# A New Control Architecture With Enhanced ARP, Burst-Based Transmission, and Hop-Based Wavelength Allocation for Ethernet-Supported IP-Over-WDM MANs

Jingxuan Liu, *Student, IEEE,* and Nirwan Ansari, *Senior Member, IEEE*

*Abstract*—This paper focuses on the control architecture and the enabling technologies for the Ethernet-supported Internet protocol-over-wavelength-division-multiplexing metropolitan area networks. We present the general architecture of an access node of such networks and propose solutions to facilitate the essential system functionalities. The aim is to render the flexible and high-capacity metropolitan network, which provides service provisioning improvement and resource utilization efficiency for the packet-dominated data traffic. Specifically, an enhanced address resolution protocol is proposed to reduce the call setup latency and the signaling overhead associated with the address probing procedure, a burst-based transmission mechanism is adopted to improve the network throughput and resource utilization efficiency, and a wavelength allocation algorithm is investigated to provide flexible bandwidth multiplexing with fairness and high scalability. Theoretical analysis and simulations are conducted to evaluate the performance of our algorithms, demonstrating that the proposed architecture and technologies deliver substantial transport performance improvement with efficient network resource utilization.

*Index Terms*—Address resolution protocol (ARP), Ethernet, traffic assembly, wavelength allocation, wavelength-division-multiplexing (WDM).

## I. INTRODUCTION

**T**HE Ethernet-supported Internet protocol (IP)-over-wavelength-division-multiplexing (WDM) ring paradigm provides a feasible solution for the new generation metropolitan optical network (MAN), enabling a graceful migration from the current voice-oriented MAN prototype into a world optimized for packets. On one hand, the WDM technology approaches the metro bottleneck problem with a clean slate. Its multichannel concurrent transport capability and the accommodation for transparent traffic delivery make the WDM-enabled network an ideal platform to deliver the traffic of mixed type applications. Meanwhile, the Ethernet holds great promise as the connectivity solution in the data traffic-dominated metropolitan area network. While the optical transport network is targeted for packet processing to reduce the complexity of multiplayer architectures, a smooth evolution is critical to success. Given that the traffic from the enterprise into the MAN is primarily Ethernet data and will continue to dominate in the future, support for standard Ethernet interfaces directly on network elements enables native service interfaces for data transport [1]–[3].

In addition to the system infrastructure considerations, efficient media access control (MAC) protocols are needed to coordinate the system resources, in order to unleash the potential of the packet-based WDM metropolitan ring network and make the most of the ring prototype in an intelligent way. This is especially true when the data packet processing in the optical domain is still not yet mature. Ethernet, while a natural fit for data traffic, has evolved to support full duplex-switched infrastructures but lacks the flexible MAC mechanisms to manage the access across multiple users in the WDM ring prototype. The foregoing challenges require innovative protocols and algorithms in the metropolitan environment that retain all the advantages of the packet-based transport mechanism, while rendering elastic bandwidth allocation and graded levels of services.

Intensive research endeavors have been devoted in the design and implementation of the MAN ([4]–[11] and references therein). The effort has been focused on either the network architectures and configurations, or the MAC protocols and the service provisioning mechanisms. For example, Hernandez-Valencia [4] presented a hybrid architecture consisting of Ethernet/time-division-multiplexing (TDM) service solutions to enable storage networking and Ethernet transport over synchronous optical network (SONET)/synchronous digital hierarchy (SDH) networks. Madamopoulos *et al.* [5] investigated the impact of different add–drop modules implementations on the system performance. Cai *et al.* [8] proposed the multitoken interarrival time (MTIT) protocol to provide the efficient bandwidth multiplexing for the WDM ring architecture. Alternative MAC protocols focusing on the fairness, scalability, and various service provisioning have also been reported.

In this paper, we architect the access nodes of the Ethernet-supported IP-over-WDM metropolitan network, and devise the enabling technologies combining the space, time, and channel domains to systematically facilitate the new network infrastructure. Specifically, we propose mechanisms comprised of the address resolution, the traffic engineering, and the wavelength allocation to render packet-optimized optical MAN with the transport performance improvement (e.g., reduced transport delay), the resource utilization efficiency

(e.g., improved network throughput), and fairness and scalability for the network resource access control. Our proposals have the following characteristics. First, the mechanisms for service improvement are devised at the access node of the Ethernet-supported IP-over-WDM MAN, where packet processing techniques are mature and affordable. Second, the transport latency reduction and signaling overhead alleviation are implemented in both the Ethernet layer and the WDM layer by taking into consideration the new network characteristics. Third, a novel wavelength allocation algorithm, in conjunction with the burst-based transmission scheme, facilitates a fair and highly scalable MAC solution to deliver efficient optical bandwidth multiplexing and network throughput. Fourth, the data traffic transmission is supported in a non-slotted single-hop fashion to benefit the dynamic data traffic with predictable transport performance after being launched into the metropolitan optical ring. Last but not least, our proposals are class-of-service (CoS)-friendly and are fully compatible with the existing technologies.

Our contributions include the following.

1) We present a simple and scalable access node architecture, which embraces the Ethernet functionality, and supports the flexible data transmission at the WDM-enabled metropolitan network.

2) We propose an enhanced address resolution protocol (E-ARP) based on the existing standard for the address probing functionality suitable for the new network prototype.

3) We advocate the burst-based transmission mechanism to improve the network throughput and facilitate the CoS provisioning at the access node.

4) We design a hop-based wavelength allocation algorithm to coordinate the parallel transmission of different data channels, and eliminate or reduce the transmission collisions caused by wavelength contentions.

Besides the theoretical analysis, extensive simulations have been conducted to evaluate the system performance.

The rest of the paper is organized as follows. Section II describes the system scenario, the control architecture of an access node, and the problem statement. Section III proposes the system mechanisms, which are essential to the Ethernet-supported IP-over-WDM network prototype, including the E-ARP scheme, the burst-based transmission scheme, and the wavelength allocation principle. The discussion for scalability, fairness, and CoS-capability are also presented. In Section IV, we investigate the system performance by theoretical analysis and simulation evaluations. The concluding remarks are made in Section V.

## II. SYSTEM ENVIRONMENT AND PROBLEM STATEMENT

This section details the network environment upon which our investigation will be conducted. The general control architecture of an access node will be presented, and the design objectives of our proposed enabling technologies will be formulated.

### A. System Environment

We consider a ring-shaped metropolitan network where $N$ access nodes are interconnected via counterrotating dual fibers



Fig. 1. Prototype of a ring-based metropolitan optical network. (a) Dual-fiber ring. (b) Access node connecting the feeder ring and the LAN.

(i.e., the feeder ring) [14], as shown in Fig. 1. The fiber ring consists of the inner ringlet and the outer ringlet, each of which makes use of the full bandwidth of the fiber, i.e., the individual wavelength can be transported concurrently in both ringlets, assuming that each access node has adequate receiver capabilities. Each fiber supports $W + 1$ wavelengths as parallel channels, of which $W$ wavelengths ($\lambda_1, \ldots, \lambda_W$) are for data channels and one for the control channel ($\lambda_0$). If necessary, the network capacity can be gradually updated by parallel fibers. The aggregated bandwidth can scale to multiterabits/second.

There are vast research results on the problem of providing bandwidth multiplexing and channel access control for the packet-oriented WDM networks. The majority of the approaches centers on the WDM layer. The implementation complexity, cost, and performance have impact on the network design. Interested readers are referred to [7]–[11] for detailed discussions and to [14] for an overview. In this paper, the efficient bandwidth sharing of the optical fiber is achieved from the perspectives of the signaling transmission and the space reuse of wavelengths.

Access to the network resources (e.g., the data transmission channel) is typically based on two alternative schemes: preallocation-based protocols and reservation-based protocols. While the former technique assigns transmission rights to different nodes in a static and predetermined manner, the reservation-based technique arbitrates the bandwidth access to the traffic demand in a real-time fashion, i.e., the resource reservation request is delivered throughout the ring layout when a data transmission is required. In our scenario, we adopt the reservation-based method as it yields flexible bandwidth utilization and is a natural fit for the dynamic data traffic. The data payload is launched into the network after the corresponding reservation request is confirmed success.

Two or more data transmissions on the same wavelength along the same section of the fiber result in a collision. Depending on the approach the network resource contention is addressed, signaling protocols discussed in the literature can be classified into two main categories which are: 1) collision-free strategies and 2) collision-and-retry strategies. These variants result in different complexity of network hardware requirements, and different optical bandwidth multiplexing efficiency. In this paper, we assume that the control packet transmission is decoupled from that of the data payload. While a control packet
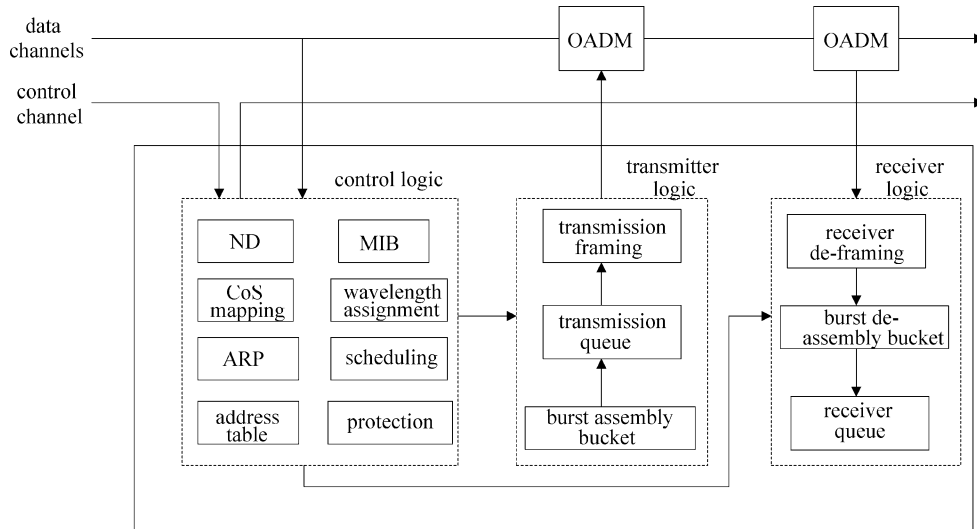
Fig. 2.    Functional architecture of an access node.

reserves the resources according to the collision-and-retry strategy, the data payload is guaranteed with a single-hop transmission without any delay or loss in the feeder ring.

In addition, we employ the destination-stripping method to extract the data traffic from the ring network. This method, together with the wavelength reuse of the data transmission (owing to the wavelength allocation algorithm as described in Section III), enables the concurrent data channel usage in disjoint source-destination pairs, yielding an improved degree of bandwidth multiplexing and resource utilization efficiency within the ring.

One or more gigabit Ethernet (GbE) local area networks (LANs) are connected to the feeder ring via the access nodes. The LAN is typically composed of enterprises or high-speed end users. Hereafter, we will refer to the $j$th router in the LAN attached to the access node $A_i$ ($i \in \{0, \ldots, N-1\}$) as $A_i.R_j$, assuming that $j \in \{0, \ldots, M-1\}$, where $M$ is the total number of routers in the attached LANs. Besides the connectivity functionality, the access node also provides MAC solutions, which are necessary to render efficient packet-oriented traffic processing and transmission at the WDM layer. The detailed architecture and functionality of the access node will be described in the next subsection. Throughout this paper, we will refer to the traffic flowing from an LAN to the metropolitan ring as the upstream traffic, while that from the metropolitan ring to the LAN as the downstream traffic.

### B. Access Node Architecture

Each access node has two interfaces: the GbE interface, which is used for the packet processing and the interaction with the associated local access networks, and the optical link interface to access the WDM ring in the optical domain. Fig. 2 illustrates the functional architecture of an access node, which mainly consists of the transmitter logic, the receiver logic, and the control logic.

The main function of the transmitter logic is to adapt the low bit-rate tributaries (up to OC-12) of the local network to the transmission granularity suitable for WDM transport media, and to forward the traffic to the destination node. In our system,

the transmitter logic assembles the upstream data packets into the larger granularity, namely, a data burst. The traffic assembly mechanism is similar to that proposed for the backbone [15], [16], and is tailored for the Ethernet-supported WDM ring topology, as will be discussed in the next section.

The transmitters emit the assembled traffic into the metropolitan network. A transmitter may be either fully tunable to all data channels, or it may be "partially tunable" to more than one but fewer than $W$ wavelengths [10]. In our system, the number of tunable transmitters is impacted by the wavelength allocation algorithm (as will be discussed in the next section). Multiple transmitters of an access node can transmit on different wavelengths concurrently with negligible tuning latency [12], [13], enabling the better exploit of the parallel transmission capability of the WDM technology. The transmitter queues are provided to accommodate the traffic waiting for the access to the data channels.

The receiver logic is similar to the transmitter logic. Each node is equipped with tunable receivers, which can be fully tunable or partially tunable just like the transmitters. The downstream traffic is disassembled into packets before being transported to the individual host in the LAN. The receiver queue is provided to order the processing of the traffic according to their CoS.

The control logic coordinates the traffic transmission and processing, and facilitates the traffic engineering functionalities. Of particular interests in this paper are the address resolution protocol (ARP) module, the burst assembly control module, and the wavelength allocation module.

The ARP module is designed to support the address probing mechanism proposed in this paper. Each access node is embodied with two tables: the local address mapping (LAM) table registering the $\langle$protocol type, protocol address, physical address$\rangle$ triplet of all the routers in the subordinated LANs, and the remote address mapping (RAM) table recording the address translation results obtained from the executed address inquiry procedures.

The traffic assembly module controls the traffic assembly/disassembly procedure. To enable the connectivity services with

graded levels of performance, the CoS mapping function should be incorporated in this module to map the incoming traffic flow into a specific transport class.

The wavelength allocation module, in conjunction with the scheduling module and the traffic selection module, enables the coordination between the control packet and the data payload, as well as resolves resource contentions. The implementation of such control functionalities requires the management information base (MIB) and other functional components to keep track of the network status, such as the data channels availability (full or empty) and the configuration status of the data channel connections.

Fig. 2 also highlights some functional modules which are necessary for the multifacet network infrastructure. For example, the node discoverer (ND) module maintains the topology information of all the nodes in the ring, and monitors the fiber cut or node failure events by periodically sending out-of-band signaling message to its neighboring nodes, both of which are essential to deploy the protection/restoration strategies. Due to the space limitation, however, the implementation of such components is beyond the scope of this paper.

The designed access node architecture has the following properties.

1) It facilitates the Ethernet-supported IP-over-WDM integration to support the dynamic data traffic, thus overcoming the inefficiency and inflexibility of the current SONET-dominated infrastructure and its circuit-based provisioning model.

2) It allows the deployment of mechanisms in both the Ethernet and WDM perspectives, thus enabling the better synergy of both mature electronic protocols and advanced optical technologies.

3) It is flexible and cost effective. The individual functional component is a skeleton based on which a variety of algorithms and devices can be adopted and developed according to the preference of the network management. Meanwhile, the service provisioning requires no SONET overhead and synchronization among access nodes, nor the dedicated protection bandwidth.

Theoretical analysis and simulation results will demonstrate that our access node architecture, in tandem with the proposed enabling technologies, yield improved network performance with reduced system cost.

### C. Problem Statement

Based on the aforementioned service requirements and network architecture, we can formulate our enabling technology design issue as follows. Given the Ethernet-supported IP-over-WDM system supporting asynchronous, variable length packets, the problem we are facing now is to design protocols and algorithms, which are indispensable for the new generation MAN to deliver improved service provisioning and network performance with reduced system cost. A brief summary of the design objectives is as follows.

1) Improve the service provisioning for the application streams in terms of the transport latency, including both the call setup delay and the transmission delay in the metropolitan network.

2) Improve the system efficiency in terms of the resource utilization, the network throughput, and the signaling transmission requirements.

3) Provide fairness and scalability control among the traffic between the access nodes, as well as differentiated classes of services for the multitype applications.

We facilitate these requirements with mechanisms developed at both the data link layer and the medium layer, based on the presented architecture.

### III. ENABLING TECHNOLOGIES

In this section, we speculate the principle of our proposed enabling technologies, and discuss their impacts to the Ethernet-supported IP-over-WDM metropolitan network. Notations are defined in Table I to simplify our description.

### A. Enhanced Address Resolution Protocol (E-ARP)

Typically, an Ethernet-supported network employs the ARP [17] to translate the network layer address (i.e., the IP protocol address) into the link-layer one (i.e., the hardware address). While very simple and well-suited to the LAN-hardened Ethernet (which is broadcast in nature), the original ARP cripples the address probing procedures in the metropolitan optical network. Consider routers $A_1 \cdot R_1$ and $A_1 \cdot R_2$, both of which need to communicate to router $A_2 \cdot R_2$. Using a conventional ARP, both senders execute individual call setup procedures on the overlapped route from $A_1$ to $A_2$, and $A_2$ to $A_2 \cdot R_2$ [see Fig. 3(a)], resulting in longer call setup delay and higher consumption of network resources. A savvy address inquiry function is highly desirable for the new network scenario.

We propose an enhanced ARP, called E-ARP, to reduce the call setup latency and the gratuitous ARP packet transmissions. Our basic idea is twofold: the address translation is rendered as early as possible, and the address mapping obtained from previous ARP transmissions is retrievable for the subsequent inquiries. We describe the E-ARP by concentrating on its distinctive characteristics as compared with the conventional ARP.

1) The access nodes incorporate the address translation function with the packet forwarding function. Upon the reception of an upstream ARP request (i.e., an ARP packet destined for a router in a remote LAN), the access node either directly replies it, or broadcasts it in the MAN, depending on whether or not the access node can find the required address mapping information in the RAM table. Meanwhile, the access node updates—when necessary—the hardware address field of the sender's entry in the LAM table with the information in the packet, or adds a new entry if the sender does not exist in the LAM. A downstream ARP request (i.e., an ARP packet broadcasted in the metropolitan ring) is received by all access nodes, and is replied only by the access node, which connects the destination router to the MAN (i.e., the access node whose LAM table has the entry indexed by the destination protocol address).

2) The access node re-edits the upstream ARP request before broadcasting it into the metropolitan ring, replacing

TABLE I
NOTATIONS FOR THE PROPOSED ENABLING TECHNOLOGIES ($i \in \{1, \ldots, N\}, l \in \{1, \ldots, +\infty\}, k \in \{1, \ldots, h_{\max}\}$)

| Term | Explanation |
|---|---|
| $pro$ | Protocol type of the data packet |
| $tha$ | Hardware address of the target of this packet |
| $tpa$ | Protocol address of the target of this packet |
| $sha$ | Hardware address of the sender of this packet |
| $spa$ | Protocol address of the sender of this packet |
| $t_d^a$ | The time when a data burst begins to assemble at the access node. |
| $t_d^s$ | The time when an upstream data burst is sent into the metropolitan network. |
| $t_c^s(l)$ | The time when a control packet is sent into the metropolitan network for the $l$-th time. |
| $t_c^r$ | The time when an access node receives a control packet. |
| $t_r(l)$ | The random time when a control packet is delayed at the source access node after the $l$-th reservation attempt fails. |
| $\tau_a$ | The burst assembly duration |
| $T_w$ | The data burst delay owing to signaling transmissions. |
| $T_c$ | The round-trip time for a control packet to be processed in the MAN. |
| $R$ | The maximum number of resource reservation attempts before a data burst is dropped. |
| $t_e$ | The time the specific data channel will be available at the destination access node. |
| $h_{\max}$ | The maximum number of hops a data burst is transported in the ring. |
| $H$ | The total number of hops to support concurrent transmissions between any pair of access nodes. |
| $h_k$ | The number of hops for a transmission to propagate from the source to the destination nodes. |
| $C_t$ | The number of data channels required to support the concurrent transport between any pair of access nodes. |
| $C$ | The total number of data channels available in the system. |
| $S_k$ | The $k$-th data channel subset shared by the traffic which traverses $k$ hops. |
| $|S_k|$ | The number of data channels in the data channel subset $S_k$. |
| $\omega_k^j$ | The $j$-th data channel in the subset $S_k$ ($j \in \{1, \ldots, |S_k|\}$). |
| $\omega_k^{pref}(A_i)$ | The most preferred data channel for node $A_i$ to transport the traffic requiring $k$ hops. |

the sender hardware address (*sha*) field with its own hardware address. The access node is also responsible to generate the uni-cast ARP reply packet. Note that the traffic is exchanged in the MAN ring according to the hardware address of the access node.

3) The address mapping information obtained from the ARP reply packet is registered in the RAM table of the access node, and is retrievable by the subsequent address translation requests, which are sent by other local routers. The redundant ARP transport on the metropolitan ring is thus avoided until the RAM table ages out the entry corresponding to the remote router.

Fig. 3(b) depicts the ARP packet flow when our E-ARP is adopted, given the same address resolution requests as that in Fig. 3(a). The detailed E-ARP procedure is shown in Fig. 4. Note that the enhanced ARP also supports the basic address information management [17], e.g., the address update and address age-out for RAM and LAM tables. Such functionalities are omitted in Fig. 4 for simplicity.

Our E-ARP mechanism features several advantages. First, the RAM table in the access node facilitates the information reuse of address inquiries, thus prompting the call setup procedure and reducing the unnecessary ARP packet transmissions in the metropolitan optical network. Second, the LAM table provides the access node with the address information of its local routers. Therefore, an access node can reply an ARP request without going further into the LAN. Third, our E-ARP maps the protocol address of a router to the hardware address of its associated access node. This way, data packets are transported in the MAN according to the hardware address of the access node, and are delivered to the ultimate routers by the local access node. This solution ushers in the decoupling of the traffic transmission in the WDM domain from that in the Ethernet domain, which is consistent with the line of thought to operating in the individual network independently. Meanwhile, addressing the traffic according to the access node also benefits our system with reduced complexity for traffic management and traffic engineering (e.g., the traffic assembly procedure which will be explained in the next subsection). Another significant merit of our E-ARP is the performance improvement achieved without requiring new protocols. The address resolution enhancement is implemented by simply equipping the access node with the RAM and the LAM tables.

### B. Burst-Based Transmission Mechanism

After the address translation procedure, the subsequent data traffic can be transported according to the obtained hardware
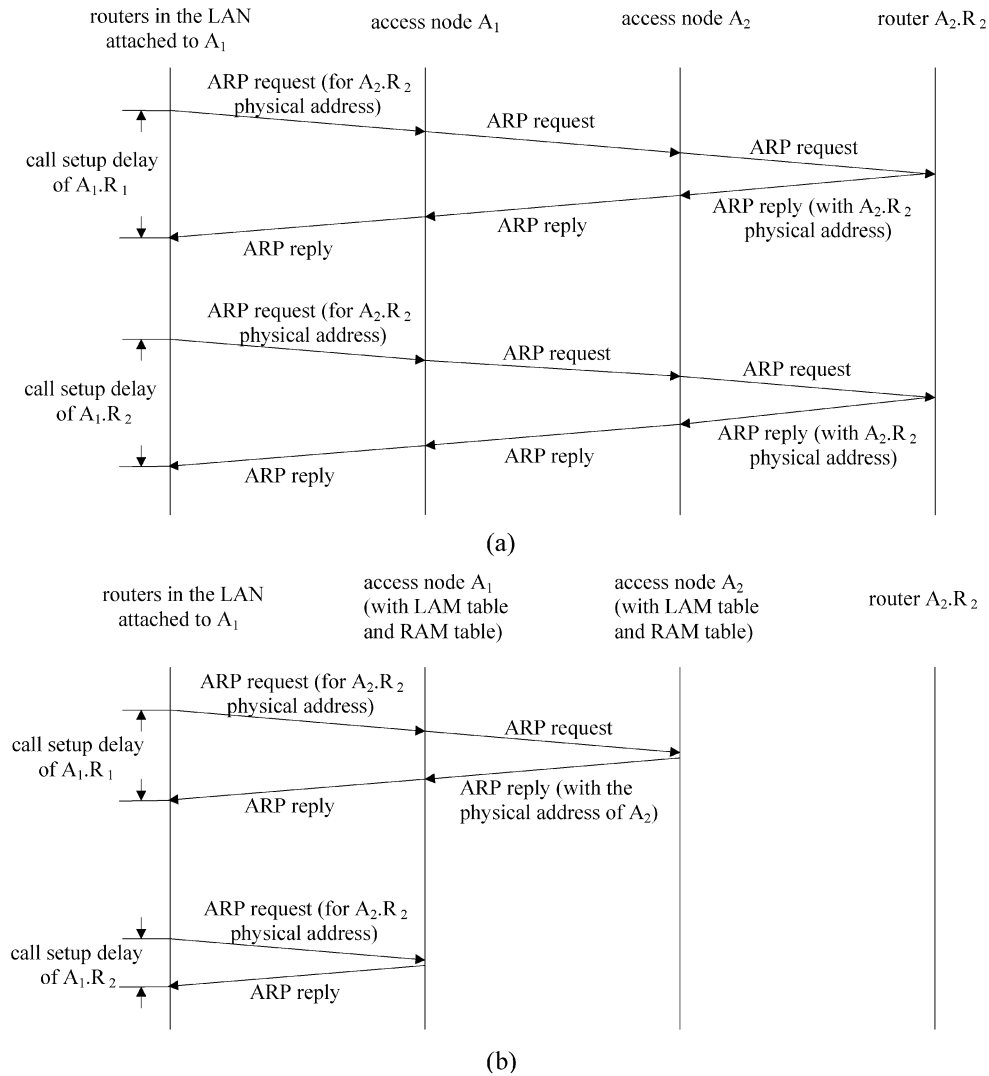
Fig. 3. ARP and the E-ARP mechanisms. (a) Packet flow of the ARP packets. (b) Packet flow of the E-ARP packets.

address. The traffic transmission mechanism plays an important role in architecting the efficient WDM-enabled integration. Motivated by the mismatch between the transmission capacity of WDM fibers and the fine traffic granularity of data packets, our proposed system adopts a modified burst-based transmission mechanism to improve the system throughput and reduce the signaling overhead. We describe our transmission mechanism in terms of the traffic assembly procedure, the signaling scheme, and the burst-based transmission benefits.

We base our burst assembly procedure on the one proposed in [16] and [18], which has proved to be an efficient paradigm in the long-haul network. Different from their proposal, however, we assemble packets at the access node according to the Ethernet address of the destination access node, a mechanism enabled by the proposed E-ARP. The individual access node is equipped with $N-1$ assembly units, each corresponding to one destination access node on the ring topology. To enable differentiated services for the incoming traffic, each unit has one or more assembly buckets corresponding to different CoS requirements. Broadcast service can be easily supported by our burst assembly mechanism in alternative ways: either add one specific assembly unit at each access node to assemble the broadcast

packets so that a single copy of each packet is sent throughout the metropolitan network, or $N-1$ copies of the individual unicast packet are replicated, each inserted in one of the $N-1$ assembly buckets.

The burst assembly interval ($\tau_a$) is subject to the constraints imposed by the round-trip time of the control packet (i.e., the cycle latency for a control packet to travel throughout the WDM ring, denoted as $T_c$) and by a predetermined maximum burst length. In the metropolitan network adopting the reservation-based channel access mechanism, $\tau_a$ should be no less than $T_c$ to upper bound the traffic generation rate by the traffic processing rate. By limiting the maximum burst length with a certain threshold, a simple fairness control is rendered to avoid one transmission occupying a certain channel for an excessive long time and potentially starving the other transmission requirements competing for the same channel or destination node.

The signaling protocol in our burst-based transmission mechanism distinguishes from its counterpart in the backbone in the two-way signaling scheme. When the data burst is fully assembled, the associated control packet is generated and transported into the network. After being processed at each access node and attempting to reserve resources at the destination node, the

If the triplet $< pro, tpa, tha >$ exists in the RAM table,

   returns the *tha* to the sender

else

   { re-writes the ARP request packet by replacing *spa* and *sha* fields

       with the respective values of the current access node;

   broadcasts the re-written ARP request packet on the WDM ring}

(a)

If the *tpa* carried in the ARP request packet is not in my LAM table,

   discards the ARP request packet

else

   {Add the triplet $< pro, spa, sha >$ to the RAM table

   Swap *spa* and *sha* in the ARP request packet with its *tpa* and *tha*, respectively,

      putting the hardware address of the current access node in the *sha* field;

   mark the ARP packet with '*REPLY*'

   send the ARP reply packet to the ring;

   }

(b)

Fig. 4. ARP mechanism. (a) Access node receives an upstream ARP request packet. (b) Access node receives a downstream ARP request packet.



Fig. 5. Transmission coordination between the control packet and the data payload when the resource reservation is invoked after the burst assembly procedure is completed. (a) Reservation succeeds. (b) Reservation fails.

control packet returns to the source node with a reservation acknowledgment. Depending on whether the corresponding reservation succeeds or fails, the data payload is either emitted into the ring network, or it is delayed at the transmission queue. Padding may be required if a minimum burst length is imposed [16]. The unsuccessful control packet will be retransmitted after a random time delay until the reservation succeeds, or the maximum transmission attempts $(R)$ are performed, with the random delay between two consecutive retransmissions uniformly distributed from 0 to $T_c$. A data burst is discarded when its control packet fails to reserve the resource after $R$ attempts. $R$ should be properly engineered based on the tradeoff consideration between the traffic drop probability and the delay allowance.

When the fast service provisioning is of essence to the network management, our signaling protocol incorporates the parallel execution of the burst assembly procedure and the resource reservation, whereby the control packet transmission is triggered as soon as a burst begins assembling. The concurrent execution of both delay contributors reduces the inherent artificial delay of data traffic at the access node. Figs. 5 and 6 illustrate the time line when the resource reservation and the burst assembly procedure are performed sequentially and concurrently, respectively, and Fig. 7 the principle of resource reservations.

Maintaining the advantages of the typical burst-based transmission mechanism, our two-way signaling protocol avoids the data traffic retransmissions in the WDM domain. Once a data



Fig. 6. Transmission coordination between the control packet and the data payload when the resource reservation and the burst assembly procedure are performed in parallel. (a) Reservation succeeds. (b) Reservation fails.



Fig. 7. Resource reservation at the access node. (a) Reservation fails. (b) Reservation succeeds.

payload is launched into the ring topology, it is guaranteed to deliver without further delay or loss caused by resource contentions. Meanwhile, the signaling scheme accommodates the data payload with a single-hop transmission, which provides security and privacy transport services, and is convenient for CoS guarantees. Data payloads are propagated to the destination access node on the ringlet with the shorter hops. This way, the maximum number of hops a data burst is transported in the ring is $h_{\max} = \lceil (N-1)/2 \rceil$. Moreover, the reservation acknowledgment also carries the information of the network status which can be exploited for dynamic adjustment of the system parameters (e.g., the burst assembly interval $\tau_a$, the delay period between consecutive signaling retransmissions, and maximum resource reservation attempts $R$), thus enhancing the traffic engineering capability.

Our burst-based transmission mechanism is simple and efficient in that no complex determination for the offset time is involved, and that the value-added delayed reservation [19] mechanism can be easily implemented at the destination access node. Both system parameters (the offset time and the delayed reservation interval) are inherently equal to the cycle latency of the control packet. Such simplification is benefited from the ring topology of the metropolitan network. Meanwhile, the burst-based transmission mechanism enables us to engineer the data traffic at the access node according to the Ethernet address of the access node, the CoS requirements, and the multicast service. This solution complies with the *de facto* trend that only simple and straightforward processing is performed in the WDM layer, while most of the intelligence of the network is provided in the electronic domain. In other words, the burst-based transmission and the Ethernet-supported WDM metropolitan network are dual-benefited mechanisms.

### C. Wavelength Allocation Algorithm

Besides the service provisioning and traffic engineering functionalities, the proposed network also features a hop-based wavelength allocation algorithm for efficient bandwidth utilization and contention resolution.

Our basic idea is to partition the bandwidth capacity of $W$ data wavelengths into the disjoint subsets $S_k$ ($k = 1, \ldots, h_{\max}$), each containing a group of data channels, and being shared among the transmission demand with the same hop numbers. For example, the traffic sourced from access node $A_i$ ($i \in \{0, \ldots, N-1\}$) and destined for $A_{i+k}$ ($k \leq h_{\max}$) share the same subset of data channels with that sourced from access node $A_{i+1}$ and destined for $A_{i+k+1}$. Herein, the data channel consists of one wavelength or a fraction of a wavelength. Assembling the data burst based on the hardware address of the destination access node enables the source access node to easily determine the number of hops that the traffic needs to be propagated.

By definition, we have $S_k = \left\{ \omega_k^j \middle| k = 1, \ldots, h_{\max}, j = 0, \ldots, |S_k| \right\}$. We also observe that on the individual ringlet, the total number of sessions required for all pairs of access nodes to communicate is

$$H = \begin{cases} \sum_{k=1}^{h_{\max}} k, & N \text{ is an odd integer} \\ \sum_{k=1}^{h_{\max}} k - \left\lfloor \frac{h_{\max}}{2} \right\rfloor, & N \text{ is an even integer} \end{cases} . \quad (1)$$

Therefore, the minimum number of data channels required to support concurrent transmissions between all pairs of access nodes is $C_t = H$. Our wavelength allocation principle features the following characteristics.

1) The number of data channels allocated to the subset $S_k$ is determined based on the associated transport distance (in hops) and the total available data channels, defined by

$$|S_k| = k \cdot \frac{C}{C_t}. \quad (2)$$

The contention-free wavelength allocation is theoretically achievable if $C \geq C_t$, i.e., $|S_k| \geq k$.

2) The data channel assignment can also take into consideration the reservation contention tolerance, the transport latency constraints, and the estimated traffic demand of the individual node, in order to facilitate the service differentiation.

3) Among the data channels of the subset $S_k$, the one chosen for the access node $A_i$ to communicate with $A_{i+k}$ can be determined via a variety of algorithms. For example, it can be selected randomly with a probability of $(|S_k|)^{-1}$. Alternatively, the data channel can be determined in a cyclic fashion, as shown in Fig. 8. The most preferred channel for $A_i$ to transport traffic requiring the hop number $k$ is determined by

$$\omega_k^{\text{pref}}(A_i) = \omega_k^{i \bmod |S_k|}. \quad (3)$$

Table II illustrates one implementation of our wavelength allocation algorithm for a system scenario, where $N = 9, W = 10$, and the data channel is equal to the individual wavelength, i.e., $C = W$.

The proposed algorithm provides a fair and scalable MAC mechanism. The transport demands requiring the same number

```
set pickedChannel = −1;
if preferredChannel available;
        set pickedChannel = ω_i^pref;
else
    For ( l = 0;  l < |S_k|;  l++) {
        pickedChannel = (preferredChannel + l) mod |S_k|;
        if pickedChannel available
        return pickedChannel
    }
if pickedChannel = −1;
return (No channel available).
```

Fig. 8.   Channel selection algorithm at the source access node.

TABLE II
EXAMPLE IMPLEMENTATION OF THE WAVELENGTH ALLOCATION ALGORITHM

| subset | $S_1$ | $S_2$ | $S_3$ | $S_4$ |
|---|---|---|---|---|
| Channel assignment | $\omega_1^1 = \lambda_1$ | $\omega_2^1 = \lambda_2$ | $\omega_3^1 = \lambda_4$ | $\omega_4^1 = \lambda_7$ |
| | | $\omega_2^2 = \lambda_3$ | $\omega_3^2 = \lambda_5$ | $\omega_4^2 = \lambda_8$ |
| | | | $\omega_3^3 = \lambda_6$ | $\omega_4^3 = \lambda_9$ |
| | | | | $\omega_4^4 = \lambda_{10}$ |

of hops of propagation equally share the same group of data channels, regardless of the index of the source or the destination access nodes. Meanwhile, our subset-based wavelength allocation is advantageous in terms of computational simplicity. Its time complexity for data channel selection is $O(\max(|S_k|))$, $k = 1, \ldots, h_{\max}$.

## IV. PERFORMANCE ANALYSIS AND SIMULATION RESULTS

The system performance is evaluated via theoretical analysis and simulation results. Interested performance metrics include the network throughput, the resource reservation blocking probability, the burst drop probability, and the transport latency. We will investigate their dependency on the number of access nodes, the maximum burst size, and the traffic intensity.

Simulations are based on the uniform traffic scenario, whereby the upstream traffic at each access node is destined for the other nodes with equal probabilities. The interarrival time of the data packets flowing into the access nodes are exponentially distributed. The metropolitan network consists of a regional size ring with a circumference of 200 km. The cycle latency of a control packet is approximately 1000 $\mu$s. Table III summarizes the notations we will use in the analysis.

### A. Network Throughput

Among many others, the property and efficiency of a system design are characterized by the network throughput, which is defined as the average data traffic (in bits) successfully transmitted by all the access nodes in a unit time. We analyze the impact of the burst-based transmission mechanism and the hop-based wavelength allocation on the network throughput. In the presence of the balanced network traffic scenario, the transport

TABLE III
NOTATIONS FOR PERFORMANCE ANALYSIS
$(i \in \{1, \ldots, N\}, l \in \{1, \ldots, +\infty\})$

| Term | Explanation |
|------|-------------|
| $D_c$ | The transport latency owing to the call setup procedure. |
| $B_0$ | The basic bandwidth capacity of one wavelength. |
| $X_{\max}$ | The maximum throughput achievable by the system. |
| $X$ | The achievable network throughput. |
| $X_o$ | The network throughput normalized to the transport capacity of a MAN adopting SONET with the same configuration. |
| $P_f^k$ | The probability that a control packet fails to reserve resources for a transport requirement with hop number $k$. |
| $P_f$ | The reservation blocking probability. |
| $P_d$ | The burst drop probability. |
| $d_i$ | The depth of the LAN attached to the access node $A_i$. |
| $t_l$ | The propagation time to transport a packet on one link of the LAN. |
| $t_m$ | The propagation time to transport a packet on one link of the MAN. |
| $D_c$ | The transport distance owing to the call setup procedure when the system adopts the E-ARP. |
| $D_c$ | The transport distance owing to the call setup procedure when the system adopts the traditional ARP. |
| $i$ | The latency reduction capability enabled by the E-ARP. |
| $\overline{T}_w^s$ | The average signaling delay for a transmission requirement of $k$ hops when the control packet and the burst assembly procedure are executed in sequence. |
| $\overline{T}_w^p$ | The average signaling delay for a transmission requirement of $k$ hops when the control packet and the burst assembly procedure are executed in parallel. |
| $D_b^s$ | The average data burst delays induced by the burst-based transmission mechanism when the resource reservation and the burst assembly procedure are executed in sequence. |
| $D_b^p$ | The average data burst delays induced by the burst-based transmission mechanism when the resource reservation and the burst assembly procedure are performed in parallel. |

request of different hop numbers is uniformly distributed. Theoretically, our system enables a throughput of

$$X_{\max} = 2 \cdot \sum_{k=1}^{h_{\max}} \left( B_c \cdot |S_k| \cdot \frac{N}{k} \right) \quad (4)$$

where $B_c$ is the bandwidth capacity of a data channel given by $B_c = (W \cdot B_0)/H$.

The actually achievable network throughout is impacted by two factors: the number of the concurrently transported data burst limited by the signaling mechanism, and the signaling/data transmission ratio. The network throughput imposed by these constraints can be derived to be

$$X = N \cdot (N - 1) \cdot \beta_1 \cdot \beta_2 \quad (5)$$

where $\beta_1$ and $\beta_2$ represent the ratio of the average reservation length ($\overline{L}$) over the cycle latency of the control packet on the ring ($T_c$), and that of the transported data bursts over the total reservation attempts, respectively. In a system with $W = H$, the network throughput normalized to the transport capacity of an MAN adopting SONET with the same number of wavelengths for data transport can be expressed as

$$X_o = \begin{cases} 8 \cdot \frac{N}{N+1} \cdot \beta_1 \cdot \beta_2, & N \text{ is an odd integer} \\ 8 \cdot \frac{N-1}{N} \cdot \beta_1 \cdot \beta_2, & N \text{ is an even integer} \end{cases} \quad (6)$$
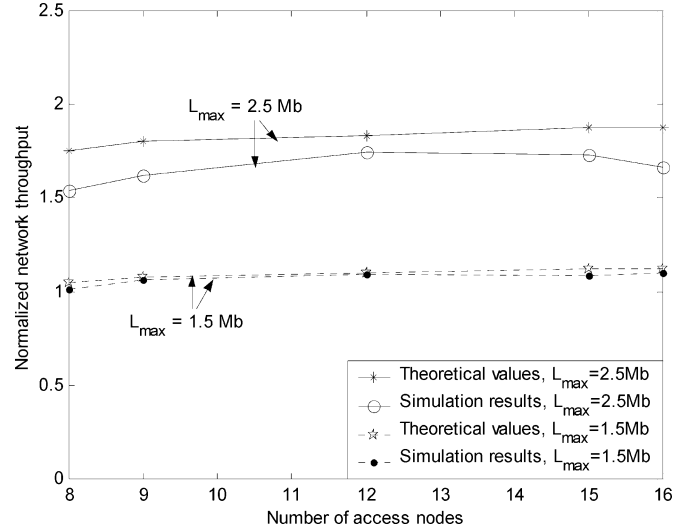


Fig. 9. Normalized achievable network throughput versus the number of access nodes.

provided the total number of sessions in the network is $H = (1/8) \cdot (N^2 - 1)$ and $H = (1/8) \cdot N^2$ for the odd and the even integer of $N$, respectively [based on (1)].

Fig. 9 shows the achievable network throughput versus the number of access nodes when the maximum burst size varies. The theoretical values are obtained based on (6), assuming an ideal implementation (i.e., the resource reservation succeeds at the first signaling attempt). We conduct simulations based on two scenarios, whereby the data bursts are transported based on the loss-free fashion with the maximum burst size ($L_{\max}$) being 2.5 Mb and the loss-subjective fashion with $L_{\max} = 1.5$ Mb, respectively. It appears that increasing the maximum burst size improves the achievable network throughput (e.g., for $N = 15$, $X_o = 1.13$ and 1.73, when $L_{\max} = 1.5$ and 2.5 Mb, respectively), and that while the experimental $X_o$ matches the theoretical values very well in the loss-subjective transport scenario, the two curves slightly diverge from each other in the loss-free alternative, revealing the impact of reservation blockings.

The simulation results in both scenarios imply the following conclusions. First, our system delivers high throughput despite the large size of the ring or the large number of access nodes, indicating the efficient resource utilization capacity and high traffic volume supportability. Second, while incurring no burst loss by retransmitting the control packets until the reservation succeeds, too many re-reservation attempts result in the large signaling response time at the access node, as well as the head-of-line (HOL) delay to the subsequent data bursts, both of which contribute to the lower network throughput. Dimensioning the maximum resource re-reservation attempts entails a tradeoff between the burst loss probability and the network throughput. Third, assembling the input packets into the larger transport granularity improves the throughput capacity with reduced signaling transmission requirement. This is consistent with the conclusion that the larger burst size within a certain range, the higher network throughput [20]. When the burst size keeps growing, however, the performance improvement slows down because the negative impact of the resource contention increases. This is also demonstrated in Fig. 10, which illustrates
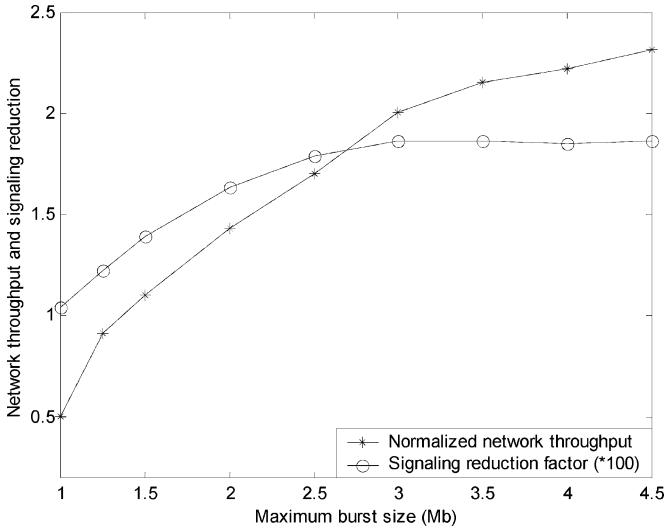
Fig. 10.  Impact of the maximum burst size (*megabit*) on the normalized achievable network throughput ($X_o$) and the signaling reduction factor ($f$), when $N = 16$ and $R = 3$.

the impact of $L_{\max}$ on the network throughput and the signaling reduction factor (i.e., the ratio of the signaling transmission requirements and the control packet processing complexity in a system without the assembly procedure over those in the proposed system), given that the maximum number of attempts for signaling retransmission is limited. Fig. 10 reinforces our previous conclusion that the network throughput is improved when $L_{\max}$ increases within a reasonable range.

### B. Reservation Blocking Probability

The reservation blocking occurs when the data channel carried by a control packet is not available at the intermediate nodes between the source and destination access nodes. Assuming that the probability density function (PDF) for the interarrival time of the reservation requests received at the access node is $f(t)$, and that the data channel is randomly selected from each subset with equal probability, the reservation blocking probability at an intermediate node is $P_c = \int_0^{L_{\max}} f(t)dt$. Therefore, the resource reservation for a data burst propagating for $k$ hops is blocked with the probability of

$$P_f^k = 1 - \left(1 - \frac{1}{|S_k|} \cdot P_c\right)^{k-1}. \qquad (7)$$

By averaging the blocking probability over all the possible distances, we have the average resource reservation blocking probability of

$$P_f = 1 - \frac{\sum_{k=2}^{h_{\max}} \left(1 - \frac{1}{|S_k|} \cdot P_c\right)^{k-1}}{h_{\max} - 1}. \qquad (8)$$

Fig. 11 presents the performance of $P_f$ and the burst drop probability ($P_d$) versus the number of access nodes. The concurrent execution of the burst assembly and the control packet transmission is adopted (see Fig. 6), wherein the resource reservation length (in time) is fixed to be the time to support the maximum burst size ($L_{\max}$). Our system possesses very low reservation blocking probability with marginal burst drop probability.
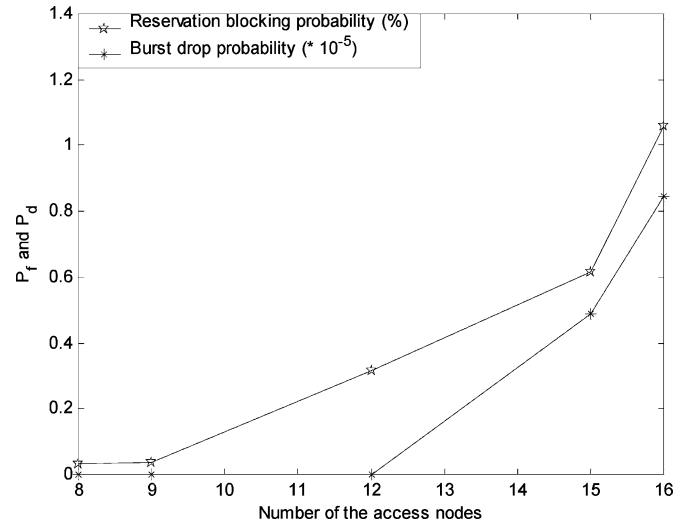


Fig. 11.  Reservation blocking probability ($P_f$) and burst drop probability ($P_d$) versus the number of access nodes. $R = 3$, $L_{\max} = 1.5$, the input packet intensity is 0.6, and the burst assembly and the control packet transmission are conducted concurrently.
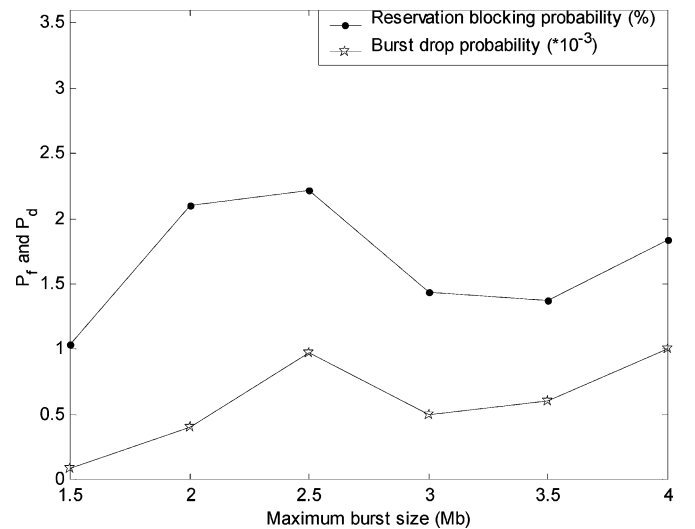


Fig. 12.  Impact of the maximum burst size on reservation blocking probability ($P_f$) and burst drop probability ($P_d$). $R = 3$, $N = 16$, the input packet intensity is 0.6, and the burst assembly and the control packet transmission are conducted concurrently.

For example, $P_d = 4.9 * 10^{-5}$ and $P_f = 0.6\%$, when $N = 15$ and $R = 3$. Note that the data burst is subjected to drop only at the access nodes of the metropolitan network, where the retransmission is controlled by the data link layer. In addition, since the control packet reserves resources for the duration enough to support the maximum burst size, our simulation results can be viewed as the conservative evaluation of the system performance.

$P_f$ and $P_d$ also present high dependence on $L_{\max}$, as shown in Fig. 12. In our system scenario, the larger maximum burst size allowance results in the decreased resource reservation requirements, and the increased reservation length (in time), both of which have the contradictive impact on the reservation blocking probability. Algorithms to determine the optimal maximum burst size, combining other constraints such as the
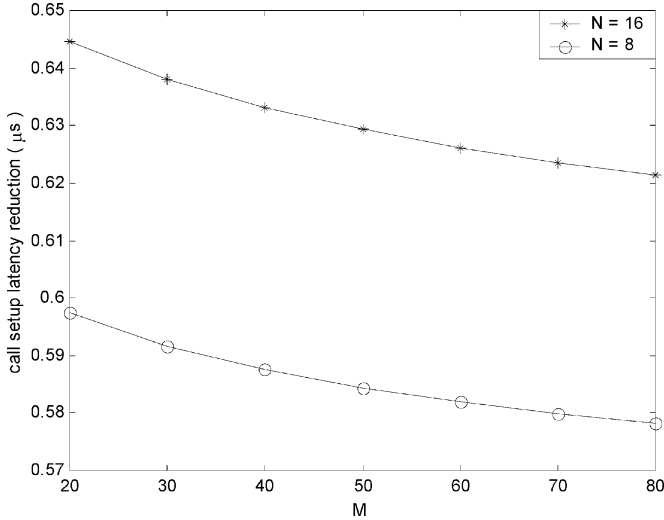
Fig. 13. Latency reduction capability of the E-ARP with respect to the ARP.

data channel availability, are critical issues and need further investigation.

### C. Transport Latency

In this paper, the transport latency involves the delay experienced by the call setup procedure, and that by the data traffic at the source access node. To focus on the effect of the E-ARP and the burst-based transmission mechanism, we do not consider the delay due to the control packet generation and data channel determination.

*1) Call Setup Latency:* To analyze the latency reduction capability of the E-ARP, we assume that the LAN has a binary tree topology, wherein the root is the access node connecting the tree to the MAN and the leaf is typically a router from a corporate or campus network. Assuming $M_i \geq M_j$ and the distance (in hops) between $A_i$ and $A_j$ is $h_k$, the total distance (in hops) required for call setup packets to build up the full communication between all the routers in the LAN $A_i$ and those in the $A_j$ is

$$D_c = 2 \cdot [(d_i + h_k) + d_i \cdot (M_i - 1)] \cdot M_j. \qquad (9)$$

When the conventional ARP is adopted, the distance for the same communication requirement is

$$D_c' = 2 \cdot [(d_i + h_k + d_j) \cdot M_i] \cdot M_j. \qquad (10)$$

Without loss of generality, we assume each LAN has the same depth ($d_i = d_j$) and that the ARP packet propagation on one link in LAN is equal to that on one link in the ring (i.e., $t_l = t_m$). The average latency reduction delivered by the E-ARP for node $A_i$ to communicate with all other nodes is approximately

$$\alpha_i = \frac{\sum_{k=1}^{h_{\max}} \left( 1 - \frac{k + M_i \cdot \log_2 M_i}{(2 \cdot \log_2 M_i + k) \cdot M_i} \right)}{h_{\max}}. \qquad (11)$$

Our E-ARP yields significant latency reduction for the call setup procedures as compared with the conventional ARP (Fig. 13). This is especially true when $N$ is large, or when $M$ is relatively small. By making the address mapping information retrievable

for the subsequent inquiries at the access nodes, our E-ARP obtains the larger $\alpha_i$ as $N$ increases, and reduces the signaling overhead for address translations (i.e., the transmission requirements (in hops) of the call setup packets) at the WDM feeder ring by $(M_i - 1)/M_i$, as implied in (10) and (11). In our system scenario, the latency reduction percentage decreases as the size of the LAN gets larger, when $M_i \gg N$, $\alpha_i \to 1/2$.

*2) Data Burst Latency:* Our proposed network architecture introduces two main sources that will cause the data burst delay at the access node: the delay caused by the burst assembly procedure, which is given by $D_a = (1/2) \cdot \tau_a$, and the potential delay because of unsuccessful resource reservations. Since the $i$th control packet retransmission induces a burst delay of $(T_c + t_r(i))$, the average signaling delay for a transmission requirement of $k$ hops is

$$\overline{T}_w^s = T_c + \sum_{i=1}^{R} i \cdot (T_c + t_r(i)) P_f^k(i) \qquad (12)$$

for a system wherein the control packet and the burst assembly procedure are executed in sequence, and

$$\overline{T}_w^p = \sum_{i=1}^{R} i \cdot (T_c + t_r(i)) P_f^k(i) \qquad (13)$$

when the two procedures are performed in parallel. The average data burst delays are, thus

$$D_b^s = \frac{1}{2} \cdot \tau_a + T_c + \frac{1}{h_{\max} - 1} \cdot \sum_{k=2}^{h_{\max}} \sum_{i=1}^{R} i \cdot (T_c + t_r(i)) P_f^k(i) \qquad (14)$$

and

$$D_b^p = \frac{1}{2} \cdot \tau_a + \max\{0, T_c - \tau_a\} + \frac{1}{h_{\max} - 1} \cdot \sum_{k=2}^{h_{\max}} \sum_{i=1}^{R} i \cdot (T_c + t_r(i)) P_f^k(i) \qquad (15)$$

respectively. In an ideal implementation, the average data burst delay is simply $D_b^s = D_a + T_c$ or $D_b^p = D_a$ when $\tau_a \geq T_c$.

To evaluate the burst delay performance, we examine both signaling mechanisms: the sequential transmission of the control packet and the data burst (see Fig. 5), wherein the resource is reserved according to the actual data burst length, and the parallel alternative (see Fig. 6) with the reservation length fixed to be the maximum burst length. The simulation results with respect to the number of access nodes and the input traffic intensity are shown in Figs. 14 and 15, respectively. Given that in the sequential signaling scheme, the total signaling response time should also account for the cycle latency of the control packet which is approximately 1000 $\mu$s (not shown in Fig. 15), the parallel signaling scheme is more favored for applications with the stringent time constraint. Our system yields satisfactory transport delay for today's voice and video interactive applications with a stringent end-to-end latency bound requirement of tens of milliseconds. Besides the data burst latency at the access node, Fig. 15 also presents the other performance metrics interested in this paper.
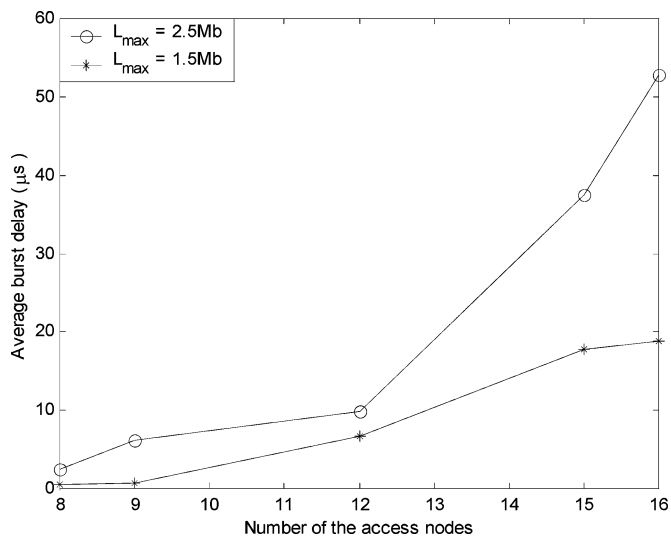
Fig. 14. Average data burst delay versus the number of access nodes when the burst assembly and the control packet transmission are conducted sequentially. The reservation length is fixed to be the maximum burst size. $R = 3$, $N = 16$, and the input packet intensity is 0.6.
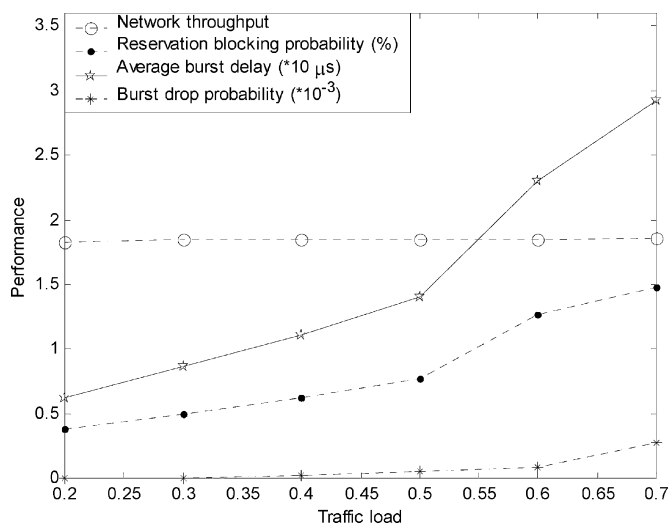


Fig. 15. Performance figures of merits versus the traffic load when the burst assembly and the control packet transmission are executed in parallel. The control packet reserves the resources according to the actual data burst length. $R = 3$, $N = 16$, and $L_{max} = 2.5$.

## V. CONCLUSION

This paper has addressed the control architecture and enabling technologies for the Ethernet-support IP-over-WDM metropolitan network. We have proposed a set of mechanisms by jointly considering both the Ethernet layer and the WDM layer that provide network service improvement to the real-time variable-length traffic. The proposed technologies include the enhanced address resolution protocol (E-ARP), the burst-based transmission scheme, and the hop-based wavelength allocation algorithm. CoS differentiation and broadcast services implementation have been discussed.

Theoretical analysis and simulations exhibit encouraging results. The E-ARP protocol significantly reduces the transport latency and the transmission requirements for the call setup procedures. The burst-based traffic transmission is proved to be a feasible and effective paradigm for the metropolitan optical network. In addition, we have investigated a novel wavelength allocation algorithm to fairly arbitrate the channel access among all the access nodes. The enhancement of the Ethernet services, in tandem with the innovative mechanism in the WDM domain, facilitates a flexible and efficient prototype for the new generation metropolitan optical network dominated by the packet traffic.

There are many important issues remaining open and worthy of further investigation to improve the Ethernet-supported IP-over-WDM metropolitan network. For example, the important parameters for the burst assembly procedure (including the maximum burst size and the burst assembly interval) should be carefully dimensioned for traffic classes with different CoS requirements. The service differentiation can also be facilitated in the data channel allocation scheme by adjusting the number of channels assigned to each traffic class [see (2)]. Another issue which will be investigated is the impact of the long-range dependence traffic characteristics on the system performance. These topics and the related work are the focus of our future research.

## REFERENCES

[1] E. Hernandez-Valencia, "Hybrid transport solutions for TDM/data networking services," *IEEE Commun. Mag.*, vol. 40, pp. 104–112, May 2002.

[2] G. Kramer and G. Pesavento, "Ethernet passive optical network (EPON): Building a next-generation optical access network," *IEEE Commun. Mag.*, vol. 40, pp. 66–73, Feb. 2002.

[3] Y. Xue. (2002) Optical network service requirements. IETF. [Online]. Available: http://www.ietf.org/internet-drafts/draft-ietf-ipo-carrier-requirements-05.txt

[4] M. Scholten, Z. Zhu, E. Hernandez-Valencia, and J. Hawkins, "Data transport applications using GFP," *IEEE Commun. Mag.*, vol. 40, pp. 96–103, May 2002.

[5] N. Madamopoulos, D. C. Friedman, I. Tomkow, and A. Boskovic, "Study of the performance of a transparent and recondifurable metropolitan area network," *J. Lightwave Technol.*, vol. 20, pp. 937–954, June 2002.

[6] D. Stoll, P. Leisching, H. Bock, A. Richter, and S. Ag, "Metropolitan DWDM: A dynamically configurable ring for the KomNet field trial in Berlin," *IEEE Commun. Mag.*, vol. 39, pp. 106–113, Feb. 2001.

[7] I. Chlamtac, V. Elek, A. Fumagalli, and C. Szabo, "Scalable WDM access network architecture based on photonic slot routing," *IEEE/ACM Trans. Networking*, vol. 7, pp. 1–9, Feb. 1999.

[8] J. Cai, A. Fumagalli, and I. Chlamtac, "The multitoken interarrival time (MTIT) access protocol for supporting variable size packets over WDM ring network," *IEEE J. Select. Areas Commun.*, vol. 18, pp. 2094–2104, Oct. 2000.

[9] K. Bengi and H. R. Van As, "Efficient QoS support in a slotted multihop WDM metro ring," *IEEE J. Select. Areas Commun.*, vol. 20, pp. 216–227, Jan. 2002.

[10] A. C. Kam and K. Siu, "Supporting bursty traffic with bandwidth guarantee in WDM distribution networks," *IEEE J. Select. Areas Commun.*, vol. 18, pp. 2029–2040, Oct. 2000.

[11] M. A. Marsan, A. Bianco, E. Keonardi, M. Meo, and F. Meri, "MAC protocols and fairness control in WDM multirings with tunable transmitters and fixed receivers," *J. Lightwave Technol.*, vol. 14, pp. 1230–1244, June 1996.

[12] I. P. Kaminow *et al.*, "A wideband all-optical WDM network," *IEEE J. Select. Areas Commun.*, vol. 14, pp. 780–799, June 1996.

[13] K. Zhu and B. Mukherjee, "Traffic grooming in an optical WDM mesh network," *IEEE J. Select. Areas Commun.*, vol. 20, pp. 122–133, Jan. 2002.

[14] C. S. Jeiger and J. M. H. Elmirghani, "Photonic packet WDM ring networks architecture and performance," *IEEE Commun. Mag.*, vol. 40, pp. 110–115, Nov. 2002.

[15] J. S. Turner, "WDM burst switching for petabit data networks," in *Proc. Optical Fiber Commun. Conf.*, vol. 2, 2000, pp. 47–49.32.

[16] Y. Xiong, M. Vandenhoute, and H. C. Cankaya, "Control architecture in optical burst-switched WDM networks," *IEEE J. Select. Areas Commun.*, vol. 18, pp. 1838–1851, Oct. 2000.

[17] D. C. Plummer, "Ethernet address resolution protocol: Or converting network protocol addresses to 48 bit Ethernet address for transmission on Ethernet hardware," IETF, Request for comments (Standard) RFC826, 1982.

[18] A. Ge, F. Callegati, and L. S. Tamil, "On optical burst switching and self-similar traffic," *IEEE Commun. Lett.*, vol. 4, pp. 98–100, Mar. 2000.

[19] M. Yoo, C. Qiao, and S. Dixit, "Optical burst switching for service differentiation in the next-generation optical Internet," *IEEE Commun. Mag.*, vol. 39, pp. 98–104, Feb. 2001.

[20] J. Cai and A. Fumagalli, "An analytical framework for performance comparison of bandwidth reservation schemes in WDM rings," in *Proc. IEEE INFOCOM 2002*, vol. 1, June 23–27, 2002, pp. 41–47.

**Jingxuan Liu** (S'98) received the B.S. and M.S. degrees from Beijing University of Posts and Telecommunications, Beijing, China, in 1996 and 1999, respectively. She is currently working toward the Ph.D. degree in computer science at the New Jersey Institute of Technology (NJIT), Newark.

Her research interests include queueing theory and its applications to computer communications, network architecture, and design, control, and performance analysis of IP-over-optical networks.

Ms. Liu is a Student Member of the IEEE Communications Society. She received second place in the IEEE North Jersey Section Student Presentation Contest.

**Nirwan Ansari** (S'78–M'83–SM'94) received the B.S.E.E. (*summa cum laude*) degree from the New Jersey Institute of Technology (NJIT), Newark, in 1982, the M.S.E.E. degree from the University of Michigan, Ann Arbor, in 1983, and the Ph.D. degree from Purdue University, West Lafayette, IN, in 1988.

He joined the Department of Electrical and Computer Engineering, NJIT, as an Assistant Professor in 1988, and has been a Full Professor since 1997. He coauthored *Computational Intelligence for Optimization* (Norwell, MA: Kluwer, 1997) with E. S. H. Hou, and it was translated into Chinese in 2000, and coedited *Neural Networks in Telecommunications* (Norwell, MA: Kluwer, 1994) with B. Yuhas. He has frequently been invited to give talks and tutorials. He is a Technical Editor of the *IEEE Communications Magazine*, as well as the *Journal of Computing and Information Technology*. His current research focuses on various aspects of multimedia communications and high-speed networks.

Dr. Ansari was a Distinguished Speaker at the 2004 Sendai International Workshop on Internet Security and Management, and a Keynote Speaker at the IEEE/ACM International Conference on E-Business and Telecommunication Networks (ICETE2004). He initiated (as the General Chair) the First IEEE International Conference on Information Technology: Research and Education (ITRE 2003), was instrumental, while serving as its Chapter Chair, in rejuvenating the North Jersey Chapter of the IEEE Communications Society, which received the 1996 Chapter of the Year Award and a 2003 Chapter Achievement Award, served as the Chair of the IEEE North Jersey Section and in the IEEE Region 1 Board of Governors during 2001–2002, and currently serves in various IEEE committees. He was the 1998 recipient of the NJIT Excellence Teaching Award in Graduate Instruction, and a 1999 IEEE Region 1 Award.