

Throughput-Competitive Advance Reservation with Bounded Path Dispersion

Reuven Cohen¹, Niloofar Fazlollahi² and David Starobinski²

¹Dept. of Mathematics, Bar-Ilan University, Ramat-Gan 52900, Israel

²Dept. of Electrical and Computer Engineering

Boston University, Boston, MA 02215

Email: reuven@macs.biu.ac.il, {nfazl,staro}@bu.edu

Abstract—In response to the high throughput needs of grid and cloud computing applications, several production networks have recently started to support advance reservation of dedicated circuits. An important open problem within this context is to devise advance reservation algorithms that can provide provable throughput performance guarantees, independently of the specific network topology and arrival pattern of reservation requests. In this paper, we first show that the throughput performance of greedy approaches, which return the earliest possible completion time for each incoming request, can be arbitrarily worse than optimal. Next, we introduce two new on-line, polynomial-time algorithms for advance reservation, called `BatchAll` and `BatchLim`. Both algorithms are shown to be throughput-optimal through the derivation of delay bounds for $1 + \epsilon$ bandwidth augmented networks. The `BatchLim` algorithm has the advantage of returning the completion time of a connection immediately as a request is placed, but at the expense of looser delay performance than `BatchAll`. We then propose a simple approach that limits path dispersion, i.e., the number of parallel paths used by the algorithms, while provably bounding the maximum reduction factor in the transmission throughput. We prove that the number of paths needed to approximate any flow is quite small and never exceeds the total number of edges in the network. Through simulation for various topologies and traffic parameters, we show that the proposed algorithms achieve reasonable delay performance, even at request arrival rates close to capacity bounds, and that three to five parallel paths are sufficient to achieve near-optimal performance.

I. INTRODUCTION

Current networking technologies are typically based on the best-effort TCP/IP packet switching paradigm, wherein packets from various sources share links. While this paradigm makes perfect sense for most standard network applications, due to the benefits of statistical multiplexing, this is not the case for scientific and other resource intensive applications [1]. Indeed, it has been observed that, when using packet switching, a large gap prevails between the theoretical network capacity and the maximum end-to-end throughput achieved by those applications [2]. This gap, which is mainly attributed to the shared nature of Internet traffic, has become increasingly problematic for modern grid and cloud computing applications requiring the transfer of petabyte file sizes and beyond.

To support high-end applications, the scientific community has been devoting significant efforts to develop a new protocol stack based on the concept of *advance reservation* [2, 3]. The most important property of advance reservation is to offer hosts and users the ability to reserve in advance *dedicated* paths to connect their resources. Advance reservation services have been successfully tested in a number experimental projects, such as UltraScienceNet [2] and OSCARS [1, 4], and are now offered by network providers, such as Internet2 and ESnet4 [5], in complement to their traditional IP-based routing services.

Several protocols and algorithms have been proposed in the literature to support advance reservation (cf. Section II). However, to the authors' knowledge, none of them provides throughput performance guarantees. Instead, most are based on heuristics or greedy approaches, whereas each request is allocated a path guaranteeing the earliest completion time at the time the request is placed.

In this paper, we first uncover fundamental limitations of greedy algorithms. Specifically, we show that there exists network topologies and request patterns for which the maximum throughput of these algorithms can be arbitrarily smaller than the optimal throughput, as the network size grows.

Next, we propose a new polynomial-time advance reservation algorithm, called `BatchAll`, that provably achieves maximum throughput for *any* network topology and pattern of request arrivals. Instead of immediately reserving a path for each incoming request as in a greedy algorithm, `BatchAll` accumulates several arrivals in a batch and assigns a more efficient set of flow paths to the whole batch, based on a maximum concurrent flow optimization. We prove the throughput-optimality of `BatchAll` by providing a bound on the ratio of the maximum delay experienced by any request in an $1 + \epsilon$ augmented resource network to the maximal delay experienced by any request in the original network using the optimal algorithm.

The `BatchAll` algorithm does not return the connection completion time to the user at the time a request is placed, but only when the connection actually starts. We therefore propose another throughput-optimal algorithm, called `BatchLim`, which provides the completion time immediately as a request is placed, but at the expense of a slightly looser bound on the maximum delay ratio. This algorithm operates by limiting the

A shorter, preliminary version of this paper appeared in the proceedings of the High-Speed Networking Workshop at IEEE INFOCOM 2008.

length of each batch.

Our model assumes that the network infrastructure supports path dispersion, i.e., multiple paths in parallel can be used to route traffic, which has recently been shown to be feasible on the ESnet network [6]. Nevertheless, a too large path dispersion may be undesirable, as it may entail fragmenting a file into a large number of segments and reassembling them at the destination. To address this issue, we present a simple approach, based on the max-flow min-cut theorem, that limits the number of parallel paths while bounding the maximum reduction factor in the transmission throughput. We prove that this bound is tight. We then, propose two algorithms `BatchAllDisp` and `BatchLimDisp`, based upon `BatchAll` and `BatchLim` respectively. These algorithms perform similarly to the original algorithms, in terms of the batching process. However, after filling each batch, the algorithms will limit the dispersion of each flow. Although these algorithms are not throughput-optimal anymore, they are still throughput-competitive.

Finally, we provide simulation results illustrating the delay performance of the various proposed algorithms as a function of the arrival rate of requests. The results are compared to a capacity bound, whose value represents an upper bound on the maximum request arrival rate for which the delay of requests is bounded. We show that `BatchAll` approaches the capacity bound at a reasonable delay value. Further, we propose and evaluate a modified version of the algorithm, called `BatchAll+`, whose performance is even closer to the bound. With respect to path dispersion, we show that excellent performance can be achieved with as few as five or so parallel paths per connection.

The rest of this paper is organized as follows. In Section II, we scan related work on advance reservation, competitive approaches, and path dispersion. In Section III, we introduce our model and define performance measures of interest. In section IV, we present natural greedy approaches and demonstrate their inefficiency. Then, in Section V, we describe the `BatchAll` and `BatchLim` algorithms and prove their throughput-optimality. Our method for bounding path dispersion is described and analyzed in Section VI. In Section VII, simulation results evaluating the performance of the algorithms for different network topologies and traffic parameters are presented. We conclude the paper in Section VIII.

II. RELATED WORK

There exists a rich literature on advance reservation of network resources. Thus, we focus on algorithmic work, which is most relevant to this paper. We refer the interested reader to [7], for further discussion on architectures and services.

Several papers consider the problem of joint routing and scheduling of file transfers. Ref. [8] introduces a scheduling algorithm for large file transfers over paths with varying bandwidth. Ref. [9] analyzes the problem of offline scheduling and routing of file transfers from several users, each storing multiple files, to a single receiver node. Ref. [10] considers a similar model but also proposes algorithms to address the problem of rescheduling connections that have not completed.

Ref. [11] proposes greedy algorithms guaranteeing earliest completion to each incoming request. Ref. [3] introduces an algorithmic framework, called GCR, that assigns grades to paths according to a general optimization objective, e.g., shortest path, widest path, earliest path, or a combination of those, and then selects the highest graded path. GCR also supports path switching, which is shown to enable significant performance improvement.

Advance reservation plays an important role in optical networks, such as those based on optical flow switching (OFS) architectures [12, 13]. Thus, a number of papers investigate joint routing-scheduling optimization with the additional constraint of maintaining wavelength continuity along a path. For instance, [14] focuses on the design of effective routing and wavelength assignment heuristics and [15, 16] propose various load balancing approaches to allocate lightpaths.

In contrast to the above work, the algorithms that we propose in this paper provide theoretical guarantees on throughput performance. These results are obtained under the following relaxations: (i) job requests do not have strict bandwidth requirements; (ii) no additional constraints (e.g., wavelength continuity) are imposed. Hence, our model is appropriate for elastic applications, such as bulk transfer, operating under a single lambda service [2].

Next, we summarize work related to competitive algorithms for advance reservation. Ref. [17] discusses the on-line ftp problem, where reservations are made for channels for the transfer of different files. The presented algorithm provides a 4-competitive algorithm for the makespan (the total completion time). However, [17] focuses on the case of fixed routes. When routing is also to be considered, the time complexity of the algorithm presented there may be exponential in the network size. Ref. [18] proposes an off-line approximation algorithm for accommodating advance reservation of job requests in some specific topologies, e.g., lines and trees. Ref. [19] studies a related problem, whereby the path selection is based on several choices supplied by a user. Ref. [20] focuses on the problem of joint routing and scheduling in packet switched network in an adversarial setting and provides delay bounds. The model is for packet scheduling, rather than file (job) scheduling as in our paper. Moreover, [20] puts some restrictions on the adversary (e.g., the maximum amount of packets that it can inject during any time interval), which are unnecessary in our paper.

Most of the other work on competitive routing algorithms mainly focuses on call admission, without the ability of advance reservation [21, 22]. In [23], a combination of call admission and circuit switching is used to obtain a routing scheme with a logarithmic competitive ratio on the total revenue received. A competitive routing scheme in terms of the number of failed routes in the setting of ad-hoc networks is presented in [24]. A competitive algorithm for admission control and routing in a multicasting setting is presented in [25].

In [26], a throughput-optimal scheme for packet switching at the switch level is presented. The scheme is based on a convergence to the optimal multicommodity flow using delayed decision for queued packets. Their results somewhat

resemble our BatchAll algorithm. However, their scheme requires packet lengths to be independent and applies to specific switch architectures, whereas our scheme addresses the full routing-scheduling question for general graphs and does not impose specific statistical requirements (e.g., independence) on the arrival process or the packet length. In [27], a queuing analysis of several optical transport network architectures is conducted. It is shown that, under certain conditions imposed on the arrival process, some of the schemes can achieve the maximum network rate. This paper does not address the full routing-scheduling issue either. Another difference with the above two papers is that we provide an algorithm, BatchLim, guaranteeing the completion time of a job at the time of its arrival.

Many papers have discussed the issue of path dispersion and attempted to achieve good throughput with limited dispersion. A survey of some results in this field is given in [28]. In [29, 30], heuristic methods of controlling multipath routing and some quantitative measures are presented. As far as we know, our work proposes the first formal treatment allowing the approximation of a flow using a limited number of paths at any desired level of accuracy.

III. NETWORK MODEL

A. Notation and assumptions

We consider a *general* network topology, represented by a graph $G(V, E)$, where V is the set of nodes and E is the set of links connecting the nodes. The graph G can be directed or undirected. The capacity of each link $e \in E$ is $C(e)$. Any pair of nodes may request a connection at any time. A connection request, also referred to as job, contains the tuple (s, d, f) , where $s \in V$ is the source node, $d \in V - \{s\}$ is the destination node, and f is the file size.

Accordingly, an advance reservation algorithm computes a starting time at which the connection can be initiated, a set of paths used for the connection, and an amount of bandwidth allocated to each path. Our model supports path dispersion, i.e., multiple paths in parallel can be used to route data traffic.

In subsequent sections, we will make frequent use of multicommodity functions. The multicommodity flow problem is a linear planning problem returning true or false based upon the feasibility of fully transmitting a list of jobs between given source-destination pairs concurrently during a given time window, such that the total flow through each link does not exceed its capacity. It is solved by a function $\text{multicomm}(G, L, T)$, where L is a list of jobs, each containing a source, a destination, and a file size, and T is the time window.

The maximum concurrent flow is calculated by the function $\text{maxflow}(G, L)$. It returns T_{\min} , the minimum value of T such that $\text{multicomm}(G, L, T)$ returns true. Both the multicommodity and maximum concurrent flow problems are known to be computable in polynomial time [31, 32].

We assume that no queuing delays arise in relay nodes in the network. Therefore, the instantaneous departure rate of information at each relay node equals the arrival rate. It is a reasonable assumption for the networks under consideration

which reserves dedicated bandwidth to each flow. As such, the average transmission rate between all pairs in the network for any time interval T must be a feasible multicommodity flow.

B. Performance metrics

The main performance metric of interest in this paper is the *saturation throughput*, defined as follows. Suppose requests generate an average demand over time of $\lambda_{ij} = \alpha_{ij}\lambda$ bits/second from each node i to node j , where α_{ij} is a fixed parameter and $j \neq i$. Define the *delay* of a request as the amount of time elapsing from the point where the request for a connection set-up arrives till the point where the corresponding connection completes. Then, the saturation throughput, denoted λ^* , is the maximum value of λ (or the supremum, if the maximum does not exist) for which the delay experienced by any request is finite almost surely.

Determining analytically the value of λ^* is a difficult task in general as the answer depends on the graph topology and the statistics of the arrival process. Yet, even without knowing λ^* , one can devise algorithms that achieve optimal or competitive saturation throughput, as described later.

A simple upper bound on λ^* can be obtained as follows. Consider a set of feasible multicommodity flows g_{ij} assigned to each pair of nodes i and j , such that $g_{ij} = \alpha_{ij}g$. The maximization of g is a linear planning problem known as the *Maximum Concurrent Flow Problem (MCFP)* [31]. The maximum value of g in MCFP, denoted g^* , provides an upper bound on the saturation throughput λ^* . This result holds true because for every achievable value of λ there must exist an equally feasible value g .

IV. GREEDY ALGORITHMS

In this section, we present greedy algorithms and show that their saturation throughput can be arbitrarily far from optimal as the network size grows.

A. The Greedy algorithm

A seemingly natural way to implement advance reservation is to follow a greedy procedure, where, at the time a request is placed, the request is allocated a path (or set of paths) guaranteeing the earliest possible completion time. We refer to this approach as Greedy and explain it next.

The Greedy algorithm divides the time axis into slots delineated by events. Each event corresponds to a set-up or tear-down instance of a connection. Therefore, during each time slot the state of all links in the network remains unchanged. In general, the time axis will consist of n time slots, where $n \geq 1$ is a variable and slot i corresponds to time interval $[t_i, t_{i+1}]$. Note that $t_1 = t$ (the time at which the current request is placed) and $t_{n+1} = \infty$.

Let $W_i = \{b_i(1), b_i(2), \dots, b_i(|E|)\}$ be the vector of reserved bandwidth on all links at time slot i where $i = 1, \dots, n$, and $b_i(e)$ denote the reserved bandwidth on link e during slot i , with $e = 1, \dots, |E|$.

For each slot $i \in L$, construct a graph G_i , where the capacity of link $e \in E$ is $C_i(e) = C(e) - b_i(e)$, i.e., $C_i(e)$ represents the available bandwidth on link e during slot i .

In order to guarantee the earliest completion time, Greedy repeatedly performs a maximum flow allocation between nodes s and d , for as many time slots as needed until the entire file is transferred. This approach ensures that, in each time slot, the maximum possible number of bits is transmitted and, hence, the earliest completion time feasible at the arrival instance of the request is achieved. The Greedy algorithm can thus be concisely described with the following pseudo-code:

- 1) Initialization
 - Set initial time slot: $i=1$.
 - Set initial size of remaining file: $r = f$.
- 2) If $\text{maxflow}(G_i, s, d, r) \leq t_{i+1} - t_i$ (i.e., all the remaining file can be transferred during the current time slot) then,

Exit step:

 - Update W_i by subtracting the used bandwidth from every link in the new flow and, if the file transfer of the new flow completes before t_{i+1} , create a new event t_{i+1} corresponding to the end of the connection.
 - Exit procedure.
- 3) Else,

Non-exit step:

 - Update W_i by subtracting the used bandwidth from every link in the new flow.
 - $r = r - r \times (t_{i+1} - t_i) / \text{maxflow}(G_i, s, d, r)$ (update size of remaining file).
 - $i = i + 1$ (advance to the next time slot).
 - Go back to step 2.

The complexity of this algorithm depends on the number of slots examined before completely serving an arriving job. Noting that at each one of the examined slots a run of maxflow is carried out, this algorithm is expensive to implement in practice since the number of slots can be very large, especially at high load.

B. The Greedy_shortest algorithm

The Greedy_shortest algorithm is a variation of Greedy, where only shortest paths (in terms of number of hops) between the source s and destination d are utilized to route data. To implement Greedy_shortest, we employ exactly the same procedure as in Greedy, except that we prune all the links not belonging to one of the shortest paths using breadth first search. Note that for a given source s and destination d , the pruned links are the same for all the graphs G_i .

C. Inefficiency Results

We next constructively show that there exist certain arrival patterns, for which the saturation throughput achieved by both Greedy and Greedy_shortest is significantly lower than optimal. Specifically, we present cases where this measure for Greedy and Greedy_shortest is $\Omega(|V|)$ times smaller than the optimal value.

Theorem 1: For any given vertex set with cardinality $|V|$, there exists a graph $G(V, E)$ such that the saturation throughput of Greedy is $|V|/2$ times smaller than the optimal saturation throughput.

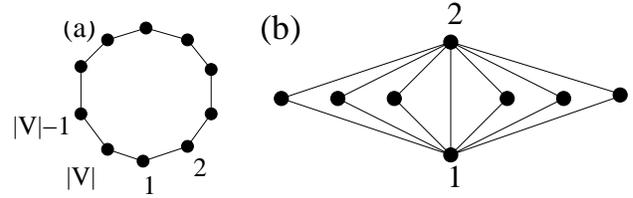


Fig. 1. Examples of graph topologies for which greedy algorithms perform far from optimal.

Proof: Consider the ring network shown in Fig. 1(a). Suppose that every link is an undirected 1 Gb/s link and requests arrive in the following order (we assume negligible delay between arrivals of requests): 1Gb request from node 1 to node 2, 1Gb request from node 2 to node 3, 1Gb request from node 3 to node 4, ..., 1Gb request from node $|V|$ to node 1.

The Greedy algorithm will allocate the maximum flow to each request, meaning that it will split the data flow to two paths, half flowing directly along a one hop path, and half flowing through the alternate path along the entire ring. Thus, the first job requires $1/2$ second for completion. Since the first job occupies all network resources, the second job has to wait until the completion of the first job before being served. Likewise, each new request will have to wait for the previous one to end, resulting in a completion time of $|V|/2$ seconds.

On the other hand, the optimal time is just one second using the direct link between each pair of nodes. Assuming the above pattern of requests repeats periodically, an optimal algorithm can support up to one request per second between each pair of neighboring nodes before reaching saturation, while Greedy can support at most $2/|V|$ request per second between each pair of nodes, hence proving the theorem. A similar proof can be shown for directed graphs. ■

We next show that restricting routing to shortest paths does not solve the inefficiency problem.

Theorem 2: For any given vertex set cardinality $|V|$, there exists a graph $G(V, E)$ such that the saturation throughput of Greedy_shortest (or any other algorithm using only shortest path routing) is $|V| - 1$ times smaller than the optimal saturation throughput.

Proof: Consider the network depicted in Fig. 1(b), where all requests are from node 1 to node 2, and only the direct path is used by the algorithm. In this scenario, an optimal algorithm would use all $|V| - 1$ paths between nodes 1 and 2. Hence, the optimal algorithm can achieve a saturation throughput $|V| - 1$ times higher than Greedy_shortest. ■

Both theorems imply that:

Corollary 3: There exist networks for which the delay of each job is bounded using the optimal algorithm, but is unbounded using either the Greedy or Greedy_shortest algorithm, even if the capacity of each link is multiplied by $|V|/2 - 1$.

V. THROUGHPUT-OPTIMAL ALGORITHMS

In this section, we present on-line, polynomial-time algorithms based on the idea of *deferred bandwidth allocation*.

Thus, arriving requests are not assigned bandwidth as long as transmissions go on in the network. Once existing transmissions complete, pending requests are served altogether. The proposed algorithms guarantee that the maximum delay experienced by any request in a network with augmented resources lies within a finite multiplicative factor of the value achieved with an optimal off-line algorithm (i.e., with full knowledge of the future and unlimited computational resources) in the original network. Hence, the proposed algorithms reach the maximum throughput achievable in the sense that for any network and any optimal algorithm that can support a list of jobs with a bounded delay, the proposed algorithms also guarantee a bounded delay at the price of augmenting the capacity of each link by a multiplicative factor of $(1 + \epsilon)$, for any arbitrarily small constant $\epsilon > 0$. This result stands in contrast to the results of Corollary 3, which indicates that, even with high augmentation, the greedy approaches are not competitive with the optimal algorithm.

A. The BatchAll Algorithm

Our first algorithm, called BatchAll, cumulatively performs bandwidth allocation for a group of requests. The group of all pending requests forms a *batch* that is served once existing transmissions complete. Batches of jobs are not predetermined since they are formed by on-line aggregation of random arrivals. We denote by L the list of jobs assigned to the next batch, by t_c the end time of the currently running batch of jobs, and by clk the value of the clock (initially set to 0). BatchAll can then concisely be described as follows:

- 1) Set t_c to 0.
- 2) When a request $l = \{s, d, f\}$ arrives at $clk = t$, give an immediate connection starting time and a connection ending time of $t_c = t + \text{maxflow}(G, l)$.
- 3) Set $L \leftarrow \text{null}$
- 4) While $clk < t_c$,
 - If a request $l = \{s, d, f\}$ arrives when $clk = t$:
 - Set $L \leftarrow L \cup l$ (i.e. add the job to the waiting batch)
 - Mark t_c as its connection starting time
- 5) At time $clk = t_c$,
 - If L is empty (i.e., there is no pending request) go back to step 2.
 - Else assign a connection ending time $t_c = t_c + \text{maxflow}(G, L)$ to all requests in the batch L and go back to step 3.

As an example, one can verify that BatchAll achieves optimal throughput for the scenarios used in the proofs of Theorems 1 and 2, which showed the inefficiency of greedy algorithms.

Next, we compare the delay performance of the BatchAll algorithm in an *augmented* network to that of the optimal off-line algorithm in the original network. The augmented network is similar to the original network, other than each link e has capacity $(1 + \epsilon)C(e)$ instead of $C(e)$, where $\epsilon > 0$. Alternatively, one may compare the performance of BatchAll to that of the optimal off-line algorithm in a lower capacity

network, allowing the maximum rate of only $(1 + \epsilon)^{-1}C(e)$ for each link e .

Theorem 4: Suppose that we augment the resources of a network such that every edge has $(1 + \epsilon)$ times its original capacity, for any $\epsilon > 0$. Then, for all requests arriving up to any time t^* in such a network, the maximum delay using BatchAll is no more than $2/\epsilon$ times the maximum delay for all requests arriving up to t^* using the optimal algorithm in the original network.

Proof: Consider the augmented resource network. Take the maximum length batch, say i , that accommodates requests arriving before time t^* and mark its length by T . Since the batch before this one was of length at most T , and all requests in batch i were received during the execution of previous batch, denoted by $i - 1$, then the total waiting time of each of these requests was at most $2T$. Thus, $2T$ is an upper bound on the maximum delay in the augmented network achieved using BatchAll up to time t^* . Since BatchAll uses a maximum concurrent flow procedure to allocate bandwidth in each batch, the total time for handling all requests received during the execution of batch $i - 1$ must have been at least T in the augmented network, or $(1 + \epsilon)T$ in the original network. Since all of these requests arrived during a time interval of length at most T and the completion time for all of them together is at least $T + \epsilon T$, one of them must have waited at least ϵT until job completion under the optimal algorithm. Thus, the maximum delay in the original network using an optimal offline algorithm up to any time t^* is at least ϵT . Therefore, the ratio between the maximum waiting time is at most $2/\epsilon$. ■

The throughput-optimality of BatchAll immediately follows from the previous theorem:

Corollary 5: For any given network, the saturation throughput of BatchAll is identical to that of the offline optimal algorithm because, for any arbitrarily small $\epsilon > 0$, the delay of any request is guaranteed to be, at worse, a finite multiplicative factor larger than the maximum delay of the optimal algorithm in the reduced resources network.

B. The BatchAll+ algorithm

We propose next a simple modification of BatchAll to improve its delay performance. The main idea is as follows. Whenever a job request arrives, we first check if we can allocate it to the current batch without affecting the completion time of the batch. If yes, we immediately allocate bandwidth to the job. If no, we proceed as in BatchAll, and add the current job to the list of jobs waiting for the next batch. To implement this idea, we keep track of the available bandwidth on each link of the graph during the current batch. We denote the resulting graph by G' . When a request arrives, we run maxflow over G' to examine whether the network can immediately accommodate the new job without disturbing the completion time of currently running jobs. Thus, BatchAll+ modifies step 4 of BatchAll as follows:

- While $clk < t_c$,
 - If request $l = \{s, d, f\}$ arrives when $clk = t$:
 - * If $\text{maxflow}(G', l) \leq t_c - t'$ then,

- Add the job to the running batch such that it completes simultaneously with the rest of running jobs.
- Deduct all bandwidth assigned to l from G' .
- * Otherwise,
 - Set $L \leftarrow L \cup l$ (i.e. add it to the waiting batch)
 - Mark t_c as its connection starting time

Note that G' must also be updated in step 2 and step 5 in the BatchAll pseudo-code, after maxflow is run on G , to reflect the available bandwidth on each link. Since computing maxflow is feasible in polynomial time, the computation complexity per request remains polynomial after this modification. It is not difficult to realize that Theorem 4 continues to hold after the modification. Therefore, BatchAll+ also achieves optimal saturation throughput.

C. The BatchLim Algorithm

The BatchAll and BatchAll+ algorithms return the starting time of a connection at the time of the request but the completion time is computed only when the connection actually starts. We next present another algorithm, called BatchLim, with a slightly larger delay ratio guarantee but which always returns the completion time at the arrival time of the request. Note that since the delay ratio guarantee remains finite, BatchLim is also throughput-optimal.

For each arriving request, the algorithm maintains a list of times, t_i , $i = 1, \dots, n$ where $t_1 = t$ (the time at which the current request is placed) and each time interval $[t_i, t_{i+1}]$ correspond to a batch of previously allocated jobs, which are either currently running (for the case $i = 1$) or scheduled to start in the future (for $i > 1$). When a new request between a source s and a destination d arrives at time t , two cases may happen: (i) if all previous jobs have completed, a new interval $[t, t + M]$ where $M = \text{maxflow}(G, s, d, f)$ is created and the job is assigned to it; (ii) otherwise, attempts are made to add the job to one of the scheduled time intervals $[t_i, t_{i+1}]$, where $2 \leq i \leq n - 1$, by using a multicommodity flow calculation. If the attempts fail, a new interval $[t_n, \max(2t_n - t, t_n + M)]$ is appended to the time list and the job is assigned to this interval.

We next provide a detailed description of how the algorithm handles a new request arriving at time t . We use the tuple $l = \{s, d, f\}$ to denote the job request and the list L_i to denote the set of jobs already assigned to interval $[t_i, t_{i+1}]$:

- 1) If no job is currently running, then:
 - Set $M = \text{maxflow}(G, l)$.
 - Set $t_2 = t + M$.
 - Assign job l to interval $[t_1, t_2]$ and exit.
- 2) Else check if job l can be assigned to an existing interval:
 - Set $i = 2$.
 - While $i \leq n - 1$:
 - If l can be fitted into interval $[t_i, t_{i+1}]$ (i.e. $\text{multicomm}(G, L_i \cup l, t_{i+1} - t_i) = \text{true}$) then assign job l to $[t_i, t_{i+1}]$ and exit.
 - Else set $i = i + 1$.
- 3) Else create new interval:

- Set $M = \text{maxflow}(G, l)$.
- If $t_n - t < M$, then $t_{n+1} = t_n + M$.
- Else $t_{n+1} = 2t_n - t$.
- Assign job l to interval $[t_n, t_{n+1}]$ and exit.

Fig. 2 illustrates runs of the BatchAll and BatchLim algorithms for the same set of requests.

The following lemma shows that the delay of any request arriving at some time t and for which a new interval $[t_n, t_{n+1}]$ is created is at most twice the length of that interval.

Lemma 6: For any $n \geq 1$, $(t_{n+1} - t) \leq 2(t_{n+1} - t_n)$, where t is the time at which the interval $[t_n, t_{n+1}]$ is formed.

Proof: Suppose a request $l = \{s, d, f\}$ arrives at time t and is assigned a new interval $[t_n, t_{n+1}]$. Then, by definition $t_{n+1} - t_n = \max(M, t_n - t)$ where $M = \text{maxflow}(G, l)$. Thus $t_n - t \leq t_{n+1} - t_n$. Adding $t_{n+1} - t_n$ to each side of the previous inequality, we derive the lemma. ■

The following lemma states that the delay experienced by any request is at most twice the length of the interval to which it is allocated.

Lemma 7: Suppose a request arrives at time t and is assigned to interval $[t_i, t_{i+1}]$, where $i = 1, 2, \dots, n$. Then $(t_{i+1} - t) \leq 2(t_{i+1} - t_i)$.

Proof: Suppose first that the request, call it l , is assigned to a new interval $[t_n, t_{n+1}]$, then the proof follows directly from Lemma 6. Next suppose that request l is assigned to an existing interval $[t_i, t_{i+1}]$, where $i < n$. Then, it must be that request l arrived after the request, call it l' , for which this interval was formed. Thus, the delay of request l is smaller than the delay of request l' and the proof follows again from Lemma 6. ■

The next theorem provides a guarantee on the maximum delay ratio of BatchLim in an $(1 + \epsilon)$ -augmented network versus the optimal delay in the original network. The theorem implies that BatchLim is also throughput-optimal.

Theorem 8: Suppose that we augment the resources of a network such that every edge has capacity $(1 + \epsilon)$ times the original capacity, for any $\epsilon > 0$. Then for requests arriving up to any time t^* in such a network, the maximum delay for a request using BatchLim is at most $4/\epsilon$ times the maximum delay for requests arriving up to t^* using the optimal algorithm in the original network.

Proof: Consider the maximum length time interval formed up to t^* by BatchLim in the $(1 + \epsilon)$ augmented network. Suppose that this interval was created at time t and denote it $[t_n, t_{n+1}]$. When creating this interval there are two possibilities:

- 1) $t_{n+1} = t_n + M$, where M is the minimum time for the new job completion. In this case, the minimum delay using any algorithm in the original network would have been $(1 + \epsilon)M$. By Lemma 7, the total delay for any job using BatchLim is at most $2M$, and therefore the ratio of waiting times is $2/(1 + \epsilon) < 4/\epsilon$.
- 2) $t_{n+1} = t_n + (t_n - t)$. Consider the previous interval $[t_{n-1}, t_n]$, which was formed at some time $t' < t$ when another request arrived that could not be fitted into any previous intervals. By the algorithm definition $t_n \geq t_{n-1} + (t_{n-1} - t')$. By maximality of $[t_n, t_{n+1}]$ it follows that $t_{n+1} - t_n = t_n - t \geq t_n - t_{n-1}$, and thus,

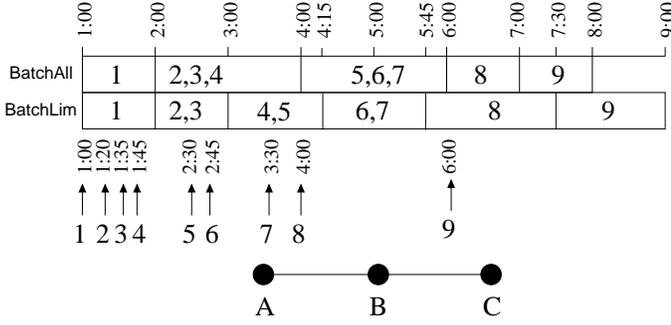


Fig. 2. Illustration of the batching process by the algorithms `BatchAll` and `BatchLim` over a simple network topology consisting of three nodes A , B and C connected via links of identical capacity. Connection requests are depicted by arrows and their arrival time. The odd numbered jobs request transmission between nodes A and B and the even numbered jobs request for transmission between B and C . Each file size corresponds to one hour of transmission at full link capacity. In `BatchAll`, a new batch is created for all requests arriving during the running of a previous batch. In `BatchLim`, for each new request an attempt is made to add it to one of the existing windows, and if it fails, a new window is appended at the end.

$t \leq t_{n-1}$. Note that not all the requests arriving during the interval $[t', t]$ (including the request arriving at time t , which initiated the creation of interval $[t_n, t_{n+1}]$) could have been completed during the interval $[t_{n-1}, t_n]$ or otherwise the new interval $[t_n, t_{n+1}]$ would not have been created at time t . Thus, the requests arriving during the interval $[t', t]$ require a time of at least $t_n - t_{n-1}$ for all to complete according to the multicommodity flow calculation. Thus, the time it would take to complete all jobs that arrived in $[t', t]$ under any possible algorithm is at least $t_n - t_{n-1}$ in the augmented network, and therefore, at least $(1 + \epsilon)(t_n - t_{n-1})$ in the original network. Since no job can start before time t' , the ending time for all these jobs in the original network must be at least $t' + (1 + \epsilon)(t_n - t_{n-1}) \geq t' + (t_{n-1} - t') + \epsilon(t_n - t_{n-1}) = t_{n-1} + \epsilon(t_n - t_{n-1})$. Since the requests have arrived between $[t', t]$, the last request has arrived at time $t \leq t_{n-1}$ at the latest. Thus, one of the requests in this batch must have waited for at least $t_{n-1} + \epsilon(t_n - t_{n-1}) - t_{n-1} = \epsilon(t_n - t_{n-1})$. In other words, the maximum delay using the optimal algorithm in the original network is at least $\epsilon(t_n - t_{n-1})$. On the other hand, according to Lemma 7, the maximum delay of any request in the augmented network using `BatchLim` is at most $2(t_{n+1} - t_n)$. Since by assumption $t_{n+1} - t_n = t_n - t$ and since $t_n - t \leq t_n - t' \leq 2(t_n - t_{n-1})$, where the last inequality follows from Lemma 6, then the maximum delay of any request in the augmented network using `BatchLim` is at most $4(t_n - t_{n-1})$. Therefore, the ratio of the maximum delays is at most $4/\epsilon$. ■

VI. BOUNDING PATH DISPERSION

The algorithms presented in the previous sections do not limit the *path dispersion*, that is, the number of paths simultaneously used by a connection. In practice, it is desirable to

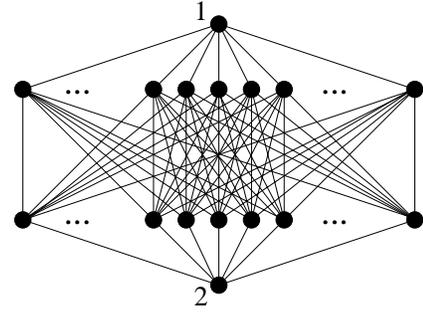


Fig. 3. A network demonstrating the optimality of Lemma 9. The edges incident to nodes 1 and 2 have infinite capacity, all other edges have capacity 1.

minimize the number of such paths due to the cost of setting up many paths and the need to split a connection into many low capacity channels. The following suggests a method of achieving this goal.

Lemma 9: In every flow of bandwidth F between two nodes on a directed graph $G(V, E)$ there exists a path of bandwidth at least $F/|E|$ between these nodes.

Proof: Remove all edges carrying flow with bandwidth strictly less than $F/|E|$ from the graph. The total flow carried by these edges is smaller than $F/|E| \times |E| = F$. By the Max-Flow–Min-Cut theorem, the maximum flow equals the minimum cut in the network, which is of capacity at least F . Therefore, since the flow bandwidth is F , there must remain at least one path between the nodes after the removal. All edges in this path have flow bandwidth of at least $F/|E|$. Therefore, the path bandwidth is at least $F/|E|$. ■

Lemma 9 is asymptotically tight in order of both $|E|$ and $|V|$ as can be seen for the network shown in Fig. 3. In this network each path amounts for only $O(|E|^{-1}) = O(|V|^{-2})$ of the total flow between nodes 1 and 2. Therefore, to obtain a flow that consists of a given fraction of the maximum flow at least $O(|E|) = O(|V|^2)$ of paths must be used. If multigraphs are allowed, the lemma is exactly tight for a graph with 2 nodes and $|E|$ equal capacity edges between them.

The following theorem establishes the maximum number of paths needed to achieve a throughput with a constant factor of that achieved by the original flow.

Theorem 10: For every flow of bandwidth F between two nodes on a directed graph $G(V, E)$ and for every integer $n \leq |E|$ there exists a set of at most n paths achieving a flow of bandwidth at least $nF/|E|$.

Proof: We prove by induction on n . For $n = 1$, there exists, by Lemma 9, a path of capacity at least $F/|E|$. Take the edge (i, j) with minimum capacity of the path and decrease the flow through all edges in the path by its weight w_{ij} . Now, the flow at at least one edge is decreased to zero, and this edge may be removed without affecting the residual flow. Thus, we are left with a graph with flow $F - w_{ij}$ and with at most $|E| - 1$ edges.

Now assume that by using n paths one can achieve a flow of $f_n \geq nF/|E|$. The residual flow is $F - f_n$, and the number of edges left is at most $|E| - n$. The flow of the widest left path, p_n , is, by Lemma 9, at least $p_n \geq (F - f_n)/(|E| - n)$.

Thus the total flow captured in the first $n + 1$ paths is

$$\begin{aligned} f_n + p_n &\geq f_n + \frac{F - f_n}{|E| - n} = \frac{(|E| - n - 1)f_n + F}{|E| - n} \\ &\geq \frac{(|E| - n - 1)n \frac{F}{|E|} + F}{|E| - n} \\ &= \frac{(|E| - n)(n + 1)F}{|E|(|E| - n)} = \frac{n + 1}{|E|} F. \end{aligned}$$

This completes the induction. \blacksquare

Theorem 10 provides both an algorithm for reducing the number of paths and a bound on the throughput loss. To approximate the flow using a limited number of paths, consider only links carrying flow with bandwidth greater than $F/|E|$ and find a path p among those (if multiple paths are available, select the widest). Suppose path p carries a flow f , then remove the bandwidth used by flow f from all the links on path p and iteratively repeat the same procedure until a satisfactory approximation of the flow is obtained. Note that the maximum number of paths needed to achieve an approximation of the flow to within a constant factor is linear in $|E|$ and the number of paths needed to obtain the full flow is exactly $|E|$, while the number of possible paths may be exponential in $|E|$. Again, this theorem is asymptotically tight, and exactly tight for multigraphs.

We illustrate the bounding path dispersion procedure with a simple example. Figure 4 shows the maximum flow from node 10 to node 3 on an 11-node topology. Links are assumed to be full duplex and link capacities are represented by surrounded numbers next to each link. The maximum flow calculated using a `maxflow` formulation from node 10 to node 3 is $F = 20$ Gb/s. Since the total number of edges is $|E| = 28$ in this graph, we have $F/|E| = 0.71$ Gb/s.¹ Therefore, at the first step, the algorithm only considers links 10-1, 1-2, 2-3, 10-9, 9-8, 8-7, 7-3 and 8-2 carrying a flow greater than 0.71 Gb/s and ignore all other links. Among the considered links, there exists three different paths between node 10 and 3, marked as $P1$, $P2$ and $P3$, and carrying flows of bandwidth 10, 5 and 5 Gb/s respectively. The algorithm selects the widest path $P1$ first. Subtracting the flow of path $P1$ from F , the remaining flow is 10 Gb/s. The links 10-1 and 1-2 are now saturated. Thus we are left with a residual flow of $F = 20 - 10 = 10$ Gb/s and a remaining $|E| = 26$ non saturated edges. We repeat the procedure and get $F/|E| = 0.38$ Gb/s. This time, only links 2-3, 10-9, 9-8, 8-7, 7-3 and 8-2 are considered. The algorithm can select either $P2$ or $P3$. Thus, with a bound of at most two paths per connection we can carry 75% of the full flow from node 10 to node 3 using paths $P1$ and $P2$ (or $P1$ and $P3$). If the bound is relaxed to three paths, then we can route the entire flow between the given nodes. Our simulation results in section VII-B reveals that with this topology, bounding flows with at most 3 paths per connection does indeed closely approximate full flow on average.

Using the above flow approximation in conjunction with the competitive algorithms, we can devise two new algorithms `BatchAllDisp` and `BatchLimDisp`, based upon `BatchAll`

¹In fact, the entire flow can be routed through 8 directed edges. Thus, one can also use $|E| = 8$ and obtain $F/|E| = 2.5$ Gb/s.

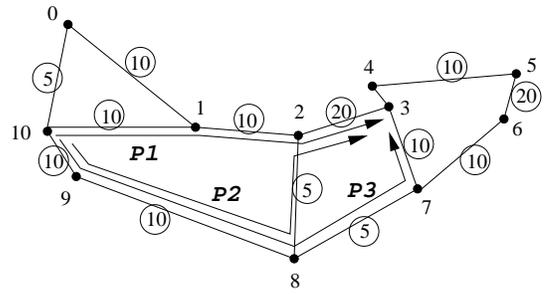


Fig. 4. Illustration of the path dispersion bounding procedure. Links are assumed to be full duplex and link bandwidths are displayed by surrounded numbers next to each link. The three paths connecting node 10 to node 3 are marked by $P1$, $P2$ and $P3$ with flow bandwidth of 10, 5 and 5 Gb/s respectively.

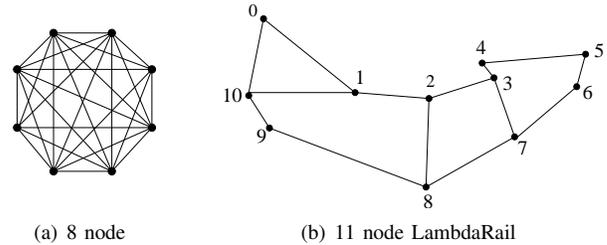


Fig. 5. Simulation topologies.

and `BatchLim` respectively. These algorithms perform similarly to `BatchAll` and `BatchLim`, in terms of the batching process. However, after filling each batch, the algorithms will limit the dispersion of each flow, and approximate the flow. To achieve a partial flow, each time we select the widest path from remaining edges and reserve the total path bandwidth. We repeat this procedure until either the desired number of paths is reached or the entire flow is routed.

VII. SIMULATIONS

In this section, we present simulation results illustrating the performance of the algorithms described in this paper. While the emphasis of the previous sections was on the throughput optimality (or competitiveness) of the proposed algorithms, here we are also interested in evaluating their delay performance. The main points of interest are as follows: (i) how do the competitive algorithms fare with respect to each other and the capacity bound of section III-B? (ii) what value of path dispersion is needed to ensure good performance?

A. Simulation Set-Up

We have developed our own simulator in C++. The simulator uses the COIN-OR Linear Program Solver (CLP) library [33] to solve multicommodity optimization problems and allows evaluating our algorithms under various topological settings and traffic conditions. The main simulation parameters are as follows:

- *Topology*: our simulator supports arbitrary topologies. In this paper, we consider the two topologies depicted in Figure 5. One is a fully connected graph (clique) of eight

nodes and the other is an 11-node topology, similar to the National LambdaRail testbed [34]. Each link on these graphs is full-duplex and assumed to have a capacity of 20 Gb/s.

- *Arrival process*: we assume that the aggregated arrival of requests to the network forms a Poisson process (this can easily be changed, if desired). The mean rate of arrivals is adjustable. Our delay measurements are carried out at different request arrival rates in units of requests per hour.
- *File size distribution*: We consider two models for the file size distribution:
 - 1) Pareto:

$$F(x) = 1 - \left(\frac{x_m}{x - \gamma} \right)^\beta, \text{ where } x \geq x_m + \gamma.$$

In the simulations, we consider two cases: i) $\beta = 2.5$, $x_m = 1.48$ TB (terabyte) and $\gamma = 6.25 * 10^{-3}$ TB; and ii) $\beta = 2$, $x_m = 1.23$ TB and $\gamma = 6.25 * 10^{-3}$ TB. In both cases, the average file size is 2.475 TB.

- 2) Exponential:

$$F(x) = \exp(-\delta x), \text{ where } x \geq 0.$$

In the simulations, $1/\delta = 2.475$ TB.

- *Source and Destination*: for each request, the source and destination are selected uniformly at random, except that they must be different nodes.

Each simulation point represents an average taken over at least 10^6 measurements.

Performance Metrics. Our simulations report results for two metrics related to the delay experienced by requests. The first is the *average delay*. The second is the probability that the delay of a request exceeds a certain threshold. We refer to this metric as *tail probability* [35].

Note that the average delay is defined only if $\beta > 2$. For the special case of a graph consisting of a single link, the performance of BatchAll is identical to that of a gated $M/G/1$ queue [36]. Denote by λ the mean arrival rate of requests, by $\mathbb{E}[S]$ and $\mathbb{E}[S^2]$ the first and second moments of the service time (the service time is the file size divided by the link capacity), and by $\rho = \lambda\mathbb{E}[S]$ the link utilization. Then, the mean waiting time $\mathbb{E}[W]$ of a job from the point that it arrives till the start of its batch is [36, Eq. 16]:

$$\mathbb{E}[W] = \frac{\lambda\mathbb{E}[S^2]}{2(1 - \rho^2)}.$$

As in the standard $M/G/1$ queue, the mean waiting time is bounded only if the second moment of the service time distribution (or equivalently that of the file size distribution) is finite. The same result holds also in general networks. Since measurements [37–40] have shown that the distribution of the size of Internet files has infinite variance, we also provide simulation results for power-law parameter $\beta = 2$ using tail probability as the performance metric.

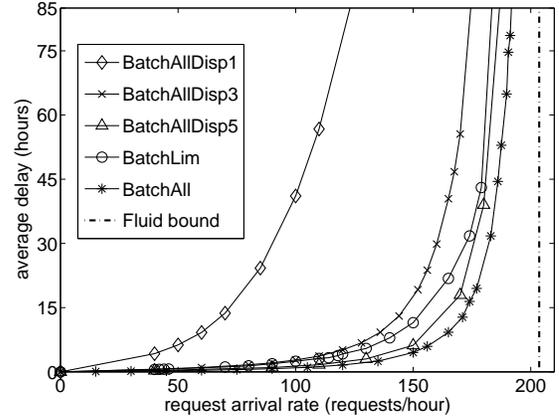


Fig. 6. Average delay performance evaluation of algorithms BatchAll, BatchLim and BatchAllDisp for the 8-node clique topology and exponentially distributed file size. Algorithm BatchAllDisp is plotted with path dispersion bounds of 1, 3 and 5 paths.

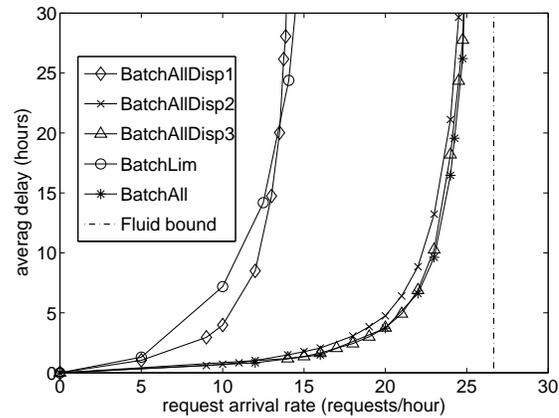


Fig. 7. Average delay performance evaluation of algorithms BatchAll, BatchLim and BatchAllDisp for the 11-node topology and Pareto distributed file size with $\beta = 2.5$. Algorithm BatchAllDisp is plotted with path dispersion bounds of 1, 2 and 3 paths per connection.

B. Average Delay Performance

In this section, we use average delay as the performance metric. When considering the Pareto distribution, we assume $\beta = 2.5$. We first compare the performance of BatchAll, BatchAllDisp and BatchLim algorithms in Figures 6 and 7. Both figures show BatchAllDisp with multiple bounds on path dispersion. The capacity bounds are also depicted in dashed lines as a performance benchmark for each scenario. Figure 6 corresponds to the 8-node clique topology with exponentially distributed file sizes and Figure 7 corresponds to the 11-node topology of Figure 5(b) with Pareto distributed file sizes.

According to Figure 6, BatchAll approaches the capacity bound at a reasonably low average delay value. Interestingly, a path dispersion of at most five per connection (corresponding to $\alpha = 0.089$) is sufficient for BatchAllDisp to closely approximate BatchAll. It is worth mentioning that, in this topology, there exist 1957 possible paths between any two nodes. Thus, with five paths, BatchAllDisp uses only 0.25% of the total paths possible. The figure also demonstrates the

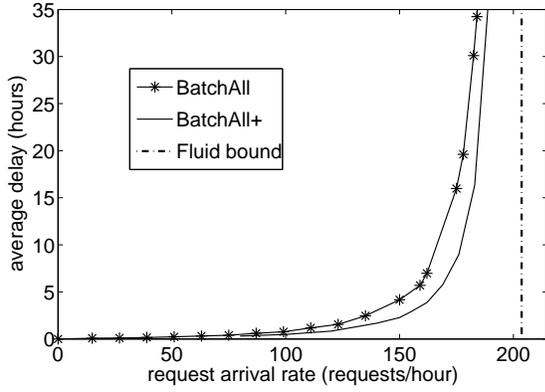


Fig. 8. Average delay performance comparison of algorithms BatchAll and BatchAll+ for the 8-node clique and Pareto file size distribution with $\beta = 2.5$.

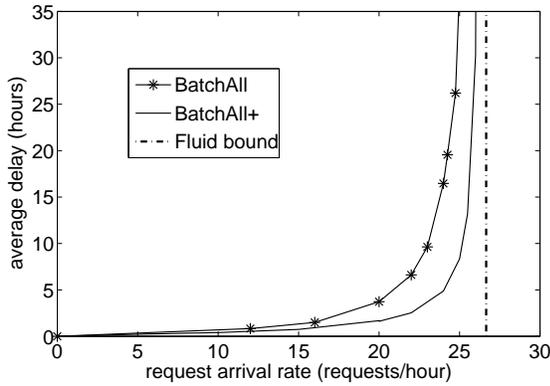


Fig. 9. Average delay performance comparison of algorithms BatchAll and BatchAll+ for the 11-node topology and Pareto file size distribution with $\beta = 2.5$.

importance of multi-path routing: the performance achieved using a single path per connections is far worse. Figure 6 also shows that BatchLim is less efficient than BatchAll in terms of average delay, especially at low load. This result is somewhat expected given that BatchLim uses a less efficient batching process and its delay ratio guarantee is looser.

Figure 7 depicts the performance of the various algorithms and the capacity bound for the 11-node topology of Figure 5(b). In this case, we observe that BatchAllDisp with as few as 3 paths per connection (or $\alpha = 0.107$) approximates BatchAll very closely. Since this network is sparser than the previous one, it is reasonable to obtain good performance with a smaller path dispersion. In this scenario with Pareto file size distribution, BatchLim performs significantly worse than BatchAll for the load range illustrated. We suspect that the performance of BatchLim is more sensitive to the file size distribution than BatchAll, because the batches formed by BatchLim are limited in time duration compared to those of BatchAll.

Figures 8 and 9 compare the performance of BatchAll and BatchAll+ for the 8-node and 11-node topologies BatchAll+, with Pareto file size distribution. The figures show the superiority of BatchAll+ over BatchAll in both cases, which is due to optimization in Step 4 of the algorithm.

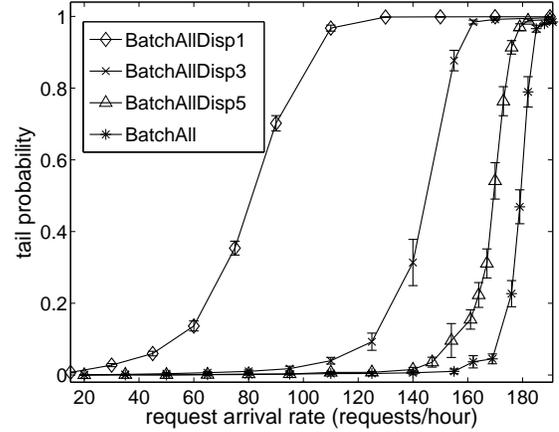


Fig. 10. Tail probability performance of algorithms BatchAll and BatchAllDisp for the 8-node topology, Pareto distributed file size with $\beta = 2$, and delay threshold of 48 hours. Algorithm BatchAllDisp is plotted with path dispersion bounds of 1, 3, and 5 paths per connection.

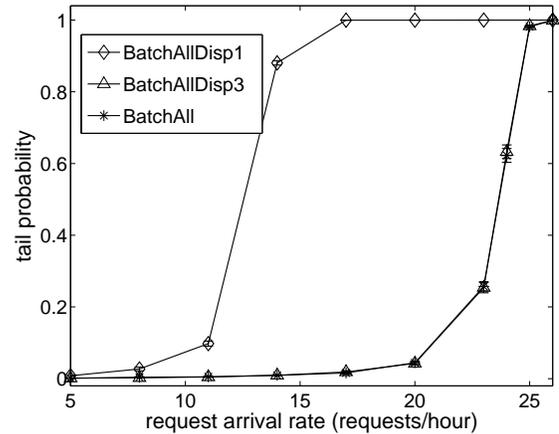


Fig. 11. Tail probability performance of algorithms BatchAll and BatchAllDisp for the 11-node topology, Pareto distributed file size with $\beta = 2$, and delay threshold of 24 hours. Algorithm BatchAllDisp is plotted with path dispersion bounds of 1 and 3 paths per connection.

C. Tail Probability Performance

In this section, we use the tail probability of the delay as the performance metric. All our simulations assume Pareto distributed file sizes with $\beta = 2$. Due to the high variability of the delay experienced by different requests, each simulation point in the following graphs represents an average taken over seven iterations, where each iteration is run for $4 \cdot 10^6$ requests. The graphs also plot 95% confidence intervals.

Figure 10 shows the probability that the delay of a request exceeds 48 hours, in the 8-node topology, as a function of the arrival rate of requests. We observe that the performance of BatchAllDisp with 5 paths per connection is close to that of BatchAll. Next, Figure 11 shows the probability that the delay of a request exceeds 24 hours, in the 11-node topology, as a function of the the arrival rate of requests. In that case the performance of BatchAllDisp with 3 paths per connection is virtually identical to that of BatchAll. Overall, the results are qualitatively similar to those obtained in the previous section.

VIII. CONCLUSION

In this paper, we considered the problem of devising throughput-optimal and throughput-competitive algorithms for networking architectures supporting advance reservation, a problem of particular relevance to modern grid and cloud computing applications. After showing the limitations of greedy approaches, we proposed two new on-line algorithms for advance reservation, called BatchAll and BatchLim, that are guaranteed to achieve optimal throughput performance. Specifically, we proved that both algorithms bounds the ratio of the maximum delay of any request in $1 + \epsilon$ bandwidth augmented networks to the maximum delay of any request in the original network using the optimal off-line algorithm. While BatchLim has a slightly looser delay ratio guarantee than that of BatchAll (i.e., $4/\epsilon$ instead of $2/\epsilon$) and, based on our simulations, inferior average delay performance, it has the distinct advantage of returning the completion time of a connection immediately as a request is placed.

We observed that path dispersion is essential to achieve full network utilization. However, splitting a transmission into too many different paths may render a flow-based approach inapplicable in many real-world environments. Thus, we presented a rigorous, theoretical approach to address the path dispersion problem and presented a method for approximating the maximum multicommodity flow using a limited number of paths. Specifically, while the number of paths between two nodes in a network scales exponentially with the number of edges, we showed that throughput competitiveness up to any desired ratio factor can be achieved with a number of paths scaling linearly with the total number of edges. In practice, our simulations indicate that three paths (in sparse graphs) to five paths (in dense graphs) should be sufficient.

We also provided a simple and practical mechanism for enhancing the delay performance of our main algorithm, BatchAll. Specifically, the enhanced version, called BatchAll+, attempts at adding any arriving request to the current batch of running requests pending that the completion time of the batch remains unchanged. Such an improvement can also be applied to the algorithm with bounded path dispersion, BatchAllDisp.

The algorithms proposed in this paper can be either run in a centralized fashion (a reasonable solution in small networks) or using link-state routing and distributed signaling mechanism, such as enhanced versions of GMPLS [41]. Distributed approximations of the multicommodity flow have also been discussed in the literature [42]. An important research area open for future work will be to carefully investigate these implementation issues.

We conclude by noting that advance reservation architectures relying on flow (circuit) switching may provide an additional benefit by saving energy [43]. Indeed, flow switching avoids the need of managing each packet individually [44]. Quantifying the potential energy gain of advance reservation architectures represents an important topic for future research.

ACKNOWLEDGMENT

This research was supported in part by the US Department of Energy under ECPI grant DE-FG02-04ER25605 and the

US National Science Foundation under CAREER grant ANI-0132802 and grant CCF-0729158.

REFERENCES

- [1] C. Guok, D. Robertson, E. Chianotakis, M. Thompson, W. Johnston, and B. Tierney, "A User Driven Dynamic Circuit Implementation," in *Proceedings of IEEE GLOBECOM*, December 2008, New Orleans, LA, USA.
- [2] N. Rao, W. Wing, S. Carter, and Q. Wu, "UltraScience Net: network testbed for large-scale science applications," *IEEE Communications Magazine*, vol. 43, pp. 12–17, 2005.
- [3] R. Cohen, N. Fazlollahi, and D. Starobinski, "Path switching and grading algorithms for advance channel reservation architectures," *IEEE/ACM Transactions on Networking (TON)*, vol. 17, no. 5, pp. 1684–1695, October 2009.
- [4] "On-demand secure circuits and advanced reservation systems," <http://www.es.net/oscars/index.html>.
- [5] W. Johnston, "ESnet: advanced networking for science," *SciDAC review*, 2007.
- [6] C. Guok, J. Lee, and K. Berket, "Improving the bulk data transfer experience," *International Journal of Internet Protocol Technology*, vol. 3, no. 1, pp. 46–53, 2008.
- [7] L. Burchard, "Networks with advance reservations: Applications, architecture, and performance," *Journal of Network and Systems Management*, vol. 13, pp. 429–449, 2005.
- [8] H. Lee, M. Veeraraghavan, H. Li, and E. Chong, "Lambda scheduling algorithm for file transfers on high-speed optical circuit," in *IEEE International Symposium on Cluster Computing and the Grid (CCGrid 2004)*, April 2004, Chicago, USA.
- [9] A. Banerjee, N. Singhal, J. Zhang, D. Ghosal, C.-N. Chuah, and B. Mukherjee, "A Time-Path Scheduling Problem (TPSP) for aggregating large data files from distributed databases using an optical burst-switched network," in *Proc. ICC*, 2004, Paris, France.
- [10] A. Banerjee, W. Feng, B. Mukherjee, and D. Ghosal, "Routing and scheduling large file transfers over lambda grids," in *Proc. of the 3rd International Workshop on Protocols for Fast Long-Distance Networks (PFLDnet'05)*, February 2005, Lyon, France.
- [11] R. Guerin and A. Orda, "Networks with advance reservations: the routing perspective," in *Proceedings of INFOCOM'00*, March 2000, Tel-Aviv, Israel.
- [12] V. Chan, G. Weinchenberg, and M. Médard, "Optical flow switching," in *Proc. 3rd International Conference on Broadband Communications, Networks and Systems*, October 2006, pp. 1–8, San Jose, USA.
- [13] B. Ganguly, "Implementation and modeling of a scheduled optical flow switching (ofs) network," Ph.D. dissertation, Massachusetts Institute of Technology, 2008.
- [14] J. Zheng and H. T. Mouftah, "Routing and wavelength assignment for advance reservation in wavelength-routed WDM optical networks," in *Proc. IEEE ICC*, vol. 5, 2002, pp. 2722–2726.
- [15] S. Figueira, N. Kaushik, S. Naiksatam, S. Chiappari, and N. Bhatnagar, "Advance reservation of lightpaths in optical-network based grids," in *Proc. ICST/IEEE Gridnets*, October 2004, San Jose, USA.
- [16] N. Kaushik and S. Figueira, "A dynamically adaptive hybrid algorithm for scheduling lightpaths in lambda-grids," in *Proc. IEEE/ACM CCGRID/GAN'05-Workshop on Grid and Advanced Networks*, May 2005, Cardiff, UK.
- [17] A. Goel, M. Henzinger, S. Plotkin, and E. Tardos, "Scheduling data transfers in a network and the set scheduling problem," in *Proc. of the 31st Annual ACM Symposium on the Theory of Computing*, 1999, pp. 189–199.
- [18] L. Lewin-Eytan, J. Naor, and A. Orda, "Admission control in networks with advance reservations," *Algorithmica*, vol. 40, pp. 293–304, 2004.
- [19] T. Erlebach, "Call admission control for advance reservation requests with alternatives," ETH, Zurich, Tech. Rep. 142, 2002.
- [20] M. Andrews, A. Fernandez, A. Goel, and L. Zhang, "Source routing and scheduling in packet networks," *Journal of the ACM (JACM)*, vol. 52, no. 4, pp. 582–601, July 2005.
- [21] Y. Azar, J. Naor and R. Rom, "Routing Strategies in Fast Networks," *IEEE Transactions on Computers*, vol. 45, no. 2, pp. 165–173, February 1996.
- [22] B. Awerbuch, Y. Azar, A. Fiat, S. Leonardi, and A. Rosen, "On-line competitive algorithms for call admission in optical networks," *Algorithmica*, vol. 31, no. 1, pp. 29–43, 2001.

- [23] S. A. Plotkin, "Competitive routing of virtual circuits in ATM networks," *IEEE J. on Selected Areas in Communication (JSAC)*, vol. 13, no. 6, pp. 1128–1136, 1995.
- [24] B. Awerbuch, D. Holmer, H. Rubens, and R. D. Kleinberg, "Provably competitive adaptive routing," in *INFOCOM*, 2005, pp. 631–641.
- [25] A. Goel, M. R. Henzinger, and S. A. Plotkin, "An online throughput-competitive algorithm for multicast routing and admission control." *J. Algorithms*, vol. 55, no. 1, pp. 1–20, 2005.
- [26] Y. Ganjali, A. Keshavarzian, and D. Shah, "Input queued switches: Cell switching vs. packet switching," in *Proceedings of INFOCOM'03*, March 2003, San-Francisco, CA, USA.
- [27] G. Weichenberg, V. Chan, and M. Médard, "On the capacity of optical networks: A framework for comparing different transport architectures," *IEEE J-SAC*, vol. 25, pp. 84–101, 2007.
- [28] E. Gustafsson and G. Karlsson, "A literature survey on traffic dispersion," *IEEE Network*, vol. 11, no. 2, pp. 28–36, March-April 1997.
- [29] J. Shen, J. Shi, and J. Crowcroft, "Proactive multi-path routing," *Lect. Notes in Comp. Sci.*, vol. 2511, pp. 145–156, 2002.
- [30] R. Chow, C. W. Lee, and J. Liu, "Traffic Dispersion Strategies for Multimedia Streaming," in *Proceedings of the 8th IEEE Workshop on Future Trends of Distributed Computing Systems*, 2001, p. 18.
- [31] F. Shahrokhi and D. W. Matula, "The maximum concurrent flow problem," *Journal of the ACM (JACM)*, vol. 37, no. 2, April 1990.
- [32] T. Leighton and S. Rao, "Multicommodity max-flow min-cut theorems and their use in designing approximation algorithms," *Journal of the ACM (JACM)*, vol. 46, no. 4, pp. 787 – 832, November 1999.
- [33] "Coin-Or project," <http://www.coin-or.org/>.
- [34] "National LambdaRail Inc." <http://www.nlr.net/>.
- [35] W. Whitt, "The impact of a heavy-tailed service-time distribution upon the M/G/1's waiting-time distribution," *Queueing Systems*, vol. 36, no. 1/3, pp. 71–87, 2000.
- [36] B. Avi-Itzhak and S. Halfin, "Response times in gated M/G/1 queues: the processor-sharing case," *Queueing Systems*, vol. 4, pp. 263–279, 1989.
- [37] V. Paxson and S. Floyd, "Wide-area traffic: the failure of poisson modelling," *IEEE/ACM Transactions on Networking*, vol. 3, no. 3, pp. 226–244, June 1995.
- [38] M. E. Crovella, M. S. Taqqu, and A. Bestavros, *Heavy-Tailed Probability Distributions in the World Wide Web*. Cambridge, MA, USA: Birkhauser Boston Inc., 1998.
- [39] M. F. Arlitt and C. L. Williamson, "Internet web servers: Workload characterization and performance implications," *IEEE/ACM Transactions on Networking*, vol. 5, no. 5, pp. 631–645, 1997.
- [40] A. Williams, M. Arlitt, C. Williamson, and K. Barker, *Web Workload Characterization: Ten Years Later*. Springer, August 2005.
- [41] "Dynamic resource allocation via GMPLS optical networks," <http://dragon.maxgigapop.net>.
- [42] B. Awerbuch, R. Khandekar, and S. Rao, "Distributed algorithms for multicommodity flow problems via approximate steepest descent framework," in *Proceedings of the ACM-SIAM symposium on Discrete Algorithms (SODA)*, 2007.
- [43] J. Baliga, R. Ayre, K. Hinton, W. V. Sorin and R. S. Tucker, "Energy consumption in optical IP networks," *Journal of Lightwave Technology*, vol. 27, no. 13, pp. 2391–2403, July 2009.
- [44] M. Zukerman, "Back to the future," *IEEE Communications Magazine*, vol. 47, no. 11, pp. 36–38, 2009.



Nilofar Fazlollahi received her B.S. in Electrical Engineering (2005) from the Sharif University of Technology, Tehran, Iran. Since 2005, she has been a Ph.D. student at Boston University under the supervision of Professor David Starobinski. In 2007, she received her M.S. while being a Ph.D. candidate. Her research interests are in the modeling and analysis of joint scheduling and routing schemes in ultra high-speed networks, and the application of identifying codes in efficient sensor network monitoring.



David Starobinski received his Ph.D. in Electrical Engineering (1999) from the Technion-Israel Institute of Technology. In 1999-2000 he was a visiting post-doctoral researcher in the EECS department at UC Berkeley. In 2007-2008, he was an invited Professor at EPFL (Switzerland). Since September 2000, he has been with Boston University, where he is now an Associate Professor.

Dr. Starobinski received a CAREER Award from the U.S. National Science Foundation (2002), an Early Career Principal Investigator (ECPI) Award from the U.S. Department of Energy (2004), the 2010 BU ECE Faculty Award for Teaching Excellence, and the Best Paper Award at WiOpt 2010. He is currently on the Editorial Board of the *IEEE/ACM Transactions on Networking*. His research interests are in the modeling and performance evaluation of high-speed, wireless, and sensor networks.



Reuven Cohen Received his B.Sc. in Physics and Computer Science and his Ph.D. in Physics from Bar-Ilan University, Ramat-Gan, Israel. He was a postdoctoral fellow in the Dept. of Mathematics and Computer Science at the Weizmann Institute, Rehovot, Israel and in the ECE department at Boston University and in the Dept. of Physics at MIT. Since October 2007 he has been at the Department of Mathematics at Bar-Ilan University in Israel, where he is an Assistant Professor. His research interests are random graphs, distributed algorithms and network

stability.