'97," St. Louis, 1997.
27. TORA OPNET source code supplied by V. Park and M. S. Corson, 1997.
28. O. Angin, A. T. Campbell, M. E. Kounavis, and R. R.-F. Liao, The mobiware toolkit: Programmable support for adaptive mobile networking, *IEEE Personal Communications Magazine, Special Issue on Adaptive Mobile Systems*, August 1998.
29. J. Gomez, A. T. Campbell, and H. Morikawa, Havana: Supporting application and channel dependent QOS in wireless networks, in "Proceedings, 7th International Conference on Network Protocols, Toronto, 1999."
30. S.-B. Lee and A. T. Campbell, INSIGNIA, internet draft, draft-lee-insignia-00.txt, work in progress, November 1998.
31. S.-B. Lee and A. T. Campbell, INSIGNIA: In-band signaling support for QOS in mobile ad hoc networks, in "Proceedings, 5th International Workshop on Mobile Multimedia Communications (MoMuC'98)," Berlin, October 1998.
32. A. Ephremides and T. Truong, Scheduling algorithms for multi-hop radio networks, *IEEE Trans. Comput.* **38** (1989), 1353.
33. OPNET, available at <http://www.mil3.com>.
34. Secure protocols for adaptive, robust, reliable, and opportunistic WINGs (SPARROW) project, available at <http://www.cse.ucsc.edu/research/ccrg/projects/sparrow.html>.
35. P. Sinha, R. Sivakumar, and V. Bharghavan, CEDAR: A core-extraction distributed ad hoc routing algorithm, in "IEEE Infocom'99," New York, March 1999.
36. Defense advanced research projects agency (DARPA), available at <http://www.darpa.mil>.