

## Scheduling data transfers in a network and the set scheduling problem <sup>☆</sup>

Ashish Goel,<sup>a,\*</sup> Monika R. Henzinger,<sup>b,2</sup> Serge Plotkin,<sup>c,3</sup>  
and Eva Tardos<sup>d,4</sup>

<sup>a</sup> *Departments of Management Science and Engineering and (by courtesy) Computer Science,  
Stanford University, Stanford, CA 94305, USA*

<sup>b</sup> *Google Inc., Mountain View, USA*

<sup>c</sup> *Department of Computer Science, Stanford University, Stanford, CA 94305, USA*

<sup>d</sup> *Department of Computer Science, Cornell University, Ithaca, NY 14853, USA*

Received 9 January 2001

---

### Abstract

In this paper we consider the *online ftp problem*. The goal is to service a sequence of file transfer requests given bandwidth constraints of the underlying communication network. The main result of the paper is a technique that leads to algorithms that optimize several natural metrics, such as max-stretch, total flow time, max flow time, and total completion time. In particular, we show how to achieve optimum total flow time and optimum max-stretch if we increase the capacity of the underlying network by a logarithmic factor. We show that the resource augmentation is necessary by proving polynomial lower bounds on the max-stretch and total flow time for the case where online and offline algorithms are using same-capacity edges. Moreover, we also give polylogarithmic lower bounds on the resource augmentation factor necessary in order to keep the total flow time and max-stretch within a constant factor of optimum.

© 2003 Elsevier Inc. All rights reserved.

---

<sup>☆</sup> A preliminary version of these results appeared as [Goel et al., in: 31st ACM Symposium on Theory of Computing, 1999, pp. 189–197].

\* Corresponding author.

*E-mail addresses:* ashishg@stanford.edu (A. Goel), monika@google.com (M.R. Henzinger), plotkin@cs.stanford.edu (S. Plotkin), eva@cs.cornell.edu (E. Tardos).

<sup>1</sup> Research supported by ARO Grant DAAG55-98-1-0170 and ASSERT award DAAG55-97-1-0221. This research was conducted when the author was at the Compaq Systems Research Center.

<sup>2</sup> The author was at the Compaq Systems Research Center, Palo Alto when this research was conducted.

<sup>3</sup> Supported by ARO Grant DAAG55-98-1-0170, ONR Grant N00014-98-1-0589, and NSF CCR-0113217.

<sup>4</sup> Supported by NSF Grant CCR-9357949 and ONR Grant N00014-98-1-0589.

Keywords: Scheduling; Flow time

---

## 1. Introduction

Consider the problem of sending large files (e.g., bitmap images) through a general topology network. The requests arrive online and the goal is to eventually satisfy all the requests. Since the bandwidth of the links in the network is limited, it makes sense to try to schedule the transmissions in a way that uses the available resources optimally.

In this paper we consider the *online ftp problem*, which is a formal abstraction of the above file transfer problem. We assume that each ftp request specifies source/destination nodes and the size of the file. The goal of the online algorithm is to choose a path that will be used for transmitting each file, and to decide on the transmission rate. The main difference between this model and the (well-studied) models for online routing and admission control [1,3,15,16] is that here we do not assume that the sources have pre-specified transmission rate requirements, i.e., we can deal with nonstreaming types of information. We will study the idealized case where the transmission delays are all zero, and data cannot be buffered along its route.

There are two related measures of performance that can be used to compare different algorithms for the online ftp problem. The first measure is the *total flow time*, i.e., the sum over all jobs of the time that elapses between the instant the ftp request is submitted and the time it is satisfied (including the transmission time). The other measure is the *max-stretch*, which is the maximum over all ratios of the flow time of each request and the smallest time needed to satisfy this request. The second quantity is determined by the link bandwidth and the size of the file. Both measures are useful since they are directly related to the performance of the network perceived by the end-user.

Let  $n$  be the number of requests and  $P$  the maximum ratio between the sizes of the files. Assume that the smallest file can be transferred in one time unit. Let  $F_{\text{MAX}}^*$  denote the optimum max-flow, i.e., the smallest value for the maximum time a request spends in the system. The main results of the paper are algorithms that achieve the *optimum max-stretch* and the *optimum total flow time* using resource augmentation.<sup>5</sup> For the max-stretch algorithm we need to increase capacities by a factor of  $O(\log P)$ , whereas for the total flow time algorithm, we need a factor of  $O(\log F_{\text{MAX}}^*)$  increased capacity.<sup>6</sup> The latter algorithm not only achieves the optimum total flow time, but *simultaneously* approximates many other objective functions, like the maximum flow time, the total square-of-flow-time, etc.

To justify the need for giving larger capacities to the online algorithm (i.e., resource augmentation), we show polynomial lower bounds on both max-stretch and total flow time for the case where both online and offline algorithms use the same capacities. Moreover, we show that in order to achieve a constant competitive ratio against an adaptive adversary we have to give the online algorithm an  $\Omega(\log P / \log \log P)$  factor more capacity for the

---

<sup>5</sup> Throughout this paper, when we refer to an optimum solution, we mean the optimum without any resource augmentation.

<sup>6</sup> Note that  $F_{\text{MAX}}^* \leq nP$  and therefore  $\log F_{\text{MAX}}^* \leq \log n + \log P$ .

max-stretch metric, and an  $\Omega(\sqrt{\log \gamma / \log \log \gamma})$  more capacity for the total flow metric, where  $\gamma = \min\{n, P\}$ .

In the context of machine scheduling, total flow time is known to be a hard metric to approximate [22] and it is only recently that progress has been made towards obtaining algorithms that give total flow time guarantees. In particular, logarithmic-factor resource augmentation was used in [25] to obtain optimum flow time for machine scheduling. Unlike the problem we consider, resource augmentation is not necessary to obtain good approximation ratios for minimizing flow time in the machine scheduling problem if preemption is allowed, as demonstrated by the logarithmic approximation obtained in [23]. Max-stretch was recently proposed as a good metric to measure user satisfaction [7]. Our lower bound on the amount of resource augmentation needed for max-stretch holds in the machine scheduling model as well, and therefore our upper bounds on the amount of augmentation required for max-stretch are also of interest in the machine scheduling model. Notice that these bounds are quite close to each other: the upper bound is  $\tilde{O}(\log P)$  whereas the lower bound is  $\tilde{\Omega}(\log P)$ . Without resource augmentation, the best known competitive ratio for max-stretch in the machine scheduling problem, even on a *single machine* is  $O(\sqrt{P})$  [7,8].

When proving upper bounds, we restrict our algorithms to use a single rate when transmitting a specific file, and do not allow preemption. The competitive ratio is computed against an offline algorithm that does not have these restrictions. Our lower bounds for online flow-time minimization algorithm without resource augmentation (i.e., both the online and the offline algorithms work in the same network) hold even if we remove this restriction, i.e., allow the algorithm to use a time-varying rate when transmitting a file. This contrasts with minimizing flow time for machine scheduling, where a  $\log P$ -competitive preemptive algorithm is known [23]. Also, the lower bound for total flow time is achieved using same-size files. This is in contrast to machine scheduling where the unit jobs case is trivial.

The online ftp problem is a special case of the *set scheduling problem*. In the set scheduling problem we have a set of resources and each job requires a specific subset of these resources (or one of a set of subsets). Set scheduling is a natural generalization of the machine scheduling problem that was extensively studied under several different metrics. (See [21] for a survey of offline approximation algorithms, and [2,7,11,19,20,23,25,26] for a sampling of recent results in online algorithms.) The set scheduling model is similar to the parallel jobs model studied by [14,28]. We show how to apply several techniques developed in the context of machine scheduling to the set scheduling problem (and hence the online ftp problem) for simpler metrics such as makespan and total completion time. In particular, we use the technique that allows us to convert an offline optimization algorithm that maximizes the number of scheduled jobs into an online algorithm that minimizes total completion time [19,20,26]. We also develop new techniques that help us attack more difficult metrics such as total flow time and max-stretch.

Our techniques apply only when the jobs are malleable [9,14,19,28], i.e., extra capacity/resources can be used to reduce the processing time of jobs. Two previously studied examples of such problems are the parallel jobs problem [9,14,28] and the vector scheduling problem [6,10,17]. The techniques developed in this paper can be better understood when compared to the technique of Hall et al. [19,20]. They use offline

$\rho$ -approximation algorithms for offline packing problems to construct  $O(\rho)$ -competitive online algorithms for average completion time. Our techniques allow the transformation of offline packing algorithms that achieve the optimum packing using  $O(\rho)$  resource augmentation for malleable jobs into online algorithms that achieve the optimum flow time using  $O(\rho \cdot \log F_{\text{MAX}}^*)$  resource augmentation. If the online algorithm is not required to work in polynomial time, then an optimum offline solution ( $\rho = 1$ ) can be used.

Significant recent progress has been made in recent years on flow and stretch metrics for scheduling. Muthukrishnan et al. [24] showed that the simple Shortest Remaining Processing Time heuristic gives an  $O(1)$  competitive ratio for the average stretch problem in machine scheduling. This is the same algorithm that was used by Leonardi and Raz [23] who proved that the Shortest Remaining Processing Time rule has an  $O(\log n)$  competitive ratio for minimizing flow time on parallel machines. This algorithm allows jobs to migrate. Becchetti et al. [5] later presented an algorithm which obtains  $O(1)$  competitive ratio for average stretch for the machine scheduling problem without job migration. Again, they use an algorithm which obtained an  $O(\log n)$  competitive ratio for the flow time problem, without migration [2]. This is an interesting set of results, where the same algorithm is often good for multiple measures. Our main result is also an algorithm that is simultaneously good for a large class of objective functions. Subsequent to our work, Epstein and van Stee studied flow time for nonmalleable jobs on a single machine [12]. One of the consequences of their work is a slight improvement in our lower bound on the amount of resource augmentation needed for the flow time problem.

In Section 2 we explain our models. Section 3 contains the main technical contributions of the paper—the lower and upper bounds on the performance of online algorithms using the total flow time and max-stretch metrics. In Section 4 we describe online algorithms for the ftp problem using the makespan and total completion time metrics. Not all online algorithms in Sections 3 and 4 run in polynomial time; polynomial-time online algorithms and offline approximation algorithms are discussed in Section 5. Section 5 also sketches an offline, polynomial-time algorithm for minimizing the makespan for the set scheduling problem (and hence the online ftp problem) if the rate at which a request is serviced is allowed to vary arbitrarily.

## 2. Models and definitions

In the *set scheduling problem* there are  $n$  jobs and  $m$  resources. Job  $j$  has an arrival time (release date)  $a_j$ , a processing time  $p_j$ , and a resource requirement  $S_j$  where  $S_j$  is a subset of  $S$ , the set of resources. We define  $P = \max_j p_j / \min_j p_j$ . The quantity  $P$  plays a crucial role in the analysis of our algorithms. As in traditional scheduling, both the preemptive and nonpreemptive variants are of interest. The Set Scheduling Problem can be formulated as either an offline or an online problem. As in job shop scheduling and multiprocessor scheduling, the performance of an algorithm for this problem can be studied under several different metrics—most notably makespan, total completion time, total flow time, and max-stretch. In this paper we will concentrate mainly on online algorithms. To the best of our knowledge, a systematic study of offline algorithms for the set-scheduling problem has not yet been performed, and may well be an interesting research direction.

The *online ftp problem* is defined as follows. We are given a network  $G = (V, E)$  where all edges have identical bandwidths. Assume that the transmission delay along any link is zero, and that there are no buffers in the network. Once a source starts transmitting data to another node, the other node starts receiving it immediately. Of course the rate at which the sender transmits the data is bounded by the minimum available bandwidth along the route over which the transmission is taking place. Let  $m$  be the number of links in the network, and  $n$  the number of ftp requests. Request  $j$  has an arrival time  $a_j$ , specifies file size  $p_j$ , and a route  $R_j$  over which the data needs to be transmitted. We also address the case where instead of the route, the request specifies only the source and the sink nodes. The former model is closer to the IP world, where the routes are determined by an external algorithm, while the second model is closer to the ATM world, where one can use source routing.

Let  $C_j$  be the completion time of job  $j$  in a schedule. The quantity  $F_j = C_j - a_j$  is called the flow time of job  $j$ . The makespan of a schedule is  $\max_j C_j$ ; total completion time is  $\sum_j C_j$ ; total flow time is  $\sum_j F_j$  and max-stretch is  $\max_j F_j/\tau_j$  where  $\tau_j$  is the time it would take to satisfy job  $j$  if it had the whole network to itself. We also permit jobs to have weights  $w_j$ . In the presence of weights the total completion time and total flow time metrics are defined as  $\sum_j w_j C_j$  and  $\sum_j w_j F_j$ , respectively. Traditionally, the total flow time and max-stretch metrics are considered to be the hardest. These are also perhaps the most interesting metrics as they most directly measure end user experience.

The following theorem captures the hardness of the set scheduling problem—the reduction is straightforward, but is sketched below for completeness.

**Theorem 1.** *The Vertex Color problem reduces (via polynomial-time reductions) to Minimizing Makespan for Set Scheduling in an approximation preserving fashion.*

**Proof.** Let  $G(V, E)$  be an instance of the vertex color problem. We construct an instance of the set scheduling problem with  $S = E$ . For each vertex  $v$  in the original vertex color problem, we introduce a job  $j_v$  with arrival time 0, processing time 1, and resource requirement  $I(v)$ , where  $I(v)$  is the set of all edges incident on  $v$ . Jobs  $j_u$  and  $j_v$  need a common resource iff there is an edge between vertices  $u$  and  $v$ . Thus, jobs  $j_u$  and  $j_v$  can be scheduled at the same time iff  $u$  and  $v$  can have the same color. This establishes a one to one correspondence between the makespan of the set scheduling problem and the number of colors needed for the vertex color problem.  $\square$

For the vertex color problem lower bounds are known for both the approximation ratio ( $\Omega(n^{1-\epsilon})$  unless  $P = NP$  [13]) and competitive ratio ( $\Omega(n^{1/3})$  [4]), yielding corresponding lower bounds for the set scheduling problem. The set scheduling problem trivially reduces to the file transfer problem with routes and rates given as input. Hence, the above lower bounds also hold for the online ftp problem if the routes as well as the transmission rates are given as input. Clearly, to make progress with the set scheduling/online ftp problems, we need to relax the model. The first relaxation we propose is to allow rate control for jobs. Thus each job would be assigned a start time  $s_j$  ( $s_j \geq a_j$ ) and a rate  $r_j$  by the scheduler. The job would execute from time  $s_j$  to  $s_j + p_j/r_j$  and would consume an  $r_j$  fraction of each resource in its resource set  $S_j$  during this interval. More than one jobs may use a resource at the same time. However, the total usage of a resource

at any time must be at most 1. The term *malleable jobs* is commonly used to describe this property [9,14,19,28]. This relaxation is particularly appropriate to the ftp problem: it is possible to control the rate of a TCP connection and more than one connections can use the same link. Further, a connection uses up the same bandwidth on each link along its route.<sup>7</sup>

### 3. Flow time and max-stretch using resource augmentation

#### 3.1. Upper bounds with resource augmentation but no preemption

Assume that all links have the same capacity in the original network; rescale capacities so that this capacity becomes 1. Further rescale time such that the smallest request takes four units of time to finish if it has the entire network to itself. Then the time required to service the largest request (if the request has the entire network to itself) is at most  $4P$ .

Let  $n$  be the number of requests, and  $m$  the number of links. Let  $K = 3 + \log n + \log P$ . We assume that the online algorithm can use a factor  $5K$  of resource augmentation. Thus the online algorithm pretends that the capacity of each link is  $5K$ . We will compare our online algorithm to an offline optimum solution that is only allowed to use the original capacity of 1 on each link.

Let  $w_j$  be the weight of job  $j$ . The online algorithm partitions the network into  $K$  copies,  $G_0 \dots G_{K-1}$ , each with edge capacities 5. We call this algorithm *MRHP* (Most Recent Highest Priority) since at any given time, connections which have been waiting in the system the shortest are the most likely to get scheduled. The online algorithm does its processing only at integral time instants. Scheduling decisions for the  $i$ th copy of the network are made every  $2^i$  time units. Figure 1 describes the behavior of MRHP at time  $t$  such that  $t = 2^k \cdot t'$ , where  $t'$  is odd.

Each  $i$  in Fig. 1 corresponds to the copy  $G_i$ . The same job may get scheduled by multiple copies of the network. The flow time of such a job is taken to be the smallest flow time from all its copies. All the jobs ultimately get scheduled by the online algorithm,

for  $i = 0$  to  $\min\{k, K - 1\}$

1. Let  $S_i$  be the set of requests which arrived in the interval  $[t - 2^i, t)$ ;
2. Find the largest weight subset of  $S_i$  that can be completed in the network  $G_i$  between times  $t$  and  $t + 2^i$ ; (Note: This step may not run in polynomial time in general)
3. Schedule this subset in  $G_i$  such that each request has starting time  $t$ , finishing time  $t + 2^i$ , and a uniform rate during this interval.

Fig. 1. Algorithm MRHP at time  $t = 2^k \cdot t'$ , where  $t'$  is odd.

<sup>7</sup> Instead of allowing a fixed rate  $r_j$  for each job, we could also allow the rate to vary. It turns out that our online algorithms, even though they use just one rate  $r_j$  for job  $j$ , are competitive against optimal solutions which are allowed to vary the rate. For offline algorithms it may help to vary the rate; we will delve into this a little in Section 5.

as  $G_{K-1}$  schedules over a sufficiently large interval to schedule all the jobs by itself the very first time it is invoked.

Let  $F_j$  denote the total time this job spends in the system. Let  $Q_k$  be the total weight of the requests which get scheduled in at least one of the networks  $G_0 \dots G_k$ . Let  $Q_k^*$  be the total weight of all requests that have a flow time of at most  $2^{k+2}$  in the optimum solution. Let  $q_k = Q_k - Q_{k-1}$ , and  $q_k^* = Q_k^* - Q_{k-1}^*$  (for convenience define  $Q_{-1}$  and  $Q_{-1}^*$  to be 0). Each job  $j$  which contributes to  $q_k$  must have a flow time  $F_j \leq 2^{k+1}$  in the MRHP schedule, and each job  $j$  which contributes to  $q_k^*$  must have a flow time  $F_j^* \geq 2^{k+1}$  in the optimum schedule.

**Lemma 2.** For all  $k$  such that  $0 \leq k \leq K - 1$ ,  $Q_k \geq Q_k^*$ .

**Proof.** Let  $S_k^*$  be the set of requests which contribute to  $Q_k^*$ . By definition, each of these requests has a flow time of at most  $2^{k+2}$ . Divide time into intervals of the form  $[i \cdot 2^k, (i + 1) \cdot 2^k)$  for  $i \geq 0$ . Let  $S_k^{(i)*}$  denote the set of requests from  $S_k^*$  which arrive during the  $i$ th interval, and let  $Q_k^{(i)*}$  denote their combined weight. All these jobs are scheduled by the optimum algorithm to finish before time  $(i + 1) \cdot 2^k + 2^{k+2}$ . Hence all these jobs must arrive and finish in the interval  $[i \cdot 2^k, (i + 1) \cdot 2^k + 2^{k+2})$ , which has length  $5 \cdot 2^k$ . Since  $G_k$  has 5 times the original capacity on each edge, and since it has all the jobs in  $S_k^{(i)*}$  available for scheduling during the interval  $[(i + 1) \cdot 2^k, (i + 2) \cdot 2^k)$ , it will schedule jobs with a weight of at least  $Q_k^{(i)*}$  during this interval. Summing up over all  $i$ ,  $Q_k \geq Q_k^*$ .  $\square$

Let  $g$  be any function from  $\mathfrak{N}^+$  to  $\mathfrak{N}^+$ . Let  $\mathcal{F}_g^*$  denote the optimum value of  $\sum_j w_j g(F_j)$  that can be obtained in an unaugmented network, and  $\mathcal{F}_g$  denote the corresponding value obtained by MRHP.

**Theorem 3.**  $\mathcal{F}_g \leq \mathcal{F}_g^*$ , for all nondecreasing functions  $g$  from  $\mathfrak{N}^+$  to  $\mathfrak{N}^+$ .

**Proof.** Since  $q_k^*$  is the weight of jobs whose flow time in the optimum solution belongs to the range  $(2^{k+1}, 2^k]$ , and  $g$  is nondecreasing,  $\mathcal{F}_g^* \geq \sum_{0 \leq k \leq K} g(2^{k+1})q_k^*$ . Similarly,  $\mathcal{F}_g \leq \sum_{0 \leq k \leq K-1} g(2^{k+1})q_k$ .

Let  $W = \sum_j w_j$ . We define  $P(k) = q_k/W$  and  $P^*(k) = q_k^*/W$ . Further, let  $g'(x) = Wg(2^{x+1})$ . Now,  $\mathcal{F}_g^* \geq \sum_{0 \leq k \leq K} g'(k)P^*(k)$  and  $\mathcal{F}_g \leq \sum_{0 \leq k \leq K-1} g'(k)P(k)$ .

$P$  and  $P^*$  are probability density functions, and Lemma 2 implies that  $P^*$  stochastically dominates  $P$ . By definition of stochastic dominance, it is possible to construct a random experiment which yields two variables  $X, X^*$  with the following properties:

- (1)  $P(k)$  is the probability of the event  $X = k$ ,
- (2)  $P^*(k)$  is the probability of the event  $X^* = k$ , and
- (3)  $X \leq X^*$ .

Property 1 implies that  $\mathcal{F}_g \leq \mathbf{E}[g'(X)]$ . Property (2) implies that  $\mathcal{F}_g^* \geq \mathbf{E}[g'(X^*)]$ . Observe that  $g'$  is also a nondecreasing function from  $\mathfrak{R}^+$  to  $\mathfrak{R}^+$ . Hence, property (3) implies that  $g'(X) \leq g'(X^*)$ , which in turn implies that  $\mathbf{E}[g'(X)] \leq \mathbf{E}[g'(X^*)]$ . Putting the above three statements together yields the desired result.  $\square$

Theorem 3 is particularly interesting because it shows that MRHP simultaneously optimizes a very wide class of metrics. In particular, the following results can be obtained as corollaries.

Let  $\mathcal{F}^*$  and  $\mathcal{F}$  denote the total weighted flow times of the optimum and online algorithms, respectively. Let  $F_{\text{MAX}}$  denote the maximum flow time (max-flow) in the schedule obtained by MRHP and  $F_{\text{MAX}}^*$  denote the max-flow in the optimum schedule.

**Corollary 3.1.** *MRHP guarantees that  $\mathcal{F} \leq \mathcal{F}^*$ .*

**Proof.** Let  $g$  be the identity function in the statement of Theorem 3.  $\square$

**Corollary 3.2.** *MRHP guarantees that  $F_{\text{MAX}} \leq F_{\text{MAX}}^*$ .*

**Proof.** For  $p > 0$ , define  $\mathcal{F}_p$  to be  $\sum_j w_j (F_j)^p$ .  $\mathcal{F}_p^*$  is defined analogously. Theorem 3 implies that  $\mathcal{F}_p \leq \mathcal{F}_p^*$  for all  $p > 0$ .  $F_{\text{MAX}}$  and  $F_{\text{MAX}}^*$  are the limiting values of  $(\mathcal{F}_p)^{1/p}$  and  $(\mathcal{F}_p^*)^{1/p}$  respectively as  $p \rightarrow \infty$ . Therefore  $F_{\text{MAX}} \leq F_{\text{MAX}}^*$ .  $\square$

The average stretch of a job can be mimicked using a total weighted flow time objective function by setting the weight  $w_j$  of job  $j$  to  $1/p_j$ . MRHP does not need to know  $K$  in advance—it can maintain an estimate of  $K$  and increment this estimate by one whenever the current value of  $K$  does not suffice to schedule all the requests. Let  $F_{\text{MAX}}^*$  be the optimum max-flow for the given sequence of jobs, given that the shortest job takes one unit time to finish if it has the entire network to itself. Notice that  $F_{\text{MAX}}^* \leq nP$ . The following theorem gives a sharper bound on the amount of resource augmentation needed by MRHP.

**Theorem 4.** *MRHP needs  $O(\log F_{\text{MAX}}^*)$  resource augmentation. Further,  $F_{\text{MAX}}^*$  need not be known in advance.*

The above theorem represents a significant improvement, since  $n$  can be arbitrarily large even in a well behaved system with small max-flow. Section 5 shows how to implement the algorithm in expected polynomial time with  $O(\log n + \log P + \log m)$  resource augmentation.

We now return to the max-flow metric introduced in Corollary 3.2. The max-flow metric ( $F_{\text{MAX}}$ ) is interesting primarily because it relates to the max-stretch metric. We give a simple online algorithm *MMF* (Minimum Max-Flow) that uses only a constant factor resource augmentation. More specifically, MMF uses at most five times the capacity of the original network. MMF assumes that the optimum max-flow is at least  $T$  and at most  $2T$ . (Initially,  $T$  is assumed to be the time required to complete the very first job in the original network.) At times  $t$  which are multiples of  $T/2$ , MMF looks at all requests which arrived



during the last  $T/2$  time units. It then assigns to each of these jobs a rate which is just sufficient for this job to finish in the next  $T/2$  time units. If the load on any edge exceeds five times the capacity of that edge in the original network, MMF doubles  $T$ , aborts the current phase, and waits till the current time becomes a multiple of the new value of  $T/2$ . The following theorem subsumes Corollary 3.2.

**Theorem 5.** *The maximum flow time of a job in the schedule produced by MMF is no larger than the optimum max-flow. MMF runs in time polynomial in  $n$ ,  $m$ , and  $\log P$ .*

**Proof.** Consider all jobs that arrive in the interval  $(\frac{T}{2}(i-1), \frac{T}{2}i]$ . Let the optimum max-flow be  $F_{\text{MAX}}^*$ . All these jobs must finish in the optimal schedule by the time  $(T/2)i + F_{\text{MAX}}^*$ . If  $F_{\text{MAX}}^* \leq 2T$  then the total requirement of all the jobs that need edge  $e$  can be at most  $2T + T/2 = 5T/2$  (recall that the capacity of each edge is one). Hence by *stretching* all these jobs over an interval of length  $T/2$ , no edge can exceed its capacity by more than a factor of 5. If an edge exceeds capacity by more than a factor of 5, we can conclude that  $F_{\text{MAX}}^* > 2T$ , and hence double our estimate of  $T$ .  $\square$

We are now ready to present *MMS* (Minimum Max-stretch) which uses  $O(\log P)$  resource augmentation and guarantees a max-stretch that is no worse than the optimum max-stretch. We first observe that MMF can be modified to guarantee a max-flow that is at most half the optimum value if the amount of capacity on each edge is ten times that in the original network. Let  $p_1$  be the amount of data transfer required by the first job. *MMS* bunches incoming requests into (at most  $\log P$ ) classes, with class  $i$  containing all requests which have a data requirement in the range  $[p_1 \cdot 2^i, p_1 \cdot 2^{i+1})$  ( $i$  may be negative as well). There can be at most  $2 + \log P$  classes. For requests within class  $i$  *MMS* invokes a separate copy of modified MMF. Thus the resource augmentation needed by *MMS* is  $O(\log P)$ . Note that *MMS* does not need to know  $P$  in advance. The fact that the max-flow obtained within each class is at most half the optimum max-flow for that class is sufficient to guarantee that the max-stretch obtained by *MMS* is no more than the optimum max-stretch. The following theorem summarizes the claims made in the above discussion.

**Theorem 6.** *MMS uses  $O(\log P)$  resource augmentation and obtains a max-stretch that is no more than the optimum max-stretch. Further, MMS does not need to know  $P$ . MMS runs in time polynomial in  $n$ ,  $m$ , and  $\log P$ .*

Note that neither MRHP, nor MMF, nor MMS need to get the transmission routes  $R_j$  as input.

**Theorem 7.** *MRHP, MMF, and MMS can obtain optimum values for their respective metrics even if the routes  $R_j$  are not given as input.*

If routes are not provided as input, MRHP, MMF, and MMS as described above would not run in polynomial time. See Theorem 15 for the amount of resource augmentation needed by polynomial-time algorithms.

### 3.2. Lower bounds with preemption but without resource augmentation

We show that without extra capacity, the competitive ratio of any randomized online algorithm which tries to minimize the total flow time (max-stretch, respectively) for the data transfer problem against an oblivious adversary cannot be bounded by any function of the network size. The lower bound for the competitive ratio in terms of the number of jobs,  $n$ , is  $\Omega(\sqrt{n})$  for both metrics. The quantity  $P$  is 1 for the flow-time lower bound, and  $\sqrt{n}$  for the max-stretch lower bound. The lower bounds hold even if the online algorithm is allowed to preempt jobs and use fractional capacities on links but the adversary is not.

Consider the length-3 path  $A - B - C - D$ . Assume that all 3 links have the same bandwidth,  $u$ . Each connection requests the same amount of data,  $r$ . We rescale time so that  $u = r$ , i.e., each request can be serviced in exactly one time unit.

*Total flow time:* The adversary first tosses an unbiased coin. If the outcome is “Heads” it chooses the link  $A - B$  as a special link, else it chooses  $C - D$ . During the first time step, the adversary generates  $k$  requests from  $A$  to  $C$  and  $k$  from  $B$  to  $D$ . The adversary does not do anything for the next  $k - 1$  time units. Then for the next  $k^2$  time units the adversary generates one request per time unit over the special link.

**Lemma 8.** *The expected flow time of any online algorithm on this sequence must be  $\Omega(k^3)$ , even if preemption is allowed and the online algorithm is allowed to use fractional capacities. Further, the optimum flow time for this sequence is  $O(k^2)$  even without using fractional capacities and preemption.*

**Proof.** During the first  $k$  time units, the algorithm can send only  $k$  units of data over the edge  $B - C$ . Hence,  $k$  units of data remains unsent at time  $k$ . Since the special edge is picked randomly by the adversary and not known to the algorithm, the expected amount of unsent data which needs to traverse the special edge is at least  $k/2$ . During the next  $k^2$  time units, even if the special edge is kept continuously busy, the expected amount of unsent data waiting to cross the special edge is at least  $k/2$ . Therefore the expected flow time of the algorithm is at least  $k^2 \cdot (k/2) = \Omega(k^3)$ . The adversary on the other hand will schedule all the requests that need the special edge during the first  $k$  time units to obtain a total flow time of  $O(k^2)$ .  $\square$

Since the number of jobs is  $n = 2k + k^2$ , the competitive ratio of any online algorithm must be  $\Omega(\sqrt{n})$  which does not depend on the network size.

*Max-stretch:* Again, the adversary first tosses an unbiased coin. If the outcome is “Heads” it chooses the link  $A - B$  as a special link, else it chooses  $C - D$ . During the first time step, the adversary generates 1 request from  $A$  to  $C$  and 1 from  $B$  to  $D$ , each of size  $k$ ; for the next  $k - 1$  time units the adversary does nothing. Over the next  $k^2$  time units the adversary generates one request of size 1 every time unit over the special link. The proof of the following lemma is similar to the previous one, and is omitted.

**Lemma 9.** *The expected max-stretch of any online algorithm on this sequence must be  $\Omega(k)$ , even if preemption is allowed and the online algorithm is allowed to use fractional capacities. Further, the optimum max-stretch for this sequence is 2 even without using fractional capacities and preemption.*

The ratio  $P = p_{\max}/p_{\min}$  for this sequence is  $k$ . Since the number of jobs is  $n = 2 + k^2$ , the competitive ratio of any online algorithm must be  $\Omega(\min\{P, \sqrt{n}\})$  which does not depend on the network size. A lower bound of  $\Omega(P^{1/3})$  for the competitive ratio of an online algorithm for the minimum max-stretch problem in the context of machine scheduling was presented in [7].

### 3.3. Lower bounds on the amount of resource augmentation

In this section we give lower bounds on the amount of resource augmentation needed for any randomized online algorithm to achieve a constant competitive ratio. These lower bounds require an adaptive adversary, and assume that the online algorithm is not allowed to preempt requests or change the rate at which a request is being serviced. All our upper bound algorithms work against adaptive adversaries, and do not preempt requests.

**Theorem 10.** *Against an adaptive adversary, any randomized online algorithm that achieves constant competitiveness for max-stretch must use  $\Omega(\min\{n, \log P / \log \log P\})$  resource augmentation.*

**Proof.** The adversary uses a one link network with capacity 1. Let  $u$  be the resource augmentation that the online algorithm uses and let  $k$  be a parameter chosen suitably below. The sequence of requests created by the adversary consists of subsequences  $A_0, A_1, \dots, A_f$ , for some  $f \geq 0$ . The beginning of a new subsequence  $A_i$  is called a *restart*. Initially  $i = 0$ . Each subsequence  $A_i$  consists of requests of size  $L_i$ , one every  $L_i$  time units where  $L_i = (16uk)^{3u-i}$ . Define an *i-phase* to be a time interval between the  $i$ th and the  $i+1$ st restart during which no new jobs of  $A_i$  arrive and no old jobs of  $A_i$  are completed by the online algorithm. Since the algorithm is not allowed to vary the rates, the adversary can determine at the beginning of an *i-phase* how long the *i-phase* would last if no new job arrived. The adversary also knows the bandwidth utilization of the online algorithm during the *i-phase*. If the adversary encounters an *i-phase* that would last at least  $L_i/(8u)$  time units and where jobs of  $A_i$  use more than  $1/3$  units of bandwidth, the adversary increments  $i$  and it restarts. If the adversary does not encounter such an *i-phase*, it stops when  $A_i$  consists of  $k$  jobs.

Note that whenever the adversary restarts, the bandwidth available to the online algorithm for jobs created after the restart is reduced by at least  $1/3$ . Thus the adversary restarts at most  $3u$  times, i.e.,  $f \leq 3u$ . It can be shown inductively that the optimum algorithm can schedule all jobs in  $\bigcup_{l>i} A_l$  (i.e., all jobs of size less than  $L_i$ ) in time at most  $L_i$ . Hence delaying the last job of each size by its size gives an algorithm with max-stretch at most 2.

We show next that the max-stretch of the online algorithm is at least  $k$ . Let  $L_f = (16uk)^{3u-f}$  be the size of the shortest jobs generated by the adversary. When the adversary

creates jobs of size  $L_f$  no  $f$ -phase exists of length at least  $L_f/(8u)$  where jobs of  $A_f$  use more than  $1/3$  units of bandwidth. Since  $k$  jobs of size  $L_f$  are created, there are at most  $2k$   $f$ -phases. The total amount of data of jobs in  $A_f$  transferred during  $f$ -phases where the jobs in  $A_f$  use more than  $1/3$  units of bandwidth is at most  $2k \cdot L_f/(8u) \cdot u = L_f k/4$ . We consider next  $f$ -phases where the jobs in  $A_f$  use at most  $1/3$  units of bandwidth. During the first  $2kL_f$  time units of these  $f$ -phases at most  $2kL_f/3$  data of jobs in  $A_f$  is transferred. Therefore the total amount of data of jobs in  $A_f$  transferred by the online algorithm during the first  $2kL_f$  time units since the last restart is at most  $11kL_f/12$ . Hence, there are some jobs of  $A_f$  left unfinished at time  $2kL_f$  and therefore, there must be some job with a stretch of  $k$ .

It follows that the competitive ratio is at least  $k/2$ . Note that the ratio  $P$  of the maximum job size to minimum job size is  $(16uk)^f$  and that the number  $n$  of jobs is at most  $fk$ . Since  $f \leq 3u$ ,  $n \leq 3ku$  and  $P \leq (16uk)^{3u}$ . If the competitive ratio is a constant, both  $n/3u$  and  $P^{1/(3u)}/(16u)$  must be a constant. The first condition translates to  $u = \Omega(n)$  and the second translates to  $u = \Omega(\log P / \log \log P)$ . Therefore

$$u = \Omega(\min\{n, \log P / \log \log P\}). \quad \square$$

**Theorem 11.** *Let  $\gamma = \min\{n, P\}$ . Against an adaptive adversary, any randomized online algorithm that achieves constant competitiveness for Total Flow Time must use  $\Omega(\sqrt{\log \gamma / \log \log \gamma})$  resource augmentation.<sup>8</sup>*

**Proof.** We define the adversary recursively. Choose a fixed constant  $\delta < 1/6$ , and let  $k > 2$  be a parameter whose value we will specify later. Let  $u$  be the amount of resource augmentation; we restrict our attention to those  $u$  that are multiples of  $\delta$ . The network consists of just a single link and the adversary generates requests which need to transfer 1 unit of data over this link; we call these the *long* jobs. Define a basic interval to be of the form  $[i/(2u), (i+1)/(2u))$ , where  $i$  is a nonzero integer. If the online algorithm schedules a job over the duration  $[t_1, t_2)$ , then we round up  $t_1$  and round down  $t_2$  to be multiples of  $1/(2u)$  without altering the rate which was assigned to this job by the online algorithm. This can only help the online algorithm. Define a basic interval to be “bad” if the online algorithm uses a capacity of more than  $\delta$  on the link during this interval.

The adversary generates a long request at time 0. Then, at time  $i > 0$ , the adversary generates a long request only if the previous time unit  $[i-1, i)$  did not contain a bad basic interval. The adversary stops as soon as it sees  $k\delta/u$  basic bad intervals or has generated  $k$  long requests. Also, during each basic bad interval, the adversary recursively generates the lower bound sequence for  $(u-\delta)$  resource augmentation with time scaled down by a factor of  $1/(4ku)$ .

Let  $F^*(u)$  denote the optimal flow time for the above sequence, and  $F(u)$  denote the flow time obtained by the online algorithm. Let  $N(u)$  and  $P(u)$  denote the number of jobs and the ratio of the largest to the smallest job generated by the adversary, respectively. We

<sup>8</sup> In a preliminary version of this paper [18] we claimed a slightly weaker lower bound of  $\Omega(\sqrt{\log \gamma / \log \log \gamma})$ .

use the lower bound example from Lemma 8 as a base case,<sup>9</sup> where  $u = 1$ . For the base case, choose  $N(1) = k$ . Then  $F(1) = \Theta(k^{1.5})$ ,  $F^*(1) = \Theta(k)$ , and  $P(1) = 1$ .

The optimum way to schedule all the jobs generated by the adversary is to schedule a long job during a time unit which does not contain a basic bad interval. There will never be more than one long job waiting in the system. Further, since there are never more than  $k\delta/u$  basic bad intervals, the total duration of the schedule is at most  $k + k\delta/u < 2k$ . Therefore, the total contribution of the long jobs to the flow time is at most  $2k$ ; this also ensures that all the jobs generated during the recursive procedure fit into the basic interval and do not overflow. During each basic bad interval time is scaled down by a factor of  $1/(4ku)$ , which implies the recurrence

$$F^*(u) \leq 2k + \frac{1}{4ku} F^*(u - \delta) \frac{k\delta}{u} < 2k + F^*(u - \delta).$$

Using the base case, we obtain

$$F^*(u) < 2ku/\delta + O(k) = O(ku). \quad (1)$$

Also,  $P(u) \leq (4ku)P(u - \delta)$ , which simplifies to  $P(u) = (ku)^{O(u)}$ , and  $N(u) = k + (k\delta/u)N(u - \delta) = k^{O(u)}$ .

We now provide a lower bound on  $F(u)$ . First, we study the case in which the adversary terminated before finding  $k\delta/u$  basic bad intervals. Since the capacity available to the algorithm is  $u$ , any job must originally be scheduled for at least  $1/u$  time. By rounding up the start and rounding down the end of an interval of length at least  $1/u$  to multiples of  $1/(2u)$ , the length of the interval decreases by a factor of at most 2. Hence, the total data transferred originally by the algorithm is at most twice the data transfer after the rounding. If the number of bad basic intervals is less than  $k\delta/u$  then the total data transfer during these bad intervals is at most  $k\delta$ . Further, the bandwidth used during the good basic intervals is at most  $2k\delta$  during the first  $2k$  time units. Hence, the total data transferred originally by the algorithm is at most  $2(k\delta + 2k\delta) = 6k\delta$  during the first  $2k$  time units. Therefore there is at least  $k(1 - 6\delta)$  amount of data left unsent at the end of the  $2k$  time units; the flow time because of this data must be at least  $(k(1 - 6\delta))^2/(4u) = \Omega(k^2/u)$ . If on the other hand, the adversary terminates because it found enough bad basic intervals, then the flow time must be at least  $(1/4ku)(k\delta/u)F(u - \delta) = F(u - \delta)/(4u^2)$ . Therefore,

$$F(u) = \min\left\{\frac{(k(1 - 6\delta))^2}{4u}, \frac{F(u - \delta)}{4u^2}\right\}.$$

For large  $k$ , the second term is always smaller than the first (inductively) and we obtain

$$F(u) = k^{3/2}/u^{O(u)}.$$

Now, the competitive ratio is  $C(u) = \sqrt{k}/u^{O(u)}$ . For  $C(u)$  to be constant,  $k$  must equal  $u^{O(u)}$ . We now lower bound  $u$  in terms of  $n = N(u)$  and  $P = P(u)$ , given that  $k = u^{O(u)}$ . We can increase  $k$  as long as neither of the two conditions  $N(u) = k^{O(u)}$  and  $P(u) = (ku)^{O(u)}$  is satisfied; hence for a constant competitive ratio either  $n = k^{O(u)} = u^{O(u^2)}$ ,

<sup>9</sup> Lemma 8 uses a length three path as opposed to a single edge network used in the current theorem; it is easy to see that the proofs in the current theorem go through even for a length three path.

or  $P = (ku)^{O(u)} = u^{O(u^2)}$ . In other words, if  $\gamma = \min\{n, P\}$ , then  $\gamma = u^{O(u^2)}$ , or  $u = \Omega(\sqrt{\log \gamma / \log \log \gamma})$ .  $\square$

#### 4. Online algorithms for makespan and total completion time

In contrast to the flow and stretch metrics studied in the previous section, standard techniques can be used to obtain constant competitive online algorithms for makespan and average completion time for the online ftp problem without the use of resource augmentation. We outline the details below.

*Makespan:* Define  $\lambda$  as the maximum amount of data that needs to be transferred over an edge in the network. We rescale time so that one unit of data can be transferred over a link in one unit of time. Let  $a_{\text{MAX}}$  be the time at which the last request arrives. Let  $L$  be the quantity  $\max(a_{\text{MAX}}, \lambda)$ .  $L$  is a lower bound on the makespan of any schedule. The online algorithm does the following:

It maintains a guess  $\tilde{\lambda}$  for the value of  $L$ . We assume that the first request arrives at time 0. The initial value of  $\tilde{\lambda}$  is set to  $p_1$ , the amount of data transfer needed by the first request. Each time a request arrives, the algorithm recomputes  $L$ . If  $L > \tilde{\lambda}$ ,  $\tilde{\lambda}$  is reset to  $\max(L, 2\tilde{\lambda})$ . The online algorithm schedules a newly arrived request  $j$  to execute from time  $\tilde{\lambda}$  to  $2\tilde{\lambda}$ , with a rate of  $p_j/\tilde{\lambda}$ . It is easy to see that the above algorithm does not violate capacity constraints. Let  $U$  represent the final value of  $\tilde{\lambda}$ ; by construction  $U$  is at most  $2L$ . The makespan is at most  $2U + U + U/2 + \dots < 4U$ . We can now claim the following result.

**Theorem 12.** *The above algorithm is 8-competitive.*

If routes are given as input,  $\tilde{\lambda}$ , and hence  $L$ , can be computed efficiently and the above algorithm runs in polynomial time. Therefore, it is also an offline approximation algorithm. Moreover, an offline algorithm can compute the exact value of  $L$  rather than maintain a guess. Hence, the offline algorithm can provide an approximation guarantee of 2.

If routes are not given as part of the input, computing  $\tilde{\lambda}$  is equivalent to the integer routing problem, and our approach does not result in an  $O(1)$  approximation algorithm. The above online algorithm still obtains an 8-approximation, but it cannot be implemented to run in polynomial time unless  $P = NP$ .

It is easy to obtain a lower bound of 1.25 on the competitive ratio of any algorithm for minimizing makespan for the online ftp problem, even with given routes. Consider again the three link network used in Section 3.2. The adversary randomly chooses one of the links  $A - B$  and  $C - D$  to be special. At time  $t = 0$ , the adversary generates a request from  $A$  to  $C$  and one from  $B$  to  $D$ , each requiring one unit of time to complete. At time  $t = 1$  the adversary generates another unit request, this time over the special link. The optimum makespan is 2, but the expected makespan of any online algorithm must be at least 2.5. Notice that the above lower bound assumes only an oblivious adversary, whereas the upper bound is against an adaptive adversary.

*Total completion time:* The general scaling technique outlined by Hall et al. [19,20] directly results in a 4-competitive online algorithm for the total completion time metric, regardless of whether routes are given as part of the input. Their technique requires an offline algorithm that can pack an optimum number of requests into a given interval. This problem is NP hard, and therefore, our online algorithm does not run in polynomial time. An  $O(\log m)$ -competitive polynomial-time algorithm is outlined in Section 5.

The lower bound example outlined above for makespan gives a lower bound of 1.1; the optimum total completion time for the above sequence is 5 whereas any online algorithm must have an expected total completion time of at least 5.5.

## 5. Polynomial time approximation and online algorithms

In this section we give offline approximation algorithms for total completion time, makespan, total flow time, average stretch, maximum flow time, and maximum stretch that run in polynomial time. The algorithms for total completion time and makespan approximate the optimum performance without resource augmentation. The algorithms for the remaining metrics achieve optimum performance using either a constant-factor or a polylogarithmic-factor resource augmentation. We conclude the section by giving polynomial-time algorithms with optimum makespan under two different relaxations of our model:

- (1) We relax the condition that the rate of a job has to be constant: we give a polynomial-time algorithm that varies the rates and achieves optimum makespan.
- (2) We assume that the start time  $s_j$  is part of the input and show that then the problem can be solved in polynomial time.

**Theorem 13.** *There exists an algorithm that achieves an  $O(\log m)$ -approximation of the total completion time for the online ftp problem in time polynomial in  $n$  and  $m$ , regardless of whether routes are given as part of the input.*

**Proof.** Consider the problem of maximizing the number of ftp requests (out of a given set of requests, all of which have the same arrival time) that can be scheduled over a given period of time. Let  $N^*$  denote the value of the optimum solution. We first describe a polynomial-time algorithm that can schedule at least  $N^*$  requests in an interval which is at most a factor  $O(\log m)$  larger than the original interval. The algorithm uses multicommodity flow followed by randomized rounding [27] as described below:

Obtain the multicommodity flow relaxation by allowing requests to complete fractionally. This relaxation is a linear program and can be solved in polynomial time. Let  $x_i$  denote the fraction of request  $i$  completed in the optimum fractional solution. Round the fractional solution by choosing  $\tilde{x}_i = 1$  with probability  $\min\{x_i \log m, 1\}$  and  $\tilde{x}_i = 0$  otherwise. Here  $\tilde{x}_i = 1$  denotes that request  $i$  will be chosen. If the number of requests satisfied by the integer solution is less than the fractional optimum, or if the integer solution exceeds capacity

on any edge by a factor greater than  $4 \log m$ , then repeat the rounding process, else use the current  $\tilde{x}_i$  as the integer solution and terminate.

Standard Chernoff bound arguments imply that the probability that the rounding process will be successful is  $1 - o(1)$  during each iteration; for our purposes it is sufficient to assume that the probability is at least  $1/2$ . The expected running time of the algorithm is polynomial, and when the algorithm terminates, we are guaranteed to have scheduled at least  $N^*$  requests and exceeded capacities by at most a factor  $O(\log m)$ . Since the requests are malleable, exceeding capacities by  $O(\log m)$  is equivalent to stretching the scheduling interval by the same factor.

Plugging this into the general technique of Hall et al. [19,20] results in an  $O(\log m)$ -competitive polynomial-time online algorithm for the total completion time of ftp requests.  $\square$

A polynomial-time 2-approximation for the makespan of the ftp problem when routes are given as part of the input follows from the discussion in Section 4; Theorem 13 results in an  $O(\log m)$ -approximation if routes are not provided as input. Hence, we obtain:

**Theorem 14.** *There exists an algorithm that achieves a 2-approximation for makespan in time polynomial in  $n$  and  $m$  if routes are given as part of the input, and  $O(\log m)$  otherwise.*

We now describe how to implement algorithm MRHP in polynomial time. The only step of MRHP which might take super-polynomial time is step 2, finding the largest weight subset  $A_i$  of  $S_i$  that can be completed between times  $t$  and  $t + 2^i$ . To implement it in expected polynomial time we need to add  $\log m + 2eK$  to the capacity of each edge, where  $K = \log n + \log P + 3$ .

We use first a linear programming relaxation of the problem, then round it probabilistically and finally show that with high probability no edge capacity constraint is violated:

1. The linear program uses for each job  $j$  a variable  $x_j$  and maximizes  $\sum_{j \in S_i} w_j x_j$  under the constraint that for each edge  $e$ ,  $\sum_{j \text{ uses } e} x_j p_j / 2^i \leq 1$  and that for each  $j$ ,  $x_j \geq 0$ . Let  $x_j^*$  denote the value of  $x_j$  in the solution.
2. We probabilistically round each job  $j$  for each network  $i$  such  $P(j \in A_i) = x_j^*$ . Let  $X_e$  be the random variable denoting the load of edge  $e$  in  $G$ .
3. The expected value  $\mu$  of  $X_e$  is  $\sum_{0 \leq i < K} \sum_{j \text{ uses } e} x_j^* p_j / 2^i \leq K$ . Using Chernoff bounds with  $\delta = (\log m + 2eK) / \mu - 1$  shows that

$$Pr(X_e > \log m + 2eK) \leq \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^\mu \leq \left( \frac{e}{\log m + 2eK} \right)^{\log m + 2eK} < \frac{1}{(m(nP)^{2e})}$$

Thus, the probability that one of the edge capacities overflows is at most  $1/(nP)^{2e}$  in which case we simply redo the rounding step.

Note that both MMF and MMS already run in time polynomial in  $n$ ,  $m$ , and  $\log P$  if routes are given as part of the input. The same ideas that we outlined above for total flow time also result in polynomial-time algorithms for the max-flow and max-stretch problems when routes are not given as input.



**Theorem 15.** *There exist (online and offline) algorithms that run in time polynomial in  $n$ ,  $m$ , and  $\log P$  and*

- *achieve optimum total flow-time or average stretch with an expected  $O(\log n + \log P + \log m)$ -factor resource augmentation regardless of whether routes are given as part of the input;*
- *achieve optimum maximum flow time with a constant-factor resource augmentation if the routes are given as part of the input, and expected  $O(\log n + \log m)$  resource augmentation otherwise;*
- *achieve optimum maximum stretch with an  $O(\log P)$ -factor resource augmentation if the routes are given as part of the input, and expected  $O(\log n + \log m + \log P)$  resource augmentation otherwise.*

We finally relax some of our conditions. Consider first the case that the rate of jobs can vary. Assume that the optimum makespan is  $M$ . We present a linear program that given  $M$  checks whether there exists a feasible solution. By performing a binary search over  $M$ , with  $0 < M \leq nP$  and assuming that time is rescaled so that the shortest job takes one time unit, we get a polynomial-time algorithm that finds the optimum makespan.

We assume w.l.o.g. that the first job arrives at time 0. Break the time from 0 to  $M$  into intervals whenever a new job arrives and number the time intervals from 1 to  $n$ . Let  $l_i$  be the length of interval  $i$  and let  $a_{\text{MAX}}$  be the arrival time of the last job. Note that the length of the last interval is  $M - a_{\text{MAX}}$ . There is a variable  $x_{j,i}$  for each interval  $i$  and each job  $j$ . The linear program checks whether there is a nonnegative assignment for the variables  $x_{j,i}$  such that

- (1) for each job  $j$ ,  $\sum_i x_{j,i} l_i \geq p_j$ ,
- (2) for each edge  $e$  and interval  $i$ ,  $\sum_{j \text{ uses } e} x_{i,j} \leq 1$ , and
- (3) for each job  $j$  and interval  $i$  such that  $j$  arrived after  $i$ ,  $x_{j,i} = 0$ .

## 6. Conclusions

In this paper we considered the *online ftp problem*. The goal is to service a sequence of file transfer requests given bandwidth constraints of the underlying communication network. For several metrics of interest (average flow time, max-stretch etc), it is provably hard to obtain sub-polynomial competitive ratios. Hence, it is worthwhile trying to study these problems in a resource-augmentation model. In this model, the online algorithm is given greater capacity than the offline optimum it is compared against. This corresponds to the concept of “over-engineering” in real-life networks, where extra capacity is built in to cope with inefficiencies in network protocols.

The main result of the paper is a technique that leads to algorithms that optimize several natural metrics, such as max-stretch, total flow time, max flow time, and total completion time. In particular, we show how to achieve optimum total flow time and optimum max-stretch if we increase the capacity of the underlying network by a logarithmic factor.

We also gave polylogarithmic lower bounds on the resource augmentation factor necessary in order to keep the total flow time and max-stretch within a constant factor of optimum.

## References

- [1] J. Aspnes, Y. Azar, A. Fiat, S. Plotkin, O. Waarts, On-line load balancing with applications to machine scheduling virtual circuit routing, *J. ACM* 44 (3) (1997) 486–504.
- [2] B. Awerbuch, Y. Azar, S. Leonardi, O. Regev, Minimizing the flow time without migration, in: *Proceedings of the 31st Annual ACM Symposium on Theory of Computing*, 1999, pp. 198–205.
- [3] B. Awerbuch, Y. Azar, S. Plotkin, Throughput competitive online routing, in: *34th IEEE Symposium on Foundations of Computer Science*, 1993, pp. 32–40.
- [4] Y. Bartal, A. Fiat, S. Leonardi, Lower bounds for on-line graph problems with application to on-line circuit optical routing, in: *Proc. of the 28th Symposium on Theory of Computation*, 1996, pp. 531–540.
- [5] L. Becchetti, S. Leonardi, S. Muthukrishnan, Scheduling to minimize average stretch without migration, in: *12th Annual ACM–SIAM Symposium on Discrete Algorithms*, 2000, pp. 548–557.
- [6] J.E. Beck, D.P. Siewiorek, Modeling multicomputer task allocation as a vector packing problem, in: *9th International Symposium on Systems Synthesis*, 1996, pp. 115–120.
- [7] M. Bender, S. Chakrabarti, S. Muthukrishnan, Flow stretch metrics for scheduling continuous job streams, in: *Ninth Annual ACM–SIAM Symposium on Discrete Algorithms*, 1998, pp. 270–279.
- [8] M. Bender, S. Muthukrishnan, R. Rajaraman, Improved algorithms for stretch scheduling, in: *13th Annual ACM–SIAM Symposium on Discrete Algorithms*, 2002, pp. 762–771.
- [9] S. Chakrabarti, C. Phillips, A. Schulz, D. Shmoys, C. Stein, J. Wein, Improved scheduling algorithms for minsum criteria, in: *23rd International Colloquium on Automata Languages, and Programming*, 1996, pp. 646–657.
- [10] C. Chekuri, S. Khanna, On multi-dimensional packing problems, in: *Tenth ACM–SIAM Symposium on Discrete Algorithms*, 1999, pp. 185–194.
- [11] C. Chekuri, R. Motwani, B. Natarajan, C. Stein, Approximation techniques for average completion time scheduling, *SIAM J. Comput.* 31 (1) (2001) 144–166.
- [12] L. Epstein, R. van Stee, Optimal on-line flow time with resource augmentation, in: *Proceedings of 13rd International Conference on Fundamentals of Computation Theory (FCT)*, 2001, pp. 472–482.
- [13] U. Feige, J. Kilian, Zero knowledge the chromatic number, *J. Comput. System Sci.* 57 (2) (1998) 187–199.
- [14] A. Feldmann, M. Kao, J. Sgall, S. Teng, Optimal online scheduling of parallel jobs with dependencies, *J. Combin. Optim.* 1 (4) (1998) 393–411.
- [15] J. Garay, I. Gopal, Call preemption in communication networks, in: *IEEE INFOCOM*, 1992, pp. 1043–1050.
- [16] J. Garay, I. Gopal, S. Kutten, Y. Mansour, M. Yung, Efficient on-line call control algorithms, *J. Algorithms* 23 (1997) 180–194.
- [17] M. Garofalakis, Y. Ioannidis, Scheduling issues in multimedia query optimization, *ACM Comput. Surveys* 27 (4) (1995) 590–592.
- [18] A. Goel, M. Henzinger, S. Plotkin, E. Tardos, Scheduling data transfers in a network the set scheduling problem, in: *31st ACM Symposium on Theory of Computing*, 1999, pp. 189–197.
- [19] L.A. Hall, D.B. Shmoys, A.S. Schulz, J. Wein, Scheduling to minimize average completion time: off-line on-line approximation algorithms, *Math. Oper. Res.* 22 (1997) 513–544.
- [20] L.A. Hall, D.B. Shmoys, J. Wein, Scheduling to minimize average completion time: off-line on-line algorithms, in: *7th Annual ACM–SIAM Symposium on Discrete Algorithms*, 1996, pp. 142–151.
- [21] D. Hochbaum, *Approximation Algorithms for NP-hard Problems*, PWS, 1997.
- [22] H. Kellerer, T. Tautenhahn, G.J. Woeginger, Approximability nonapproximability results for minimizing total flow time on a single machine, *SIAM J. Comput.* 28 (4) (1999) 1155–1166.
- [23] S. Leonardi, D. Raz, Approximating total flow time on parallel machines, in: *Twenty-Ninth Annual ACM Symposium on Theory of Computing*, 1997, pp. 110–119.
- [24] S. Muthukrishnan, R. Rajaraman, A. Shaheen, J. Gehrke, Online scheduling to minimize average stretch, in: *Twenty-Ninth Annual IEEE Symposium on Foundations of Computer Science*, 1999, pp. 433–442.

- [25] C. Phillips, C. Stein, E. Torng, J. Wein, Critical scheduling via resource augmentation, *Algorithmica* 32 (2) (2002) 163–200.
- [26] C. Phillips, C. Stein, J. Wein, Scheduling jobs that arrive over time, *Math. Program. B* 82 (1–2) (1998) 199–224.
- [27] P. Raghavan, C. Thompson, Randomized rounding, *Combinatorica* 7 (1987) 365–374.
- [28] J. Sgall, Randomized on-line scheduling of parallel jobs, *J. Algorithms* 21 (1996) 149–175.