# PROCEEDINGS OF SPIE

SPIEDigitalLibrary.org/conference-proceedings-of-spie

Cross-layer protocols optimized for real-time multimedia services in energy-constrained mobile ad hoc networks

William S. Hortos

# Cross-layer protocols optimized for real-time multimedia services in energy-constrained mobile ad hoc networks

William S. Hortos

Florida Institute of Technology, Orlando Graduate Center, 3165 McCrory Place, Suite 161, Orlando, FL 32803

## ABSTRACT

Mobile ad hoc networking (MANET) supports self-organizing, mobile infrastructures and enables an autonomous network of mobile nodes that can operate without a wired backbone. Ad hoc networks are characterized by multihop, wireless connectivity via packet radios and by the need for efficient dynamic protocols. All routers are mobile and can establish connectivity with other nodes only when they are within transmission range. Importantly, ad hoc wireless nodes are resource-constrained, having limited processing, memory, and battery capacity. Delivery of high quality-of-service (QoS), real-time multimedia services from Internet-based applications over a MANET is a challenge not yet achieved by proposed Internet Engineering Task Force (IETF) ad hoc network protocols in terms of standard performance metrics such as end-to-end throughput, packet error rate, and delay.

In the distributed operations of route discovery and maintenance, strong interaction occurs across MANET protocol layers, in particular, the physical, media access control (MAC), network, and application layers. The QoS requirements are specified for the service classes by the application layer. The cross-layer design must also satisfy the battery-limited energy constraints, by minimizing the distributed power consumption at the nodes and of selected routes. Interactions across the layers are modeled in terms of the set of concatenated design parameters including associated energy costs. Functional dependencies of the QoS metrics are described in terms of the concatenated control parameters.

New cross-layer designs are sought that optimize layer interdependencies to achieve the "best" QoS available in an energy-constrained, time-varying network. The protocol design, based on a reactive MANET protocol, adapts the provisioned QoS to dynamic network conditions and residual energy capacities. The cross-layer optimization is based on stochastic dynamic programming conditions derived from time-dependent models of MANET packet flows. Regulation of network behavior is modeled by the optimal control of the conditional rates of multivariate point processes (MVPPs); these rates depend on the concatenated control parameters through a change of probability measure. The MVPP models capture behavior of many service applications, e.g., voice, video and the self-similar behavior of Internet data sessions.

Performance verification of the cross-layer protocols, derived from the dynamic programming conditions, can be achieved by embedding the conditions in a reactive routing protocol for MANETs, in a simulation environment, such as the wireless extension of ns-2. A canonical MANET scenario consists of a distributed collection of battery-powered laptops or hand-held terminals, capable of hosting multimedia applications. Simulation details and performance tradeoffs, not presented, remain for a sequel to the paper.

Keywords: Cross-layer network protocol, mobile ad hoc network (MANET), real-time multimedia services, Internet protocols, quality of service (QoS), dynamic programming, application layer, MAC, multivariate point processes, energy constraints

## 1. INTRODUCTION

A mobile ad hoc network (MANET) consists of a collection of mobile nodes forming a dynamic autonomous network. Nodes communicate with each other over the wireless medium without the intervention of an infrastructure consisting of centralized access points (APs) or base stations (BSs). Hence, the nodes of the MANET form a fully mobile infrastructure; each node acts as both a router and a host. Due to the limited transmission range (and power reserves) of wireless network interfaces, multiple hops may be needed to exchange data packets between nodes in the network. The

advantage of such networks is that they can be deployed rapidly anywhere and anyplace without the presence of a fixed infrastructure of BSs and system administrators.

This paper extends previous development by this author of analytical models of the behavior of packet flows in wireless multimedia network, based on the theory of semi-martingale representations of multivariate point processes (MVPPs) with randomly modulated rates.[1] In that work, the models achieve a comprehensive representation of the controllable and observable, real-time call processing events. The MVPPs are used to represent the transient packet flows of information between nodes in a wireless packet-switched network. In the finite-population MANET model the set of nodes consists of a maximum $M_{\max}$ mobile subscribers (MSs).

The extension of the models in this paper is based on an application of Girsanov's theorem to control the MVPP rates through an absolutely continuous change of probability measure from a reference measure on network events. The controlled rates represent the mechanisms of packet entry, routing and buffering in the network, to support the end-to-end connectivity of the path between source and destination nodes. The counting processes on packet events as well as inter-event times are allowed to obey general non-stationary probability distribution functions (PDFs) for call arrivals, packet delays, multiple service loading, and service completions. Inter-event times are random stopping times, progressively measurable with respect to the $\sigma$–algebra generated by the history of network events. Random routing variables, which characterize packet admission, processing, collisions, as well as routing, depend on non-stationary path loss distributions due to fading, reflections and user mobility, traffic loading, as well as on the observable network history. The set of these random variables, which further satisfy formal mathematical assumptions, constitute the admissible class of controls that influence the MVPP models through a change of probability measure on network events.[1]

The extent of observed network history available to the routing control depends interaction of the higher layer functions of the MANET protocol. Limited or partial observations lead to the concept of stochastic filtration of the network state dynamics to effect "best-effort" routing to support the QoS requirements of the packet flows. Observations of network events can be partial, based on incomplete state information or incomplete in time; time can be either deterministic or random with a known PDF.

The MVPP models generalize Poisson point processes, exponentially distributed message processing, and other renewal processes commonly assumed to formulate queueing structures for the evaluation of the asymptotic performance of control schemes for packet-based networks at or near equilibrium. The inherent MANET dynamics, however, are random and transient, and rarely support an assumption of ergodicity or the existence of an equilibrium. Therefore, the models of MANET performance for real-time multimedia packet flows considered here are limited to a finite time interval, $[0,T], T < \infty$. The MVPP approach also overcomes limitations of conventional Markov and Bayesian models that assume successive observations are independent, thus enabling the probability of a sequence of observations to be written as the product of probabilities of individual observations. This assumption is clearly precluded by node mobility and competition for constrained network resources over the duration of a session.

The proposed MVPP models of real-time network events are not only mathematically tractable and extendible, but can encompass self-similar processes of long-range dependence (LRD) that characterize Internet traffic. The objective of the models is to provide an analytical basis for accurate evaluations of the comparative performance of adaptive cross-layer designs for MANET protocols that provide resource allocation, routing and congestion control in support of processing of real-time multimedia applications.

## 2. BACKGROUND

Due to the frequent changes in MANET topology and the lack of network resources both in the wireless medium and at the mobile nodes, the effective delivery of the real-time multimedia services, as offered in wired Internet protocol (IP) networks, becomes very challenging. Consequently, routing in MANETs for multimedia experiences link failure more often. Hence, a routing protocol that supports quality of service (QoS) requirements for ad hoc networks requires an understanding of the causes of link failure to improve its performance. Fundamentally, link failure is due to node mobility and lack of network resources, in particular, bandwidth, packet buffers, and battery power reserves. Therefore, it is essential to capture the aforementioned characteristics to identify the quality of the paths in the MANET. Moreover,

the routing protocols must be adaptive to utilize optimally the highly constrained, time-varying resources. For instance, it is possible that a route, earlier determined to satisfy certain QoS requirements and battery reserve constraints, may no longer do so due to the inherently dynamic nature of the network topology. It is thus important that the network intelligently adapts to these dynamic conditions in order to maintain service on the path.

According to the Internet Engineering Task Force (IETF) Request for Comment (RFC) 2386 on the delivery of QoS-based routing on the Internet, the term "quality of service" denotes the provision of a set of service requirements to the packet flows while routing them through the network.[2] For MANETs, with limited, time-varying resources, the ability to meet application-specific requirements, such as delay, is problematic, but remains a worthy objective for optimized protocol design. Hence, the QoS concept for wired IP networks may not be valid for MANETs; even with high-speed, fixed links, Internet service providers and administrators find it difficult to guarantee the QoS of end-to-end services.[3, 4]

The dynamic nature of the MANET makes routing and consequently QoS support in these networks fundamentally different from wired networks. Moreover, since network quality in terms of available resources, e.g., buffer space and battery states at the nodes, varies with time, present QoS models for wired networks are insufficient for such networks.[2, 5] A QoS model does not define specific protocols or implementations for a network, but rather the methodology and architecture by which certain service types can be provided. Integrated Services (IntServ) [4, 6] and Differentiated Services (DiffServ) [4, 7] are two models proposed to deliver QoS guarantees in the Internet.

The IntServ architecture allows source nodes to communicate their QoS requirements to routers and destinations on the path by means of signaling protocol, such as, RSVP.[8, 9] Hence, IntServ provides *per-flow* end-to-end QoS guarantees. IntServ defines two service classes: *guaranteed service* [10, 11] and *controlled load*,[12, 13] in addition to the *best-effort* service. Guaranteed service ensures a maximum end-to-end delay, and is intended for applications with strict delay requirements, such as voice. Controlled load, conversely, guarantees to provide a level of service equivalent to best-effort service in a lightly loaded network, regardless of network load. Controlled-load service is designed for adaptive, real-time applications. IntServ is not appropriate for MANETs, as the amount of state information increases with the number of packet flows, that is, the model is not scalable.

The DiffServ architecture avoids the problem of scalability by defining a small number of *per-hop* behaviors (PHBs) at the network edge routers and associating a different DiffServ Code Point (DSCP) in the IP header of packets belonging to each class of PHBs. Core routers use DSCPs to differentiate among different QoS classes on a per-hop basis. Hence, DiffServ is scalable, but it does not ensure service guarantees on an end-to-end basis. This is a disadvantage to the deployment of DiffServ in the Internet, and also is a disadvantage for MANETs. In DiffServ, there are three different service classes: *expedited forwarding*, *assured forwarding*, and *best effort*. Expedited forwarding provides a low-delay, low-loss rate, and assured bandwidth. Assured forwarding provides guaranteed or expected throughput for applications, while best effort offers no guarantees.

DiffServ and IntServ require accurate information on the link state, e.g., available bandwidth, packet loss rate, delay, and other features on a per-hop basis and on the topology. The transient, constrained resources of the MANET greatly inhibit the ability to maintain accurate routing information. However, a QoS model for MANETs should make some use of the concepts and features of existing models in order to construct a model that satisfies such networks.

A variant of these two architectures, called a Flexible QoS Model for MANET (FQMM) [5, 14] has been proposed for ad hoc networks. FQMM defines three type of roles for nodes: an ingress (source) mode which sends data, an interior node which forwards data to other nodes, and an egress (destination) mode. Each node may play multiple roles in the network. The FQMM selectively uses the per-flow state property of IntServ and the service differentiation of DiffServ. For applications with high priority, the per-flow QoS guarantees of IntServ are provided, while applications with lower priorities are given the per-class differentiation of DiffServ. In FQMM, both IntServ and DiffServ schemes are used separately for different priority classes. Due to its construction, the disadvantages related to IntServ and DiffServ are inherited by FQMM. Moreover, FQMM does not account for the characteristics of MANETs.

# 3. QUALITY OF SERVICE IN MANETS

Unlike wired networks, QoS support in MANETs depends not only on *available resources* in the network but also on the *mobility rate* or *transient rate* of such resources. Indeed, mobility and transient resources may cause link failures, and consequently broken paths and unwanted network partitioning. In addition, MANETs usually have less resources than fixed networks. Therefore, additional criteria are needed to characterize the quality of links between nodes. Applications must adapt to the transient, limited resources offered by the MANET. Thus, the QoS required by an application depends on the "quality" of the network. This "quality" should be a function of the availability and stability of the resources in both the wireless medium and at the mobile nodes. Hence, QoS in MANETs is viewed as the *provision of a set of resource parameter levels in order to adapt the different application packet streams, offered at sources, to the "quality" of the network,* while *routing them through the network.* Therefore, QoS routing with energy constraints is a mechanism under which paths are created in the MANET, based on the "quality" of the network state, and then selected according to the QoS requirements of the packet flows and to the residual power reserves of the paths and their constituent nodes. In energy-constrained optimization, network "quality" includes the state of the battery reserves at the nodes and on the paths. Hence, the multiple objectives of energy-efficient QoS protocols are to optimize resource utilization, including battery power, while satisfying service application requirements.

## 3.1 Cross-layer quality-of-service model

To meet the multiple objectives of the MANET protocol, a cross-layer approach has been proposed for the QoS model that groups the resource parameters and performance metrics associated with different protocol layers, that is, the application layer, network layer, link/medium access control (MAC) layer, and physical (PHY) layer, and then maps them accordingly.[6, 7, 15, 16] This approach is suggested since the QoS requirements of applications depend strictly on the "quality" of the network, i.e., on the available resources as well as the stability of those resources.

## 3.2 Service application classes and protocol layers

Service applications can be broadly classified as either real-time or non-real-time. Real-time services can have distinct constant bit rates (CBRs), such as 8-kilobits per second (kbps) and 13-kbps voice codecs, or variable bit rates (VBRs), as used in interactive video. Excessive delay or delay variation (jitter) noticeably degrades the quality of real-time services. Non-real-time services, such as file transfers, Internet accesses, e-mail and other delay insensitive services, are supported by available bit rates (ABRs) and transmitted by IP networks as high-rate bursts, characterized as "on-off" processes. For packet data services, transmission stops at the end of the data burst, since no information is generated during the unpredictable "off" intervals. Transmission of real-time services is continuously maintained during the session, while packet data services are provided to users with demand for high transmission rates, but short session duration times. Certain non-real-time packet data services differ in their tolerance of delay variation as opposed to fixed delay, e.g., many web pages include real-time video and audio clips.

While some QoS requirements are expressed in terms of the interrelated parameters of receive signal strength (RSS), signal-to-interference-plus-noise ratio (SINR), bit-error ratio (BER), and frame-error ratio (FER), others distinguish service applications by their bit-rate, i.e., CBR, VBR, ABR and undefined, requirements and their tolerance to fixed delay and delay variations. All service applications can be accommodated in the model.[17]

Class 1 services require real-time connections with very low delay-tolerance. Examples include voice, interactive video and video conferencing. Real-time multimedia applications, such as video conferencing, impose these requirements on inter-network gateways, since the traffic they generate must be delivered in a certain temporal sequence. Class 1 applications receive the highest service priority over other classes and require ensured bandwidth or transmission rates. Terms of service may be negotiated between two or more CBR alternatives based on bandwidth and other radio resource availability.

Class 2 services are non-real-time, delay-sensitive, and connection-oriented with limited delay requirements. Examples include MPEG-2 video, remote login, file transfer protocol (FTP), and similar applications associated with the transport control protocol (TCP). This class typically receives lower priority than Class 1. The transfer rate can be negotiated as a VBR between maximum and minimum acceptable limits, based on QoS latency requirements and resource availability and stability.

Class 3 services are message-oriented and delay-tolerant. Typical services are paging, e-mail, voice mail, facsimile, and data file transfer such as Web downloads. Class 3 applications can be conveyed at the earliest possible time and the transfer rate can be adjusted continuously, based on the available unused bandwidth (ABR) and other resources, after the QoS of active services from the other two classes have been met first.

The classes require different processing and queueing priorities at the nodes to ensure their delay requirements are met. Since Class 1 applications have low tolerance for delay and delay variation, they cannot be stored and forwarded in a long buffer, as can Class 3 services, nor can they be retransmitted with a feedback mechanism, used when cumulative errors cannot be corrected by codecs. Unlike Class 2 or Class 3 messages, Class 1 streams cannot be demoted in service priority at the nodes without loss of the voice over IP (VoIP) link.

It is essential to the cross-layer approach to relate the service classes to the operation of the MANET protocol layers. At the application layer, it is common to classify the service requirements into a set of QoS priority classes with their corresponding application layer parameters (ALPs). The correspondence between the Classes 1, 2 and 3 and the service requirements of the application-service priority classes is given in the following, in terms of a mapping to the appropriate performance metrics and parameters. Class 1 corresponds to applications that have strict *delay* constraints, such as, voice. This class is mapped to the *delay* metric of the ALPs. Class 2 corresponds to applications requiring high *throughput* such as video and bulk transaction processing tasks. Similarly, this class is mapped to the throughput metric of the ALPs. Class 3 is not defined by specific constraints and is mapped to *best-effort* services of the ALPs. The mapping is depicted in Figure 1.

Network layer functions include connection management, e.g., fast virtual-circuit reservation; packet store and forwarding; congestion control, e.g., low-priority packet dropping; and other aspects of packet routing and flow control. These functions are to be optimized in the cross-layer protocol design. At the network layer, the power state (on/off), queue state, and stability of the nodes characterize the quality of the network. These quantities are called the network layer parameters (NLPs). The power level represents the amount of available battery capacity over time. The queue state indicates the available, unallocated buffer space at each node. Stability refers to the transience in connectivity as measured by the connectivity variance of a node with respect to neighboring nodes over time. To compute the quality of a path, a combination of concave and additive functionals are used to represent the NLPs of a path, given the NLP values for the nodes comprising the path. The NLPs of a node can also indicate whether the node is forced to be *selfish*. In the so-called selfish mode, the node can halt its routing function and act only as a host due to its poor "quality."

The data link/MAC layer functions include error detection and recovery, congestion control on the link level, transmission scheduling and priorities, collision avoidance, and time-division scheduling. At the MAC layer, the effect of these functions on the quality of the network is represented by the signal-to-interference-plus-noise power ratio (SINR) on the links to the node and in the nodes transceiver channels. The SINR parameters are referred to as the MAC layer parameters (MLPs). The SINR determines the communication performance of the link, that is, the data rate and associated bit error probability, packet error rate, or equivalent bit error rate (BER) that can be supported by the link. Links with low SINR usually cannot support the service requirements of some applications and are thus avoided in route discovery. This, in turn, leads to partial connectivity or partitioning among the nodes in the MANET. It is also important to minimize and balance the traffic load being transmitted over the wireless interface because of the scarcity of network resources. This can be achieved through the use of coding schemes. To improve the link quality to meet the QoS requirements of various applications, coding schemes such as forward error correction (FEC) and automatic repeat request (ARQ) are applied. FEC uses a coding scheme for both error detection and correction that imposes a fixed overhead on the data bit streams. FEC is more appropriate for a high-priority class requiring very low BER, e.g., Class 1. ARQ uses only an error detecting code; when an error is detected in a packet, a packet is retransmitted. The ARQ scheme is effective provided the channel BER is not too high and retransmission delay is allowable. ARQ is more suitable for s low-priority service classification, e.g., Class 3. Hybrid ARQ/FEC methods exploit the benefits of the two schemes. If the packet error cannot be corrected by the error-correcting scheme at the receiver, a retransmission will be requested. This technique is suited for the medium priority class, e.g., Class 2. Furthermore, bandwidth conservation is weighed against the processing requirements on the mobile nodes. Consequently, the complexity of coding algorithms and processing delays are also factors in the tradeoff analysis.
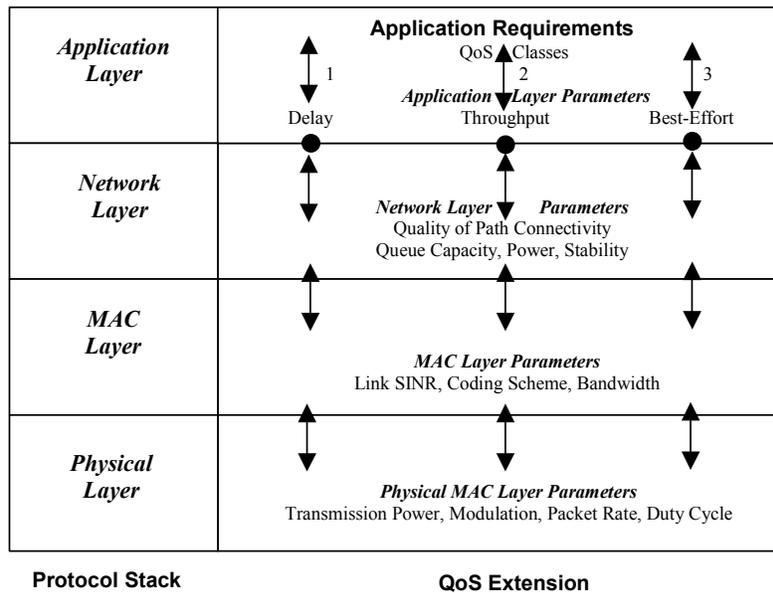
Figure 1. Layer interchange of parameters and metrics in cross-layer QoS model

At the physical (PHY) layer, network quality is represented by the parameters of transmission power , channel code selection, packet rate, modulation, and antenna modes of diversity and beam-forming for the transmissions. Moreover, the on/off characteristics and the duty cycle of transmit, receive (idle), and sleep modes of node operation are ascribed to the PHY layer. These parameters are referred to as the PHY layer parameters (PLPs) and influence the conservation of battery energy and the network lifetime. Transmission power and duty cycle directly affect the remaining (residual) battery capacity at the nodes over time, an NLP; the former influences the level of channel interference, i.e., the SINR, an MLP. Transmission power levels also control the hop count and stability of network connectivity, both NLPs. Channel coding, modulation type and antenna modes also impact the BER or SINR. The choice of bandwidth and packet rate influence throughput, associated with the application layer, while the packet rate is one driver of the queue state, an NLP.

The NLPs, MLPs and PLPs determine the quality of the network and battery utilization, thereby generating energy-efficient paths with adequate QoS support. The primary objectives of the cross-layer protocol optimization are to use these parameters to distribute the traffic and energy consumption in the network and to avoid selection of end-to-end paths with quality below that required by the highest priority application in the packet flow. Then, ALPs select the path, or paths, from the candidates with acceptable quality, that have the greatest likelihood to meet the application requirements with the lowest and best distribution of energy consumption among the nodes. The selection requires a tradeoff among multiple objectives. Thus, it is implied that applications may be required to *adapt to the available quality of the network*.[18] This concept motivates a proposed cross-layer QoS model for energy-constrained networks that responds to both MANET resource limitations and diverse application requirements. The model does not define specific protocols or implementations; however, its implementation can be effected as an embedded modification of a known reactive routing protocol for MANETs, which do not require continuous maintenance of routing tables at the nodes. Examples of these are Dynamic Source Routing (DSR) and Ad Hoc On-Demand Distance Vector (AODV) routing.

Figure 1 shows the QoS classes and energy conservation objectives along with their corresponding ALPs, NLPs, MLPs, and PLPs. Table 1 shows one proposed mapping between the service requirements and energy constraints on battery reserves and on network lifetime and the ALPs, NLPs, MLPs, and PLPs. In this model, Class 1 and Class 2 can be mapped to the queue size and hop count of the NLPs, to the link SINR and bandwidth of the MLPs, and the channel code and packet rate of the PLPs. Class 3 is mapped to the connectivity stability or variance and hop count of the NLPs and to the link SINR of the MLP, and to the modulation, packet rate and antenna mode of the PHY layer. Energy metrics of battery capacities and network lifetime can be directly mapped to transmission power and on/off duty cycle of

the PLPs, while they also indirectly depend on the throughput and delay requirements of the application layer, hop count of the NLPs, and SINR created by the channel access scheme of the MAC layer. In summary, limits on the parameters and metrics associated with the protocol layers form a set of constraints along with the energy-related constraints on the determination of the best available paths between source and destination.

Table 1. Mapping QoS requirements to protocol layer parameters

| QoS Class/Requirement | ALPs | NLPs | MLPs | PLPs |
|---|---|---|---|---|
| Class 1 | Delay | Buffer size, Hop Count, Shaping | SINR, Bandwidth Access Slot | Code, Packet Rate |
| Class 2 | Throughput, Delay Variation | Buffer size, Hop Count, Shaping | SINR, Bandwidth Access Slot | Code, Packet Rate |
| Class 3 | Best-effort | Connection stability, Hop Count | SINR, Access Slot | Antenna Mode, Modulation, Packet Rate |
| Remaining Battery Capacity at Nodes, on Paths | Throughput, Delay | Hop Count, Shaping | SINR | TX Power, Duty Cycle, Antenna Mode |

The parameters and metrics associated with the protocol layers are directly used in the construction of the multivariate point-process stochastic model developed in Section 4.

### 3.3 Example

An example is provided to explain how a service application can adapt to the corresponding ALPs and, hence, to the available quality of the network. Consider a shaping mechanism.[5] Shaping is the process of delaying or dropping packets within a flow to cause them to conform to the QoS state of the selected path.[19] To decide whether to delay or drop packets, a node checks the requirements of the application. If the application is delay-sensitive, that is, of Class 1 then packet dropping may be used. Although packet dropping implies an increase in the loss rate, the probability of path failure is reduced, thus avoiding added congestion and delay. Conversely, if the application requires a low-loss rate, typical of Class 2 applications, then delaying and buffering the stream may be more appropriate when connection stability is high, and the path can thus support the induced queueing delay. At the network layer, the routing must be adaptive to the NLP values of the nodes in the path generated between source and destination. The MAC layer, on the other hand, can adapt the coding technique, with codes received from the PHY layer, to satisfy the service application requirements, given current channel and network conditions.

## 4. ANALYTICAL ELEMENTS OF CROSS-LAYER MODELS

The following introduces the analytical basis for the representation of a wide range of the traffic, link distortions, control as well as the metrics and constraints of cross-layer protocol design in terms of the real-time MVPP models.

The MANET operations that handle the flow of multimedia packets during the input, support, and completion of a call are first discussion. The call's source can originate from any active mobile station (MS) $i$ as one of $M_{max}$ mobile nodes in the network. The uninterrupted arrival stream of packets in the MANET are commonly assumed to occur according to a sequence of random $F_t$-stopping times, $\tau_0^a, \tau_1^a, \ldots, \tau_n^a, \ldots$, such that the corresponding sequence of inter-arrival times, $\tau_1^a - \tau_0^a, \tau_2^a - \tau_1^a, \ldots, \tau_{n+1}^a - \tau_n^a, \ldots$, are independent and identically distributed (i.i.d.) random variables with the common arbitrary right-continuous PDF, $F^a(t), \tau_n^a < t \leq \tau_{n+1}^a$, for every $n$. In other words, the inter-arrival sequence forms a renewal process. In practical operation, the PDFs can change between stopping times due to transient behavior, i.e., $F_n^a(t), \tau_n^a < t \leq \tau_{n+1}^a$, for every $n$, violating the renewal assumption. In synchronous network operation or computer simulations, arrival times as well as service times can be slotted and deterministic, i.e., $\tau_n = T_n = nT_0$, for $T_0$ a

known slot time, frame time, or simulation time increment. For real-time multimedia applications in MANETs, arriving packet streams may be assumed to have a set of multiple service requirements, selected from at most $S$ service types. Each type requires a different QoS, expressed in terms of protocol parameters and metrics. Any of the $S$ integrated service applications may be active, requiring a different set of network resources to maintain their distinct QoS requirements during processing. Therefore, the sequence of service-completion times $\tau_{i,m}^s$, corresponding to service application $s, s = 1, \ldots, S$, at node $i$ and the corresponding sequence of inter-service completion times, $\tau_{i,1}^s - \tau_{i,0}^s, \tau_{i,2}^s - \tau_{i,1}^s, \ldots, \tau_{i,m+1}^s - \tau_{i,m}^s, \ldots$, may not be i.i.d. random variables and may not share a common PDF with the inter-service time sequences associated with other applications. Superposition of the sequences of the inter-service times corresponding to any two or more of the service types will not form a renewal sequence. Without the exponential assumption on inter-service times, superposition of service completions will not, in general, form a renewal process.[20]

The MVPP modeling assumptions must be able to represent the self-similar behavior of Internet traffic.[21] The only assumptions required in the general analytical development are those that support the semi-martingale decomposition of the MVPPs describing integrated multimedia traffic, as the sum of $(F_t)$-predictable, integrated, *non-explosive* rate processes and pure-jump martingales, with respect to the probability space $(\Omega; F_t; \wp)$ or controlled space $(\Omega; F_t; \wp^U)$.

## 4.1 Traffic types and probability distributions

Packet flows from the source node at random time $\tau_n^a$ carry one or more of a maximum $S$ simultaneously active service applications. The service load arriving at time $\tau_n^a$ is modeled as an embedded, discrete vector-valued, discrete-time process $B_n = (b_{1,n}, b_{2,n}, \ldots, b_{S,n})$, where $b_{s,n}$ is the processing load in packets or frames corresponding to service type $s$, $s = 1, 2, \ldots, S$, and can vary from arrival time to arrival time. The condition $b_{s,n} = 0$ indicates service application $s$ is inactive at $\tau_n^a$. Since the load process $(B_n, n \in \mathbf{Z}_+)$ at arrival times $\tau_n^a$ and the counting process on new sessions, denoted $(N_t^A, t \in [0, T])$, have different statistical properties, the combined multimedia service-arrival process may take the form of any number of hybrid MVPPs, based on these individual properties. For example, if $(B_n, n \in \mathbf{Z}_+)$ is a discrete-time, discrete-space Markov process, and $(N_t^A, t \in [0, T])$ a Poisson process with time-varying rate, $\alpha_t, t \in [0, T]$, the combined arrival process is a non-homogeneous, Markov-modulated Poisson process. Through appropriate selection of the statistical properties of service applications and arrival counting processes, many random processes commonly used to model both wired and wireless telecommunication traffic can be constructed.[20]

In general, a multi-server model is appropriate at any node $i$, with the inter-service events of the processors for each service application $s$ obeying a different PDF, $F_{i,s,t}^d, t \in [0, T]$. According to the construction in [22, 23], the corresponding conditional random rate for the class $s$ at node $i$, on the event $\{\tau_{i,n}^s \leq t < \tau_{i,n+1}^s\}$, is $\sigma_{i,s,t \wedge \tau_{i,n+1}^s} = -\dfrac{d\, F_{i,s,t \wedge \tau_{i,n+1}^s}^d / dt}{\left(1 - F_{i,s,t \wedge \tau_{i,n+1}^s}^d\right)}$,

where $\tau_{i,n}^s$ is the $n$'th service completion time at node $i$ and "$\wedge$" denotes the infimum of two stopping times. As the PDF can also change after each time $\tau_{i,n}^s$, the construction allows a marked renewal sequence with a conditional PDF $F_{i,s,n,t}^d$ between the $n$'th and $n+1$'th service completion times. Martingale representation theory for MVPPs, applied to the counting process $\breve{N}_{i,s,t}^D$ of the uninterrupted packet-processing completions for type $s$ to time $t$, leads to the result

$$E\left[\breve{N}_{i,s,t}^D\right] = E\left[\sum_n \int_{\tau_n^s}^{t \wedge t_{n+1}^s} \sigma_{i,s,v} dv\right], \text{ and } \breve{N}_{i,s,t}^D - \sum_n \int_{\tau_n^s}^{t \wedge t_{n+1}^s} \sigma_{i,s,v} dv \qquad (1)$$

is a zero-mean $(F_t, \wp)$-martingale, provided $F_t$ is a $\sigma$-algebra of the network events to time $t$ containing the history $\{\check{N}^D_{i,s,v}, 0 \le v \le t\}$. A representation can also be provided for the counting process $\check{N}^A_{i,s,t}$ on the uninterrupted number of packet arrivals to node $i$ of application $s$ to time $t$ in terms of the conditional arrival rate $\alpha_{i,s,t}$ and the sequence $\left(\tau^a_{s,n}\right)$.

*Class 1 applications*. Voice traffic and interactive video inject CBR traffic into networks. These services cannot function with less bandwidth or bit rate than the minimum application requirements, nor benefit from extra bandwidth. Effective voice over (VoIP) services require application-to-application delays of less than 150 ms and packet loss rates of less than 2%. Voice service can be modeled by a non-homogeneous, Markov-modulated Poisson process (MMPP), with one or two selectable CBRs, $\alpha_1$ and $\alpha_2$, as the Poisson intensities. These rates are modulated by a random "on-off" process, $V_A$, with a mean "on" time equal to the average talkspurt activity cycle. Packetized voice service is generally not buffered, so that if the processors at a node are busy, the connection request for live voice is blocked, routed to a neighboring node, or, as a last resort, dropped.

*Class 2 applications*. Traditional interactive data applications, such as Telnet sessions, and interactive multimedia applications, such as modern codecs and LAN TV, are more "bursty" in nature and fluctuate between low- and high-rate requirements. Researchers have previously used MMPPs to model aggregate voice, video and data VBR traffic.[24] For example, an MPEG-2 video encoder generates a bit stream which is modeled at the video frame level. MPEG-2 frames can be of type intra (*I*), predictive (*P*), or bi-directional (*B*). Here only *I* and *P* frame types are considered, with the *I* frames being generated at scene changes. The *I*-frame bit rates contribute to the largest amplitudes, while the *P* frames transmit the differential information in successive frames and result in a distribution of moderately valued frame rates. The VBR video source model is represented by an *I* state discrete-time Markov chain, with a transition probability matrix $\mathbf{P}_V$ and a rate vector $\mathbf{R}_V = [r_1, r_2, \ldots, r_I]$. The rate $r_i$ represents the number of bits generated per video frame when the process is in state $i$. The last state *I* represents the intra-frame state and the remaining states corresponding to the *P* frames. The diagonal elements of $\mathbf{P}_V$ are dominant except that $p_{II}=0$. The latter transition results from an immediate transition from an *I* frame to a *P* frame. The diagonally dominant structure signifies that VBR video is characterized by strong short-term correlations.

*Class 3 applications*. Data applications, such as file transfers or multimedia mail, function with a wide range of available bandwidth. Packets for such connection-less services can be buffered in queues at the nodes. Packetized data networks adequately support ABR services with best-effort QoS guarantees. Arrival and processing completions for ABR applications are best modeled with PDFs with parameters that can be adjusted to the available RRs of idle signal processors, transport channels, etc., such as those for marked renewal processes. Studies [21, 25] reveal that packet loss and delay is very different in simulations using actual aggregate self-similar traffic data with LRD rather than traditional MMPP models.

To apply the MVPP approach to self-similar traffic, several models have been proposed that capture the LRD behavior, e.g., M/G/∞ model with Pareto service times [26], the superposition of two-state Markov sources [27], the mixture of exponentials to fit the heavy-tail distributions, the superposition of *N* "on-off" processes with sub-exponential "on" periods [28], deterministic chaotic-maps, and self-similar (fractal) point processes.[29] In each model, it is shown that the number of arrivals over any interval (number of busy servers in an M/G/∞ model) exhibit an LRD correlation structure.

A mathematically tractable model has been developed based on a fractal construction of a basic point process (cluster process), where clusters are embedded over an infinite number of time scales.[30] The new model decomposes the self-similar process in a way that is tractable for characterization and control of packet traffic. The point process is constructed recursively as a succession of embedded "on-off" processes that contain *m* time scales. This process may be viewed as the basic process embedded in the "on" state of the *m*'th time scale. The time between visits to the "on" and "off" periods in this *m* time-scale process is exponentially distributed with parameter $\lambda q^m$, where $\lambda$ is the underlying Poisson arrival rate and, after each arrival time, a decision is made with probability $p$ to continue generating arrivals with rate $\lambda$ or with probability $1 - p = q$ to turn off for a period of time. The number of arrivals before entering an "off"

period is geometrically distributed with a mean $q^{-1}$. The probability density function of inter-arrival times for the corresponding point process is

$$
f_{m,\tau}(t) = \begin{cases} \dfrac{2\lambda e^{-\lambda t} + \sum_{i=1}^{m}(2q)^i \lambda q^i e^{-\lambda t}}{2 + \sum_{i=1}^{m}(2q)^i} & \forall t \geq 0 \\ 0 & \text{otherwise.} \end{cases} \tag{2}
$$

Let $f_\tau(t) \equiv \lim_{m \to \infty} f_{m,\tau}(t)$, which uniformly converges $\forall t \geq 0$ to the limit probability density function, when the generalized process consists of an infinite number of time-scale embeddings,

$$
f_\tau(t) = \begin{cases} \dfrac{2\lambda e^{-\lambda t} + \sum_{i=1}^{\infty}(2q)^i \lambda q^i e^{-\lambda t}}{2 + \sum_{i=1}^{\infty}(2q)^i} & \forall t \geq 0 \\ 0 & \text{otherwise.} \end{cases} \tag{3}
$$

Potential modeling problems with LRD can be avoided by observing MANET operation over a finite horizon and considering practical queues of finite size. Grossglauser and Bolot report that the amount of correlation that needs to be considered for performance evaluation depends not only on the correlation properties of the source traffic, but also on the time scales specific to the system under study.[31] For example, the time scale associated with a queueing system is a function of the maximum buffer size. Thus, if finite-size queues are available at the nodes, the impact on the loss of the correlation in the arrival process becomes negligible beyond a time scale referred to as the *correlation horizon* (CH).[31] Consequently, *any* model, among a host of available models including Markov and self-similar processes, can be chosen as long as the selected model captures the correlation structure of the source traffic *up to* the CH. *Truncated* forms of standard heavy-tailed CDFs are thus sufficient for the development of accurate models. One such example for the inter-arrival time distribution is a truncated Pareto $F_{T_C}(t)$ is a truncated Pareto PDF defined by

$$
F_{T_C}(t) = \begin{cases} \left(\dfrac{t+k}{k}\right)^{-\theta}, & \text{if } 0 < t < T_C < \infty \\ 0, & \text{otherwise,} \end{cases} \tag{4}
$$

where $1 < \theta < 2$ and the parameter $T_C$ is referred to as the *cutoff lag*. Truncated versions of (3) can also be defined. The cutoff lag eliminates correlation in the input process beyond a lag equal to $T_C$.

## 4.2 Energy metrics and constraints

Although finding routes that meet QoS requirements is an important design issue for MANETs bearing real-time multimedia services, a more universal objective is to provide energy-efficient routes, since the nodes' operational time is the most critical limiting factor. Proposed energy-efficient cross-layer protocols should minimize either the *active communication energy* required to transmit or receive packets or the *inactive energy* consumed when a node remains idle but listens to the wireless medium for any possible communication requests from other nodes. *Transmission power control* approach and *load distribution* approach belong to the former category, while *sleep/power-down mode* approach belongs to the latter category. Metrics used to determine energy-efficient routing are energy consumed/packet, network lifetime or time to network partition, variance in node power levels, energy/packet, and maximum node cost.

The *load distribution* approach is used in this paper. The goal of this approach is to balance energy usage among all mobile nodes by selecting a path with underutilized nodes rather than the shortest route. This may result in longer routes, but packets are routed only through energy-rich intermediate nodes. Protocols based on this approach do not necessarily provide the lowest energy route, but prevent certain nodes from being overloaded, thus ensuring longer network lifetime. One such protocol is the Conditional Max-Min Battery Capacity Routing (CMMBCR) protocol.[32] In this paper, modifications to CMMBCR produce a new energy-aware route maintenance protocol that can be applied to reactive routing protocols without the need for periodic route recovery. Whenever all source nodes determine their routes, the CMMBCR protocol is used to both minimize the total transmission power and prolong the life of each node by taking into account the remaining battery power. Furthermore, when the remaining battery power of an intermediate node, actively forwarding packets, reaches a critically low level prior to node failure, the protocol allows selected end-to-end connections through this node to find alternate paths to avoid overuse of the available power at the node.

In CMMBCR, when all nodes on a set of paths have sufficient remaining battery capacity, above a threshold $\theta$, a path with minimum total transmission power among these paths is selected. Since less total power is required to forward packets for each connection, the relaying load for most nodes can be reduced, and their lifetime will be extended. However, if all paths have nodes with low battery capacity, i.e., below the threshold, a path including nodes with the lowest battery reserves must be avoided to extend the lifetime of these nodes. The battery capacity for source-destination path $p_j$ at time t is $R_j(t) = \min_{\text{node } i \, \in \, p_j} C_i(t)$, where $C_i(t)$ is the remaining battery capacity of node $i$ at time $t$.

Given node $a$ and node $b$, this protocol mechanism considers two sets $A$ and $B$, where $A$ is the set of all possible paths between nodes $a$ and $b$ at time $t$, and $B$ is the set of all possible paths between any two nodes at time $t$ for which the condition $R_j(t) \geq \theta$ holds. Paths are selected as follows: if all nodes in a given path have remaining battery capacity greater than $\theta$, choose a path in $A \cap B \neq \phi$ by applying the Minimum Total Transmission Power Routing (MTTPR) protocol; otherwise, select a path $p_j$ with the maximum battery capacity, i.e., apply the Min-Max Battery Cost Routing (MMBCR) protocol.

CMMBCR considers both minimum transmission power and remaining battery capacity. However, when applied to ad-hoc on-demand routing protocols that do not perform periodic route recovery when there is no path break, CMMBCR alone cannot evenly distribute power consumption to each node. CMMBCR is applied only when setting up a new path, so it does not reflect power consumption rate of each node during active communications. For example, a node could be selected as an intermediate node for many source-destination connections at the same time. As a result, the node's power drain rate is more severe than others on the paths and it soon fails. To avoid this problem, the route maintenance procedures are revised so that the battery consumption rate of nodes can be evenly distributed.

In CMMBCR, the threshold $\theta$ must be not be fixed. If one attempts to determine $\theta$ as an absolute value, CMMBCR gets the nodes with residual battery capacities less than the threshold to participate in the MMBCR, which can cause more nodes to expend energy in a longer route. However, if nodes with residual battery capacity less than $\theta$ are allowed to participate in the MTTPR procedure when the traffic is light, more energy can be saved in the overall network than with CMMBCR. Otherwise, if $\theta$ is defined as the relative percentage of the remaining battery capacity of each node, then there is no way to efficiently find $\theta$, since no centralized controller exists that knows the energy status of all nodes.

This motivates the revisions made to CMMBCR. In the revised CMMBCR, when acquiring an initial route, the source is allowed to specify the minimum remaining battery capacity required of nodes comprising the path according to the application type of its route request, i.e., $\theta = \theta(s), s = 1,\ldots,S$. This modification directly relates energy conservation at the nodes and on the paths to the application layer. Thus, among routes for which the constituent nodes have residual battery capacities above $\theta(s)$, the route with the minimum total transmission power is selected. It is assumed that nodes have the ability to monitor their remaining battery capacity.

To perform route maintenance based on the remaining power of nodes during a data session, two battery thresholds are defined: (1) selective-victim-search zone (SVSZ), $\theta_{SVSZ}$, and (2) forced-victim-search-zone (FVSZ), $\theta_{FVSZ}$.[32] The SVSZ is used to signify that the battery of the node is running low, but remains adequate to keep the node operational. However, without promptly relieving the node of its routing tasks, the node will soon deplete its battery, because it may be forwarding packets from so many connections causing its battery capacity to fall rapidly to $\theta_{FVSZ}$. In this case, the node should attempt to select a connection for which it continues to forward packets, while it tries to find another route for the connection, if one exists. If alternate paths do not exist, the SVSZ node continues to forward packets as usual, until the remaining capacity reaches the lower threshold, $\theta_{FVSZ}$. Below $\theta_{FVSZ}$, the battery level is low enough to force the node to decline forwarding packets for other nodes. All connections passing through this node should find alternate routes around it. The remaining power of the node below $\theta_{FVSZ}$ is reserved for packets that will be generated only when the node acts as a source. In this context, if the remaining battery capacity of a node is above $\theta_{SVSZ}$, and above the battery threshold $\theta(s)$ requested by the source of a connection, it continues to relay packets. Meanwhile, if the current battery capacity of a node falls below $\theta(s)$, the source and destination should be notified to re-

route based on $\theta(s)$, because the path can no longer satisfy the requested battery capacity required by the application type.

When the battery capacity $C(t)$ of a node is such that $\theta_{FVSZ} < C(t) \leq \theta_{SVSZ}$, a connection passing through the node is selectively chosen for which the re-route is required to evenly distribute the overall power consumption. Connections that have expended the intermediate node's battery power above the average battery power consumed by a connection through the node become candidates for the victim. The rationale for this selection rule is that it allows equitable use of the resources, including battery capacity, at this node among all paths that use the node. However, the SVSZ node will continue to forward packets for the sources of the connections, while the victims try to find alternate routes.

Upon entering the SVSZ state, a node must select a connection between source and destination and notify the source to find another path to evenly distribute its power requirements to other nodes. Selection is based on two specific conditions. In "definite selection," if the current battery capacity of a node falls below the value $\theta(s)$ that a source requires when seeking a new path, the source are notified or "selected" to re-route. Definite selection also occurs when all the connections through a node with current battery capacity at or below $\theta_{FVSZ}$ are notified to re-route around the node. In "possible selection," if the current battery capacity of a node is such that $\theta_{FVSZ} < C(t) \leq \theta_{SVSZ}$, a connection through the node is selected as the "victim" to re-route based on one of the following conditions:

(1) The routing protocol monitors the source-destination pairs with some discrete port numbers in the packet to distinguish flows and to maintain a count of the packets that are serviced for each flow through a node. Connections, on which flows have been serviced above the average service rate for all flows, are possible candidates for the victim.

(2) The routing protocol monitors and records the amount of energy expended on each packet, resulting from the transmission and reception of each flow through the node with the aid of the network interface card. Thus, connections, on which flows have expended more than the average energy consumption for all flows through the node, are possible victims.

Traffic load, expended energy per packet, or a combination of these criteria can be used in the victim selection rule. Conditions based on other metrics and parameters of the protocol layers can also be used in the selection rule.

Once the victim connection is selected, the intermediate node sends a selective victim message packet to the source of the victim to alert this source that it should seek another route that does not include a node in the SVSZ state. The SVSZ node continues to relay packets for this source until the source finds another route. If an alternate path does not exist, the residual battery energy of the node will drain until it enters the FVSZ state, whereupon the node stops relaying packets for connections through it and sends a forced victim message packet to the source of each such connection.

Let $C(t)$ be the remaining battery capacity of the intermediate node at time t, and $\theta_i(s)$ be the minimum requirement for the remaining power of nodes requested by the source of path $p_i$ to support packets of application type $s$. The following pseudo-code describes the energy-aware route maintenance algorithm.[32]

## 4.3 Class of admissible routing policies

The most significant control function of the MANET protocol is routing the multimedia packet flows among nodes on source-to-destination paths in the autonomous network. Routing packets from node $i$ to node $j$ is modeled by a random vector $\boldsymbol{u}_{ij,t} = \left( u_{ij,1,t}, u_{ij,2,t}, \ldots, u_{ij,S,t} \right)$. The $s$'th component of $\boldsymbol{u}_{ij,t}$ is the indicator of the effect of a network event on application class $s$, e.g., a change in processing or bit rate, service cessation, or service interruption due to a lack of required resources at the nodes of the path. The components of $\boldsymbol{u}_{ij,t}$ are $\left( F_t \right)$-predictable (more strongly, left-continuous) indicator functions of random events and Borel-measurable. The expected values of the $\boldsymbol{u}_{ij,t}$ with respect to $\wp$ and $\Omega$, are the probabilities of packet forwarding and dropping that shape the flows from source to destination.

```
if (C(t) < θ_i(s))
        Connection on path p_i should be notified to re-route.
else if (C(t) > θ_SVSZ)
            Relay packet(s).
else if (θ_FVSZ < C(t) ≤ θ_SVSZ)
        {
        Notify sources of connections using the node to selectively re-route
                according to the Victim Selection Rule;
        Continue relaying packets.
        }
else if (θ_FVSZ ≥ C(t))
        {
         Stop relaying packets on all paths through the node;
         Notify all paths through the node to re-route;
         Node only used to transmit or receive packets within one hop.
        }
}
```

Figure 2. Energy-aware route maintenance algorithm

The $M_{max} \times M_{max} \times S$ time-varying control array,

$$U_t = (u_{ij,s,t}; i,j = 1, \ldots, M_{max}; s = 1, \ldots, S), \ t \in [0,T]$$

describes the random connectivity granted to the QoS classes, due to routing, resource reservation, and radio path distortions, over the period of observed network operation. Note that the sum of entries of $\boldsymbol{u}_{ij,t}$ over $j$ may be greater than 1 to model point-to-point, point-to-multipoint, and broadcast transmissions among the nodes. Controls are assumed to be closed under concatenation in time, that is, $U=[U_1,U_2,t]$, $t \in [0,T]$, is also a control if $U_1$ and $U_2$ are controls. Controls are also assumed to satisfy a stochastic causality property. The set of control arrays $U_t$, $t \in [0,T]$, that satisfy the above conditions, are referred to the *class of admissible routing policies, $\mathcal{U}$.*

As $t$ varies, the sample paths of the entries $u_{ij,s,t}$ form the temporal evolution of source-to-destination connectivity enabled by network conditions for both connection-less and connection-oriented services. The exponential random rates $\lambda_{i,t}$ of session duration are assumed to be much lower than the conditional rates for packet arrivals and service completions for all application types. Then, for a service request from source node $j_1$ moving to node $k$, the expectations of the entries in an indicator-function sequence, $(u_{j_1 j_2,s,\tau_1}(\omega), u_{j_2 j_3,s,\tau_2}(\omega), \ldots, u_{j_n k,s,\tau_n}(\omega))$, estimate a history of the path traversed by the packet flow bearing the service, given that $\tau_1(\omega) < \tau_2(\omega) < \cdots < \tau_n(\omega)$ are the $n$ packet arrival or completion times that occur for service application $s$ at nodes $j_1, j_2, \ldots, j_n$, respectively.

Mobility, determined by the location, speed and direction of node $i$ at time $t$, and the resource requirements of the applications active in the service request, determine which nodes can be accessed to construct the route.

## 4.4 Network state

In order to represent the state of the MANET, the candidate process must have sufficient dimensions to distinguish the services, nodes, sources and destinations of messages. Packets from Class 2 and Class 3 applications can be queued during periods of deep fading, high channel interference, blocking, connection instability, low battery capacity, or other link disturbances and to allow preemption by more delay-sensitive services. The queue of service-connected packets, both in service or awaiting service, in a buffer at node $i$ at time $t>0$, is represented as the discrete-valued vector of parallel queues, $\boldsymbol{Q}_{i,t} = (Q_{i,1,t}, Q_{i,2,t}, \ldots, Q_{i,S,t})$. Each component of $\boldsymbol{Q}_{i,t}$ has a corresponding birth-and-death equation and a semi-martingale representation in terms of a discrete-valued jump, right-continuous, zero-mean, $(F_t, \wp)$-local $((F_t, \wp^U)$-local) martingale and an integrated conditional random rate process with respect to the family of $\sigma$-algebras $(F_t, t \in [0,T]; F_t \subseteq F)^{22}$, generated by the evolution of observed network events to time $t$, i.e.,

$$Q_{i,s,t} = Q_{i,s,0} + \left( Q_{i,s,t} - \int_0^t (\alpha_{i,s,v} - 1(Q_{i,s,v-} > 0)\sigma_{i,s,v})dv \right) + \int_0^t (\alpha_{i,s,v} - 1(Q_{i,s,v-} > 0)\sigma_{i,s,v})dv \qquad (5)$$

In the single-server case, the total number of packets at node $i$ at time $t$ is the sum over the number of service application types of the components (5) of $Q_{i,t}$. The queuing array $\boldsymbol{Q}_t = (\boldsymbol{Q}_{1,t}, \boldsymbol{Q}_{2,t}, \ldots, \boldsymbol{Q}_{M_{max},t})$ is augmented with the vector $\boldsymbol{C}_t = (C_{1,t}, C_{2,t}, \ldots, C_{M_{max},t})$ of remaining battery capacities of the nodes at time $t$ to form the energy-constrained state of the MANET.

If the processing capabilities of the mobile nodes can handle only one active call at a time, the instantaneous $(F_t)$–progressively measurable conditional rates for the arrivals and departures of service application type $s$ at node $I$, where node 0 is the source node, are

$$\alpha_{i,s,t} = u_{0i,s,t-}\alpha_{s,t-} + \sum_{j \in \{\text{hearing radius of node } i\}} u_{ji.s.t-}1\left(Q_{j,s,t-} > 0\right)\sigma_{j,s,t-} \qquad (6)$$

and $\sigma_{i,s,t}$, respectively, where the exact formulae for these conditional rates depend on the PDFs in (1) for the underlying random events of packet arrival and service completion for each application type $s$, the observed network history to time $t$ on which the conditional rates are estimated, and the admissible control policy $U \in \mathcal{U}$. Indicator functions in (5) and (6) show that uninterrupted packet processing cannot occur when no packets of type $s$ are in the node. In the special case of non-homogeneous Poisson arrivals with deterministic, time-varying intensities $\alpha_{i,s,t}$, and single-stage exponential processors with deterministic transient rates $\sigma_{i,s,t}$ at the nodes, the *form* of the rates coincide with equation (6). In general, the structure of the rates in (6) is more complex and is conditioned on the time $t$ between stopping times between $\tau_{i,n}^s$ and $\tau_{i,n+1}^s$, i.e., the event $\{\tau_{i,n}^s \leq t \leq \tau_{i,n+1}^s\}$, for each service application $s$.

At nodes with extended capabilities, multiple parallel processors may handle at most $L$ simultaneous service requests, each of which may carry up to $S$ active applications. Therefore, at most $L \times S$ packet processors or processing modes can be assumed available at these nodes. Packet flows arrive at node $j$ from nodes within its listening distance, originating from source nodes. At node $j$ the random instantaneous rates of the arrivals and departures of service application $s$ are given by

$$\alpha_{j,s,t} = \sum_{k \in \{\text{source nodes in coverage area of node } j\}} u_{kj,s,t-}\alpha_{k,s,t-} + \sum_{i \in \{\text{intermediate nodes in coverage area of node } j\}} u_{ij,s,t-}1\left(Q_{j,s,t-} > 0\right)\sigma_{i,s,t-} \qquad (7)$$

and

$$\sigma_{j,s,t-} = \sum_{l=1}^{L} 1\left(Q_{j,s,t-}^l > 0\right)\sigma_{j,s,t-}^l, \qquad (8)$$

respectively, where $\sigma_{j,s,t}^l$ is the random packet completion rate of the $l$'th processor for application type $s$ at node $j$. The routing terms $u_{ij,s,t-}$ in (7) and (8) determine the flow of packets over links of the end-to-end path. Connectivity between nodes is also determined by less controllable, often unobservable events, created by resource constraints and link disturbances. The constraints are placed on buffer space, processor limitations and queueing delays, battery capacities, as well as radio path distortions due to multipath reflections, shadowing, and residual channel interference among users. Thus, the general structure of $u_{ij,s,.}$ is a product of the indicators of both controlled, observed events and constrained, indirectly observable events. As indicators of random events, the expected values of the $u_{ij,s,.}$ with respect to $\wp$ and $(F_t, t \in [0,T]; F_t \subseteq F)$ are the probabilities of the events. For example, for routing at the nodes of a stationary network, with $u_{ij,s,t} = 1$ (packets of application $s$ from node $i$ are routed to node $j$ at time $t$), $E[u_{ij,s,t}] = p_{ij,s}$, a fixed probability for any time $t>0$. Depending on the statistical characteristics of network events, i.e., ergodic, stationary, renewal, Markovian, etc., the traffic streams in the MVPP model become randomly modulated point processes of the corresponding stochastic type through the terms $u_{ij,s,.}$ that indicate occurrence of the events. Thus, it is natural to

partition each admissible routing array, $U=[U_C, U_{NC}]$, into two components, corresponding to the controlled and uncontrollable events that influence network quality and connectivity.

## 4.5 QoS and resource constraints

The QoS requirements of each service application are represented by the metrics and parameters of the protocol layers, particularly, the application layer, that shape network events and their corresponding indicator functions $(u_{ij,s,t})$.

Parameters, such as, throughput, delay, buffer size, packet loss rate, SINR, battery capacity, bandwidth, and connectivity stability are variables that explicitly determine the conditional rates of packet arrivals, processing, and routing for each service application at its required QoS. The MVPP models thus encompass *QoS-based* and *energy-constrained routing*.

For example, PHY-layer adaptive beamforming within a node's coverage area can simultaneously increase the gain factors of BER and reduce co-channel interference. If the spatial distribution of nodes is uniform, coverage-area sectorization reduces interference and increases capacity by the antenna gain factor, $G_A$. Voice activity monitoring (VAM) can be modeled either by the indicator of a gain condition, $1(G_{\pi,ij,t} \geq 1/act\%)$, where $G_v$ is the voice activity gain, or by the indicator of the random event of on-off voice activity, $1(\kappa_{ij,1,t} > 0)$, where, by convention, application type 1 denotes live voice and $\kappa_{ij,1,t}$ the voice activity on the link between node $i$ and node $j$ at time $t$. In general, the entry $u_{ij,s,t}$ is the product of factors that include $1(BER_s \leq 10^{-n})$ for service application $s$, which can, in turn, be factored into a product of indicators of protocol parameter values that together comprise the QoS requirements for type $s$. For example, $1(BER_t \leq 10^{-n}) = 1(P_{TX,i,t} \geq P_s) \cdot 1(G_{\text{coding},ij,t} \geq 10^{y/10}) \cdot 1(G_{A,j,t} \geq g) \cdot 1(P_{\text{avg. path loss},ij,t} \leq \pi) \cdot 1(C_i(t) \geq \theta(s)) \cdot 1(C_j(t) \geq \theta(s)) \cdot 1(I_{\text{co-ch .interf.},ij,t} \leq \varsigma_s) \cdot 1(BW \geq B_s)$. For fading links, the condition $P_{\text{avg. path loss},ij,t} \leq \pi$ can be replaced with the random event of the number of replicas received at node $j$ of the transmission from node $i$, according to a discrete-event distribution, with the relative path loss on each replica obeying a Rayleigh or Rician distribution. The joint distribution of the number of reflected paths and the amplitude of the reflections is used to determine the expected value of the indicator of the multipath fading event. Traffic shaping mechanisms in response to resource constraints can also be represented in the model by the $u_{ij,s,.}$, with factors of the general form $1(\text{Available } NR_{j,t} \geq \rho_s)$, where $NR_{j,t}$ is the network resource at node $j$ at time $t$ required at level $\rho_s$ to guarantee the QoS of application type $s$.

## 4.6 Controlled probability measures

A family of probability measures $\wp^U$ on the network events $\Omega$ is constructed from a reference measure and the class of admissible routing policies $\mathcal{U}$. The role of the Radon-Nikodym derivatives or likelihood ratios is fundamental to the construction. An absolutely continuous change of measure is given in terms of the local description of the MVPPs that define the packet-flow behavior, that is, $dN_{ijs,t}d\wp \rightarrow dN_{ijs,t}d\wp^U$, the change of conditional rates $p_{ojs}\alpha_t \rightarrow u_{0js,t}\alpha_{s,t}$ and $p_{ijs}\sigma_{is,t}1(Q_{is,t-} > 0) \rightarrow u_{ijs,t}\sigma_{is,t}1(Q_{is,t-} > 0)$ for $U \in \mathcal{U}$. Details of this construction are omitted, but have been presented as a variation of the result by Doleans-Dade applied to MVPPs created for the MANET. [23, 33]

## 5. CROSS-LAYER OPTIMIZATION OF THE ENERGY-CONSTRAINED MANET

The cross-layer optimization problem is formulated in terms of protocol metrics on MANET performance for admissible $U \in \mathcal{U}$, and expressed as the constrained optimization of cost functionals constructed from the MVPPs. Optimality conditions for the routing policies $U \in \mathcal{U}$ that yield the best network "quality," given the observation $\sigma$-algebras on network events, take the form of generalized backward recursive relations.

## 5.1 Performance metrics

*Throughput and queue capacity.* An application-layer performance metric is throughput, the time average of the number of packets of each service application delivered from source to destination per unit of time. Related metrics of interest include the distribution of the number of packet transmissions from each node, a PHY layer metric; the time

average of packet delay, an application layer metric; the fraction of channel capacity used for successful transmission, a network layer metric; and the probability of successful packet transmission, or goodput, another application layer measure.

In terms of the MVPP processes, the number of packets of service application $s$, $s = 1, \ldots, S$, in the network is the sum of the individual components in the buffer or queue state $\boldsymbol{Q}_{,t}$, that is, $Q_{\text{total},s,t} = \sum_{i=1}^{M_{\max}} Q_{i,s,t}$. Under the assumptions stated in Section 4, $Q_{\text{total}, s}$ is a right-continuous, $(\wp, F_t)$-supermartingale and a non-explosive MVPP with a conditional random rate that is the sum over $i$ of the rates of the $Q_{i,s,t}$ given in (5).

The *link throughput of service type s* over the subinterval $(v, t] \subset [0, T]$ for any $v < t$ from node $i$ to node $j$ is given by $\widetilde{N}_{ijs,w} = N_{ijs,w}^C - N_{ijs,w}^{NC}$, the difference between controlled and uncontrolled, or lost, packets of type $s$ transported between the nodes over the interval. Similarly, the *node i throughput of type s* over $(v, t] \subset [0, T]$ is denoted as $L_{is,t} - L_{is,v}$, where $L_{is}$ is defined for $i = 1, \ldots, M_{\max}; s = 1, \ldots, S$ and $v \in [0, T]$ by $L_{is,v} = \sum_{j=1}^{M_{\max}} \widetilde{N}_{ijs,v}$. Lastly, the *system throughput of type s* over $(v, t]$ is given by $L_{s,t} - L_{s,v}$, where $L_s$ is defined for all $v \in [0, T]$ by $L_{s,v} = \sum_{i=1}^{M_{\max}} \widetilde{N}_{id,v}$, and, once a packet of type $s$ is successfully reaches destination node $d$, it assumed that the packet cannot be returned.

The *average throughput rates*, corresponding to the stochastic throughput processes, are formed by taking the expectation with respect to the probability measure $\wp$, or the controlled probability measure $\wp^U$, divided by the subinterval length $t - v$. Observe that the expressions for the throughput metrics are merely linear combinations of the MVPPs $\left( \breve{N}_{0js,t}^A, \breve{N}_{ijs,t}^D \right)$ and so share with them the same $(\wp, F_t)$-semi-martingale structure described in Section 4.

Define $U_C([0, t])$ as the set of values of all admissible control arrays $U_C$ with entries that satisfy the conditions of admissibility over $[0, t]$. The *mean link capacity from node i to node j* at time $t \in [0, T]$ is defined as $K_{ij,t} = \max_{U_C([0,t])} E^U \left[ \widetilde{N}_{ijs,t} \right]$ while the *mean capacity at node i* at time t is defined as $K_{i,t} = \max_{U_C([0,t])} E^U \left[ L_{is,t} \right]$, where $E^U [\cdot]$ represents integration with respect to $\wp^U$ so the $N_{ijs,\cdot}$ admit $(\wp^U, F_t)$-rates that depend on the values $u_{ijs,t}(\omega), \omega \in \Omega$. in $U_C([0, t])$.

*Packet dropping and blocking.* The expected number of blocked or dropped packets, used in traffic shaping, for a given application or all application types are the expectations of the counting processes for the corresponding network events over the observation period, summed over the indices of interest. Based on the semi-martingale representations of the underlying MVPPs and an assumption that arrival and processing rates are non-explosive, the Fubini Theorem is applied to represent the expectations, with respect to the $\sigma$-finite measure $\wp$, as sums, over the application types and nodes of interest, of the integrals of the expected $(F_t)$-predictable rates of the corresponding packet flows over $[0, T]$. Thus, the number of blocked packets of type $s$ to node $j$ over $[0, T]$ is represented as the $(F_t, \wp)$-semi-martingale

$$N_{\text{blocked},j,s,[0,T]} = \left[ N_{\text{blocked}, j,s,[0,T]} - \int_0^T (1 - u_{0j,s,v}) \mathbf{1}(Q_{j,s,v-} \geq q_{j,s}) \alpha_{s,v} dv \right] + \int_0^T (1 - u_{0j,s,v}) \mathbf{1}(Q_{j,s,v-} \geq q_{j,s}) \alpha_{s,v} dv; \quad (9)$$

the expected value is

$$E[N_{\text{blocked},j,s,[0,T]}] = E\left[ \int_0^T (1 - u_{0j,s,v}) \mathbf{1}(Q_{j,s,v-} \geq q_{j,s}) \alpha_{s,v} dv \right] = \int_0^T (P(Q_{j,s,v-} \geq q_{j,s}) - P(u_{0j,s,v}(Q_{j,s,v-} \geq q_{j,s}))) \overline{\alpha}_{s,v} dv \quad (10)$$

where the second term in the integrand on the right-hand side of (10) is a joint probability, $\overline{\alpha}_{s,t}$ is the $\wp$-mean of the arrival rate of type $s$ at time $t$, and $q_{j,s}$ is buffer size at node $j$ for application type $s$. Equation (10) can be used to

represent the average number of blocked or dropped packets for any or all types, as expectations of the indicators for the events, summed over the indices of interest, then integrated over interval $[0, T]$. Instantaneous average rates of dropped packets are the integrands of the average number of the corresponding network events, at some time $t \in [0, T]$.

*Processing priorities.* Prioritization at a single-processor node can be modeled by re-ordering the processing times $\left(\tau_{i,n}^{s}, s = 1, \ldots, S\right)$ at node $i$, conditioned on the number of packets, $Q_{i,s,t}$, of each type $s$ currently at the node at time $t$ and the last processing completion time $\tau_{i,n*}^{s*}$ before $t$. The custom queueing feature of some network equipment allows reserving specific amounts of bandwidth for each type $s$ to ensure the corresponding stream a minimum allocation of bandwidth. The MVPP model allows the rates $\left(\sigma_{i,s,t}, s = 1, \ldots, S\right)$ to adapt, to a maximum total rate, i.e., $\sum_{s=1}^{S} \sigma_{i,s,t} \leq \Lambda_{i,t}$, given the number of packets of each type and the last processing completion time, to reduce backlogs at the node.

*Delay.* Sources of delay and delay variation over the end-to-end path are user mobility, blocking, channel contention due to path congestion and channel interference, packet retransmission in response to unrecoverable block errors, queueing latency, and the power failure of nodes. The queueing delay and its variation at a node can be bounded to not exceed maximum delay and variance targets, $\Delta_s$ and $\sigma_{\Delta,s}^2$, respectively, assigned at the application layer as the QoS requirements for service application $s$. The average allowable delay for application s is denoted $\overline{\Delta}_s$. Other sources of delay are assumed relatively small and due to the effects of MAC layer design and battery-based energy management.

Little's formula states that the average number of customers in a queueing system in *steady-state* is equal to the arrival rate of customers to the system, times the average time spent in the system. The result makes no specific assumptions regarding the arrival distribution or service processing distribution; nor does it depend upon the number of servers in the system or upon the queueing discipline within the system. In terms of the MVPP model of the queue of packets for application type *s*, the random arrival rate of packets of type *s* and the delay limits, the delay condition can be approximated instantaneously at time *t* or by a time average over an observation interval $[0, T)$,

$$1\left(Q_{i,s,t} < \Delta_s\left[u_{0i,s,t}\alpha_{s,t} + \sum_{j \in \{\text{neighborhood of node } i\}} u_{ji,s,t}\sigma_{j,s,t}\right]\right), \tag{11}$$

$$1\left(\int_0^T Q_{i,s,\tau}d\tau < \overline{\Delta}_s\left[\int_0^T u_{0i,s,v}\alpha_{s,v}dv + \sum_{j \in \{\text{neighborhood of node } i\}} \int_0^T u_{ji,s,w}\sigma_{j,s,w}dw\right]\right). \tag{12}$$

The delay variation condition can be expressed in terms of instantaneous variance of queueing delay at node *i* at time *t*

$$1\left(\left(Q_{i,s,t} - E[Q_{i,s,t}]\right)^2 < \sigma_{\Delta,s}^2\left[\left(u_{0i,s,t}\alpha_{s,t} + \sum_{j \in \{\text{ngbrhd of node } i\}} u_{ji,s,t}\sigma_{j,s,t} - E\left[u_{0i,s,t}\alpha_{s,t} + \sum_{j \in \{\text{ngbrhd of node } i\}} u_{ji,s,t}\sigma_{j,s,t}\right]\right)^2\right]\right)$$
$$\tag{13}$$

where expectation can be taken with respect to the probability measure $\wp$ or the controlled probability measure $\wp^U$.

## 5.2 Optimization problem for the cross-layer design

Each performance metric described in the foregoing can be the expressed as a general cost functional to be optimized over admissible routing policies $U = [U_C, U_{NC}] \in \mathcal{U}$, subject to energy constraints at the nodes and along the paths. To each admissible routing array $U$, there is a unique cost functional of the form:

$$C(U) = E^U\left[\int_0^T c(t, U_t)dt + f_T\right] \tag{14}$$

The terms in (14) are assumed to obey the following: For each $U \in \mathcal{U}$, the instantaneous cost $\left(c(t, U_t), t \in [0, T]\right)$ is a composite process $c(t, U_t(\omega), \omega) = (c \cdot U)(t, \omega)$, where $c$ is $(F_t)$-adapted for each value of $U_t(\omega) \in \mathbf{P}$, the space of values

of the routing arrays. The function $c$ is Lebesgue-measurable with respect to $t$ and continuous in the sample path values $U_t(\omega)$ for all $\omega \in \Omega$; while $c$ itself has left-continuous sample paths with finite right-hand limits at each discontinuity for each $U_t(\omega)$ for all $\omega \in \Omega$. Function $c$ is thus progressively measurable with respect to the family of $\sigma$-algebras, $(F_t, t \in [0,T]; F_t \subseteq F)$, with left-continuous sample paths for each $U \in \mathcal{U}$. The terminal cost $f_T$, a nonnegative, $F$-measurable and $\wp^U$-integrable function for $U \in \mathcal{U}$, represents the cost incurred at the end of MANET operation at $t=T$.

The cross-layer optimization problem is to determine routing policies $U^* \in \mathcal{U}$ over the interval $[0,T]$ that satisfy $C(U^*) = \inf_{U \in \mathcal{U}} C(U)$, subject to energy constraints. The minimum instead of the infimum can be taken, assuming each $U \in \mathcal{U}$ takes values in a compact set and the instantaneous and terminal costs are almost surely $\wp^U$-bounded for $U \in \mathcal{U}$. A policy $U^*$, if it exists, that satisfies the criterion is called an *optimal routing policy* of the cross-layer design.

### 5.3 Recursive optimality conditions: complete network observations

Recursive optimality conditions that characterize optimal policies for the controlled real-time MVPP models of the MANET have been developed.[23] Due to the limits of this exposition, results are presented without proof for the case of complete observations of MANET behavior. The local structure of the optimality conditions resemble the Hamilton-Jacobi-Bellman dynamic programming conditions. With complete observations, the *conditional cost function* $\phi(U_1, U_2, t)$ for the admissible control policy concatenated at time $t$ from $U_1$ and $U_2$ in $\mathcal{U}$, obeys

$$\phi(U_1, U_2, t) = E^{[U_1, U_2, t]}\left[\int_t^T c(v, U_{2,v})dv + f_T \Big| F_t\right] = E^{U_2}\left[\int_t^T c(v, U_{2,v})dv + f_T \Big| F_t\right] = \phi(U_2, U_2, t),$$ based on the causality

conditions imposed on admissible routing policies $U \in \mathcal{U}$. Hence, the *optimal cost-to-go function* $W_t^U = \inf_{U \in \mathcal{U}} \phi(U, \hat{U}, t) = \inf_{\hat{U} \in \mathcal{U}} \phi(\hat{U}, \hat{U}, t)$, that is, $W_t^U = W_t$ does not depend on the policy $U$.

Let the network state be an array-valued process, $(X_t, t \in [0,T])$, formed from the network MVPPs. In an energy-constrained MANET, the network state includes the vector of residual battery capacities at the nodes, $(C_t, t \in [0,T])$, where $C_t = (C_{i,t}, i = 1, \ldots, M_{\max})$ with $C_{i,t}$ is the remaining battery capacity at node $i$ at time $t$. The following theorem introduces functions $w_n$, generalizing the optimal cost-to-go functions used in dynamic programming conditions of optimality. The result is a special case of a theorem characterizing local optimality conditions for the optimal control of general packet-switched radio networks, when either partial or complete observations of the state are available.[23]

**Theorem 1.** Suppose the control policies have complete observations of the state $(X_t, t \in [0,T])$ and, for every $U \in \mathcal{U}$, the state has a local description in terms of the conditional ($\wp^U, F_t$)-rates. Then $U=U^*$ is optimal in $\mathcal{U}$ if and only if there exist functions, $w_n(t, t_0, x_0, \ldots, t_n, x_n)$, measurable in their arguments and absolutely continuous in $t$, such that, $\tau_n$ the $n$-th transition time of $X$, and $e_{jks}$ is the $M_{\max}+1 \times M_{\max} \times S$ array with all zero entries except 1 in the $i, j, s$ position,

$$\frac{\partial w_n(t, \tau_0, X_0, \ldots, \tau_n, X_n)}{\partial t} + \min_{U \in \mathcal{U}}\left\{\sum_s \sum_{l=1}^K \left[\alpha_{s,t}(u_{0ls,t}) \cdot (w_{n+1}(t, \tau_0, X_0, \ldots, \tau_n, X_n + e_{0ls}) - w_n(t, \tau_0, X_0, \ldots, \tau_n, X_n))\right] + \right.$$

$$\sum_s \sum_{j=1}^K \sigma_{js,t} 1(Q_{js,t-} > 0)\left[\sum_{k=1}^K (u_{jks,t}) \cdot (w_{n+1}(t, \tau_0, X_0, \ldots, \tau_n, X_n + e_{jks}) - w_n(t, \tau_0, X_0, \ldots, \tau_n, X_n))\right] +$$

$$\left. \sum_s \sum_{m=1}^K \sigma_{ms,t} 1(Q_{ms,t-} > 0) \cdot (u_{m0s,t}) \cdot (w_{n+1}(t, \tau_0, X_0, \ldots, \tau_n, X_n + e_{m0s}) - w_n(t, \tau_0, X_0, \ldots, \tau_n, X_n)) + c(t, U_t)\right\} = 0 \quad (15)$$

for $\tau_n \leq t < \tau_{n+1}$, the energy constraints on $C_t$ in Figure 2 are satisfied at $t$, and

$$w_n(t, \tau_0, X_0, \ldots, \tau_n, X_n) = f_T \text{ for } \tau_n \leq T < \tau_{n+1}. \quad (16)$$

The minimum in (16) is attained at optimal policies $U_t^*(\omega)$ a.s. $\wp^{U^*}$. Moreover, the optimal cost-to-go process at time $t \in [0, T]$ takes the form

$$W_t = \sum_{n=0}^{\infty} 1(\tau_n \le t < \tau_{n+1}) \cdot w_n(t, \tau_0, X_0, \ldots, \tau_n, X_n). \qquad (17)$$

## 5.4 Optimality conditions with Markov assumptions

The optimality conditions in (15)-(16) greatly simplify when the routing policies $U$ depend only on the last observation of the network state to time $t$, i.e., $U_t(\omega) = U(t, X_{t-}(\omega))$ for each $\omega \in \Gamma$, $\Gamma \in F_t$ and the assumptions on the MVPPs underlying the network state are such that $X$ is a Markov process. This is the special case of Markov control of a Markov process. In general, these assumptions do not accurately represent the dynamics of MANET packet flows and interdependence of events occurring at the nodes and along paths, but they do allow conditions (15)-(17) to be expressed in terms of the *optimal cost-to-go function* $W_t = V_t = V_t(X_{t-})$ to yield dynamic programming conditions where the infinitesimal generator for $V_t$ depends on the conditional ($\wp^{U^*}$, $F_t$)-rates for MVPPs given in (6)-(8).

## 6. PERFORMANCE VERIFICATION

Performance of the optimized cross-layer protocols, determined by the dynamic programming conditions in (15)-(17) and energy-aware algorithm of Figure 2, can be evaluated in a network simulation environment to establish the QoS of multimedia streams under energy constraints. Such an environment is found in the wireless extensions made by the Monarch Group at Carnegie-Mellon University to the ns-2 simulator created by the VINT Project at the University of California at Berkeley.[34, 35] A canonical MANET scenario consists of a distributed collection of battery-powered laptops or hand-held terminals, capable of hosting multimedia applications. The optimality conditions and energy algorithm are embedded in a reactive routing protocol, such as, DSR or AODV, in a MANET of $M$ nodes, distributed over an $X$ m × $Y$ m area; $M$ may vary between 10 and 100, and the dimensions $X$ and $Y$ between 500 and 1500. Each node may have one or more buffers, capable of storing $P$ packets. The velocity of each node can vary dynamically between 0 and 15 m/s, using the random waypoint model provided in the mobility extensions to the ns-2 simulator. To simulate changing traffic load, simulation runs may begin with 2 TCP connections within the MANET and be increased to 30% of the number of nodes. Linear, convex, and concave delay functions can be made selectable on each link. Multimedia packet streams are constructed using a two-state Markov (on-off) model for voice with talkspurts, from a CBR source of 64 kbps; VBR video at 128, 256, 1024 and 2048 kbps; and and ABR Internet UDP data exhibiting self-similar behavior. Packet sizes can vary from 64 to 1500 bytes. In cases when node mobility is less than 0.3 m/s, the PHY-layer and MAC-layer specifications of the IEEE 802.11b,g LAN standards are used in the simulations; for node mobility above 0.3 m/s, the PHY-layer and MAC-layer specifications of the Third-generation Partnership Project (3GPP), wideband code-division multiple access (W-CDMA) wireless standard are applied. Data rates can thus vary from 64 kbps to 11 Mbps, depending on the standard in effect. Each node is assumed to have a battery with a capacity of $J$ joules. Manufacturers' estimated power consumption in the transmit, receive, and idle modes for typical WLAN cards and 3G mobile handsets are placed in a database for computation of residual battery capacities at the nodes and the total energy expenditures of the connections created by the cross-layer design. The expended energy calculation depends on packet size, bandwidth, radio frequency, and the terminal model. Performance evaluations are based on the application-layer metrics of throughput, packet loss rate, delay, and delay variation over the simulation run; the metrics are computed for each application type on each source-destination path and compared to the QoS guarantees for the type. Network configurations can be scored based on satisfaction of the QoS needs of applications active in packet streams, the balance of energy dissipation among the nodes, and network lifetime. Simulation runs may last 400s or until end of network lifetime, whichever occurs first. Details of the simulation implementation and the performance tradeoffs are not presented, but remain for a sequel to the paper.

## 7. CONCLUSIONS

In this paper, a QoS model is introduced that is derived from models recommended for wired Internet applications, but extended to the resource-limited, transient environment of an energy-constrained MANET of battery-powered nodes.

The claim is made that QoS support in MANETs is essentially different from traditional networks due this unique environment, which does not permit pre-planned resource allocations sufficient to ensure the QoS guaranteed by wired networks. A new definition of QoS for the energy-constrained MANET, adaptable to its unique environment, is introduced. The core of a protocol implementation is a cross-layer design based on this QoS model that supports adaptation and optimization across multiple layers of the protocol.

Real-time optimization of the cross-layer design is based on MVPP models of the packet flows for service applications on end-to-end connections in the MANET. The MVPP models encompass a wide range of statistical properties ascribed to packet arrivals, service processing, traffic shaping, and routing in previous studies of wireless multimedia networks. Control of modeled MANET behavior is realized through an absolutely continuous change of probability measure on the random, real-time packet-processing events. Controlled probability measures are constructed from a reference measure on network events via likelihood ratios that depend explicitly on the entries in admissible control arrays and the protocol parameters that influence packet routing, dropping, and blocking. In turn, the conditional random rates in the semi-martingale decompositions, with respect to a family of MANET observations and controlled measures, of the point processes also depend on the admissible routing entries and protocol parameters. The optimization problem thus becomes the determination of routing policies that minimize cost functionals of the parameters associated with the application, network, MAC, and physical layers of the protocol, subject to network and energy constraints at the nodes. With complete observations of network events, backward recursive optimality conditions are derived to characterize the admissible routing policies that provide the best available QoS performance of the energy-constrained cross-layer design. The effect of Markov assumptions on the controls and network state on the optimality conditions are presented. Methods of performance verification by simulation of the cross-layer designs are indicated.

Simulations of the QoS performance achieved by the cross-layer designs and parametric tradeoff analyses remain to be conducted. The simulations will determine if the optimization approach is scalable with the number of nodes and the effects of node mobility. Moreover, it has been established in other research that multi-user channel interference, influenced by MAC-layer tasks, exhibits self-similar behavior. This suggests that the MVPP models may be extended to represent this interference as a controllable state variable, rather than a constraint on protocol design.. Additional analysis is required to establish conditions on the network, performance criteria and control policies that allow separation of the optimal policies from the estimates of network state, given partial temporal and state observations of the MANET. The results of this analysis should reduce the extent of information required to support the optimal policies in the PHY, MAC and higher layers of the cross-layer protocol design.

## ACKNOWLEDGMENTS

## REFERENCES

1. W. S. Hortos, "Real-time performance analysis of wireless multimedia networks based on partially observed, multivariate point processes," *Digital Wireless Commun.II, Proc. SPIE*, 4045-05, Orlando, FL, April 2000.

2. E. Crawley, R. Nair, B. Rajagopalan, and H. Sandick, "A framework for QoS-based routing in the internet," *RFC 2386*, Aug. 1998.

3. "Packet radio networks," www.tapr.org/tapr/html/pktf.html.

4. A. J. Goldsmith and S. B. Wicker, "Design challenges for energy-constrained ad hoc wireless networks," *IEEE Wireless Comm.*, vol. 9, no.4, pp. 8 – 27, Aug. 2002.

5. X. Xiao and L. M. Ni, "Internet QoS: A big picture," *IEEE Network*, pp. 8–18, March/April 1999.

6. R. Braden, D. Clark, and S. Shenker, "Integrated services in the Internet architecture: an overview," *IETF RFC 1633*, June 1994.

7. S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, and W. Weiss, "An architecture for differentiated services," IETF RFC 2475, Dec. 1998.

8.  K. Wu and J. Harms, "QoS support in mobile ad hoc networks," *Crossing Boundaries – the GSA J. of Univ. of Alberta*, vol.1, no. 1, pp. 92 – 106, Nov. 2001.

9.  R. Braden, L. Zhang, S. Berson, S. Herzog, and S. Jamin, "Resource reservation protocol (RSVP) – version 1 functional specification," *IETF RFC 2205*, Sept. 1997.

10.  D. D. Perkins and H. D. Hughes, "A survey on quality-of-service support for mobile ad hoc networks," *Wireless Comm. and Mobile Comp. (WCMC),* vol. 2, no. 5, 2002.

11.  S. Shenker, C. Partridge, and R. Guerin, "Specification of guaranteed quality of service," *IETF RFC 2212*, Sept. 1997.

12.  D. Chalmers and M. Sloman, "A survey of QoS in mobile computing environments," *IEEE Comm. Surveys*, 1999.

13.  J. Wroclawski, "Specification of the controlled-load network element service," *IETF RFC 2211*, Sept. 1997.

14.  H. Xiao, W. K. G. Seah, A. Lo, and K. C. Chua, "A flexible quality of service model for mobile ad-hoc networks," *Proc. IEEE VTC-2000*, pp. 445 – 449, Tokyo, Japan, May 5 – 18, 2000.

15.  N. Nikaein, C. Bonnet, Y. Moret, and I. A. Rai, "2LQoS - two-layered quality of service model for reactive routing protocols for mobile ad hoc networks," *Proc. SCI - 6th World Multiconf. on Systemics, Cybern. and Informatics*, Orlando, Florida, July 14-18, 2002.

16.  N. Nikaein and C. Bonnet, "Layered quality of service model for routing in mobile ad hoc networks," *Proc. WMAN*, 2003.

17.  A. S. Acampora and M. Naghshineh, "Control and quality-of-service provisioning in high-speed microcellular networks," *IEEE Personal Comm. Mag.*, pp. 36-43, 1994.

18.  N. Nikaein and C. Bonnet, " A glance at quality of service models for mobile ad hoc networks," *Proc. DNAC 2002: 16th Conf. on New Arch. for Comm.*, Paris, France, 2002.

19.  L. Zhang, S. Deering, and D. Estrin, "RSVP: A new resource reservation protocol," *IEEE Network Mag.*, vol.31, no. 9, pp. 8 – 18, Sept. 1993.

20.  V. S. Frost and B. Melamed, "Traffic modeling for telecommunication networks," *IEEE Comm. Mag.,* vol. 32, no.3, pp. 70-81, Mar. 1994.

21.  W. Willinger, M. S. Taqqu, W. E. Leland, and D. V. Wilson, "Self-similarity in high-speed packet traffic: analysis and modeling of Ethernet traffic measurements," *Statist. Sci.*, vol.10, no. 1, pp.67-85, 1994.

22.  P. Brémaud, *Point Processes and Queues*: *Martingale Dynamics*, Springer-Verlag, New York, 1981.

23.  W. S. Hortos Jr., *Partially Observable Point Processes and the Control of Packet Radio Networks*, doctoral dissertation, University of Michigan, UMI Dissertation Information Services, Ann Arbor, May 1990.

24.  T. V. J. G. Babu, T. Le-Ngoc, and J. F. Hayes, "Performance of a priority-based dynamic capacity allocation scheme WATM systems," *Proc. IEEE GLOBECOM '98*, vol. 4, pp. 2234-2238, Nov.1998.

25.  W. E. Leland, M. S. Taqqu, W. Willinger, and D. V. Wilson, "On the self-similar nature of Ethernet traffic (extended version)," *IEEE/ACM Trans. Networking*, vol. 2, no.1, pp.1-15, Feb.1994.

26.  M. Parulekar and A. Makowski, "Tail probabilities for a multiplexer with self-similar traffic," *Proc. IEEE INFOCOM '96*, vol. 3, pp. 1452-1459, Mar. 1996.

27.  A. T. Andersen and B. F. Nielsen, "An application of superpositions of two-state Markovian sources to the modelling of self-similar behaviour," *Proc. IEEE INFOCOM '97*, vol. 1, pp. 196-204, Apr. 1997.

28.  N. Likhanov, B. Tsybakov, and N. D. Georganas, "Analysis of an ATM buffer with self-similar ("fractal") input traffic," *Proc. IEEE INFOCOM '95*, vol. 3, pp. 985-992, Apr.1995.

29.  W. M. Lam and G. W. Wornell, "Multiscale representation and estimation of fractal point processes," *IEEE Trans. Sig. Process.*, vol. 43, no. 11, pp. 2606-2617, Nov. 1995.

30.  M. A. Krishnam, A. Venkatachalam and J. M. Capone, "Self-similar point process through a fractal construction," *Proc. IEEE INFOCOM 2000*, Tel Aviv, Israel, 8 pages.

31.  M. Grossglauser and J.-C. Bolot, "On the relevance of long-range dependence in network traffic," *IEEE/ACM Trans, Networking*, vol. 7, no.5, pp. 629-640, Oct. 1999.

32. D-K. Kim, J-W. Park, C-K Toh, and Y-H. Choi, "Power-aware route maintenance protocol for mobile ad hoc networks," *Proc. IEEE 10th Int. Conf. Telecom. (ICT) 2003*, pp. 501-506, Papeete, French Polynesia, Feb. 2003.

33. C. Doléans-Dade and P.A. Meyer, "Intégrales stochastique par rapport aux martingales locales," *Séminaire Probabilités, IV, Lecture Notes in Mathematics*, vol. 124, Springer-Verlag, Berlin, pp. 77-107, 1970.

34. "Network simulator version 2, ns-2," http://www-mash.cs.berkeley.edu/ns/.

35. "CMU Monarch Project, wireless and mobility extensions to ns-2," http://www.monarch.cs.cmu.edu/cmu-ns/html.