

# Orchestration of Network-Wide Active Measurements for Supporting Distributed Computing Applications

Prasad Calyam, *Student Member, IEEE*, Chang-Gun Lee, *Member, IEEE*, Eylem Ekici, *Member, IEEE*, Mark Haffner, *Student Member, IEEE*, and Nathan Howes, *Student Member, IEEE*

**Abstract**—Recent computing applications such as videoconferencing and grid computing run their tasks on distributed computing resources connected through networks. For such applications, knowledge of the network status such as delay, jitter, and available bandwidth can help them select proper network resources to meet the Quality-of-Service (QoS) requirements. Also, the applications can dynamically change the resource selection if the current selection is found to experience poor performance. For such purposes, Internet Service Providers (ISPs) have started to instrument their networks with Network Measurement Infrastructures (NMIs) that run active measurement tasks periodically and/or on demand. However, one problem that most network engineers have overlooked is the *measurement conflict problem*, which happens when multiple active measurement tasks inject probing packets into the same network segment at the same time, resulting in misleading reports of network performance due to their combined effects. This paper proposes enhanced Earliest Deadline First (EDF) algorithms that allow “Concurrent Executions” to orchestrate offline/online measurement jobs in a conflict-free manner. The simulation study shows that our measurement scheduling mechanism can improve the schedulable utilization of offline measurement tasks up to 300 percent and the response time of on-demand jobs up to 50 percent. Further, we implement and deploy our scheduling mechanism in a real working NMI for monitoring the Internet2 Abilene network. As a case study, we show the utility of our algorithms in the widely used Network Weather Service (NWS).

**Index Terms**—Active network probes, measurement conflict, real-time scheduling, concurrent execution, Network Weather Service.

## 1 INTRODUCTION

RECENT computing applications such as videoconferencing and grid computing utilize distributed computing resources connected through the Internet. Thus, their user-level performance relies on the status of the Internet paths that they use.<sup>1</sup> If we can measure and predict the status, we can select the computing resources and their connecting Internet paths that can soft guarantee the user-level quality of service (QoS).

Therefore, for the success of distributed computing applications, it is critical to collect Internet status measurements in an accurate and timely manner. Fortunately, Internet Service Providers (ISPs) have started to instrument their networks with Network Measurement Infrastructures

(NMIs) [3], [4], [5] for continuous monitoring and estimation of networkwide status. For this, they use active measurement tools such as Ping, Traceroute, H.323 Beacon [2], Iperf [6], Pathchar [7], and Pathload [8] that actively inject probing packets to collect useful measurements such as end-to-end delay, jitter, loss, bandwidth, and so forth. The NMIs periodically run these measurement tools on the measurement servers at strategic points to collect the periodic sampling of network status, which is essential for network status prediction [9]. They can also run the measurement tools on demand for applications that require a more detailed look about certain network paths.

When executing the periodic and on-demand measurement jobs, an important problem that we recently observed is the *measurement conflict problem*, which has been overlooked by most network engineers. Since active measurement tools consume a nonnegligible amount of network resources for injecting probing packets, if two or more measurement jobs run concurrently over the same path, they can interfere with each other, resulting in misleading reports of network status. Our experiment in Fig. 1 illustrates the measurement conflict problem. In the experiment, we connect two measurement servers by a local area network (LAN) testbed with a bandwidth of 1,500 kilobits per second (Kbps) and run one H.323 videoconferencing session at a 768 Kbps dialing speed as the background traffic. Thus, the remaining bandwidth should be approximately 732 Kbps. Given that the streaming media and videoconferencing traffic is essentially User Datagram Protocol (UDP) traffic, Iperf in the UDP mode is popularly used to measure the available bandwidth. Thus, we make the two servers occasionally initiate Iperf jobs to monitor

1. There are several ways to map the network measurements to the user-level quality. The E-Model [2] is an example that estimates user-level voice-over-IP (VOIP) quality from the measured network status.

- P. Calyam is with OARnet and the Ohio State University, 1224 Kinnear Road, Columbus, OH 43212. E-mail: pcalyam@oar.net.
- C.-G. Lee is with the School of Computer Science and Engineering, Seoul National University, Seoul, 151-742, Korea. E-mail: cglee@snu.ac.kr.
- E. Ekici is with the Department of Electrical and Computer Engineering, the Ohio State University, 320 Dreese Laboratory, 2015 Neil Avenue, Columbus, OH 43210. E-mail: ekici@ece.osu.edu.
- M. Haffner and N. Howes are with the Department of Electrical and Computer Engineering, the Ohio State University, 1224 Kinnear Road, Columbus, OH 43212. E-mail: {haffner.12, howes.16}@osu.edu.

Manuscript received 12 July 2006; revised 7 Feb. 2007; accepted 14 May 2007; published online 6 June 2007.

Recommended for acceptance by A. Zomaya.

For information on obtaining reprints of this article, please send e-mail to: [tc@computer.org](mailto:tc@computer.org), and reference IEEECS Log Number TC-0270-0706.

Digital Object Identifier no. 10.1109/TC.2007.70745.

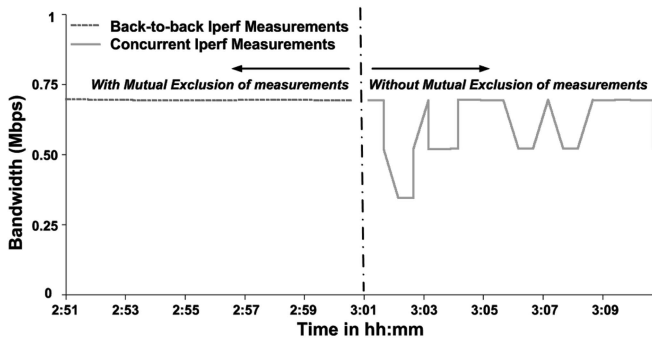


Fig. 1. Iperf test results with and without mutual exclusion of measurements.

the available bandwidth. When we make two Iperf jobs run back to back with mutual exclusion (shown in the left half of Fig. 1), their measurements are in agreement with our expectation. However, when we intentionally make two Iperf execution durations overlap (shown in the right half of Fig. 1), it causes a misrepresentation of the remaining bandwidth, which is merely due to conflicts of two Iperf jobs. This implies that, if measurement tools are initiated without being orchestrated with each other, their execution duration may overlap, resulting in misleading measurement reports.

This observation motivates a *scheduling problem of measurement jobs for orchestrating them to prevent conflicts while still providing the periodicity of periodic measurement jobs and quick response to the on-demand measurement jobs*. The nature of the problem is similar to real-time scheduling even though the time granularity of periods is much coarser (in order of minutes) than that of the classical real-time systems. The measurement scheduling problem, however, has a fundamental difference from the classical real-time scheduling problems: More than one measurement job can be scheduled at the same time on the same server and the same network path as long as they can produce the correct measurement data. We call this a “concurrent execution” of multiple jobs with “no conflict.”

This paper proposes conflict-free scheduling algorithms for measurement jobs leveraging the real-time scheduling principles and the concurrent execution. More specifically, the contributions of this paper can be summarized as follows:

- We propose an offline scheduling algorithm based on the EDF principle [10], but which allows concurrent execution if possible, which can significantly improve the schedulability of a given set of periodic measurement tasks.
- We propose an online mechanism that can steal leftover times from the offline schedule to serve on-demand measurement requests as early as possible without violating the periodicity requirements of existing measurement tasks.
- We implement an actual NMI scheduling framework equipped with the proposed scheduling mechanisms to measure an operational network, that is, the Internet2 Abilene network.

The rest of this paper is organized as follows: Section 2 summarizes the related work. Section 3 formally defines the

measurement task scheduling problem. Section 4 presents our offline and online measurement scheduling algorithms. Section 5 presents our case study for applying the proposed scheduling algorithms to the Network Weather Service (NWS) and its use for distributed computing. In Section 6, our experimental results from both simulations and actual implementation are presented. Finally, Section 7 concludes the paper.

## 2 RELATED WORK

Many of the earliest NMIs used simple *Ping* and *Traceroute* measurements without paying attention to possible overlaps of their execution durations. This is acceptable since they are neither CPU nor channel intensive, allowing overlaps without causing measurement conflicts. However, many of today’s NMIs, such as NLANR AMP [3], Internet2 E2EpiPES [5], NWS [9], Surveyor [11], and RIPE [12], employ toolkits that have several CPU and/or channel-intensive measurement tools, which may cause measurement conflict problems. Nevertheless, these NMIs use a simple scheme that creates *cron* jobs that start active measurements at the planned periodic time points, without paying attention to avoiding measurement conflicts. As a result, they can give erroneous measurement results and fail to reflect the actual network status.

To address the measurement scheduling problem, [3] and [11] use a simple round-robin approach where measurement servers take turns such that only one tool executes at a time. In NMIs such as in [5], a resource broker scheduling scheme is used. Using this resource broker scheme, multiple measurement requests are queued for scheduling and executed on a first-come, first-serve basis on a measurement server. The NWS uses a token-passing mechanism [13] in an attempt to meet the measurement periodicity requirements while obtaining accurate network status information. This mechanism allows only a single server in possession of a token to initiate measurements. The round-robin, resource broker, and token-passing mechanisms are similar in principle, that is, they allow only one instance of measurement to be executed at a time. Therefore, they cannot leverage the concurrent execution of multiple measurement jobs and, hence, limit the schedulability.

In addition to our contributions from the scheduling perspective, another significant contribution is our systematic scheduling framework that automates the whole process from the measurement specification to the runtime measurement data collection. None of the previous schemes provides such a systematic framework. As a result, existing schemes require considerable time and effort to specify distinct sampling requirements, add or delete measurement tasks, and generate measurement schedules accordingly. Furthermore, it is hard to implement the policy contracts among multiple ISPs for measurements across ISP borders. With our systematic measurement framework, however, the entire process can be automated and the manual effort can be minimized.

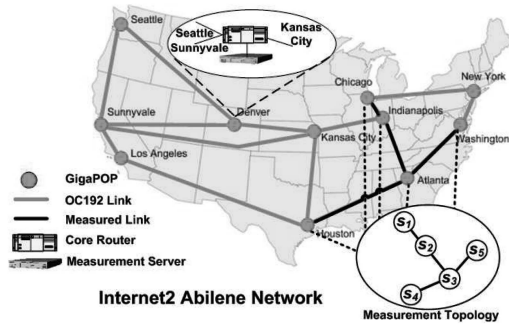


Fig. 2. Measurement topology constructed for a set of network paths.

### 3 PROBLEM DESCRIPTION AND TERMINOLOGY

An ISP deploys measurement servers at strategic points to continuously estimate the networkwide status. The measurement servers measure the network paths to other servers. They are attached to core routers, as shown in the case of the Denver core router in Fig. 2. The paths to be measured are specified by a measurement topology, which can be formally represented by a graph  $G = (N, E)$ , where  $N$  is the set of measurement servers and  $E$  is the set of edges between a pair of servers. Fig. 2 shows an example measurement topology that consists of measurement servers  $N = \{S_1, S_2, S_3, S_4, S_5\}$  and edges among  $S_1, S_2, S_3, S_4,$  and  $S_5$ .

On top of the measurement topology, a set of periodic measurement tasks is specified. Each periodic measurement task  $\tau_i$  is specified to measure a path from a source server  $src_i$  to a destination server  $dst_i$  using an active measurement tool  $tool_i$ . The measurement should be periodically repeated with period  $p_i$ . The  $j$ th instance (or *job*) of  $\tau_i$  is denoted by  $\tau_{ij}$ . The time when the  $j$ th job  $\tau_{ij}$  is released is called the *release time* and is simply given by  $(j - 1)p_i$ . The execution time of a single measurement instance is denoted by  $e_i$ . Then, a periodic measurement task can be represented using the similar notion of a real-time periodic task as follows:

$$\tau_i = (src_i, dst_i, tool_i, p_i, e_i).$$

The set of all offline specified periodic measurement tasks is denoted by

$$\Gamma = \{\tau_1, \tau_2, \dots, \tau_n\}.$$

We also define a *hyperperiod* for task set  $\Gamma$  as the least common multiple of all of the task periods in the set. Thus, each hyperperiod repeats the same pattern of release times. Therefore, the same schedule constructed for a single hyperperiod can be used repeatedly.

In addition to such offline specified measurement tasks, there can be on-demand measurement requests to quickly collect customized measurements. For example, an Internet network engineer might want to trace back the source of a denial-of-service (DoS) attack as soon as possible by running on-demand measurement jobs over suspicious paths [14]. Such an on-demand measurement request is denoted by

$$J_k = (src_k, dst_k, tool_k, e_k).$$

For such an on-demand request, a quick response is desirable. Thus, as a performance metric, we use the *response time*, which is defined as the time difference between the time when the measurement job is requested and the time when the request is finally served.

Our problem is scheduling the aforementioned offline and online measurement jobs on a given measurement topology. Unlike the OS-level schedule that determines when the OS threads and packets can be executed and transmitted, the measurement-level scheduling problem determines the start and stop times of a measurement tool whose execution can last a few minutes to have a statistically stable measure. For such measurement-level scheduling, an important constraint is a measurement conflict problem. Overlapping the execution intervals of two measurement jobs may or may not be problematic, depending on the measurement tools used. If a measurement tool is neither CPU intensive nor channel intensive, like Ping, it does not interfere with other tools. Thus, overlapping its execution interval with others on the same server and/or path can still give us correct measurement reports. Such an overlap is called a “concurrent execution” with no conflict, which is desirable to improve the schedulability. On the other hand, other active measurement tools such as Iperf [6] and Pathchar [7] are CPU intensive for sophisticated calculations and/or channel intensive due to a large amount of probing packets. Thus, overlapping their execution intervals over the same measurement server or the same channel can cause serious interference and lead to misleading reports of the network status. We define a *measurement conflict* as an execution overlap of multiple measurement jobs that results in misleading reports.

In addition to the measurement conflict issue, one additional constraint of the measurement-level scheduling problem is the *Measurement-Level Agreement* (MLA). Utilizing excessive network resources just for active measurements is not appropriate since it significantly degrades the regular user traffic performance. Thus, we need a regulation on the measurement traffic. Since an end-to-end measurement could involve analyzing data along network paths of multiple ISPs, we envision “measurement federations” in which many ISPs participate in interdomain measurements based on MLAs for reaping the mutual benefits of performing end-to-end path measurements. MLAs can specify that only a certain percentage (1 percent to 5 percent) or only a certain number of bits per second (1 or 2 Mbps) of the network bandwidth in ISP backbones could be used for measurement traffic, which can ensure that the actual application traffic is not seriously affected by measurement traffic.<sup>2</sup> We use the notation  $\psi$  to denote the MLA specification in an NMI. In the measurement-level scheduling problem, the sum of the bandwidth usage by concurrent measurement jobs over the same channel should be less than  $\psi$  at all times.

From the above inputs and constraints, the measurement-level scheduling problem can be formally described as follows.

2. Since most active measurement tools have options to specify the packet sizes and bandwidth usage of a measurement test, simple calculations can be used to determine how much of a network’s bandwidth will be used by a given set of active measurements over a certain period of time.

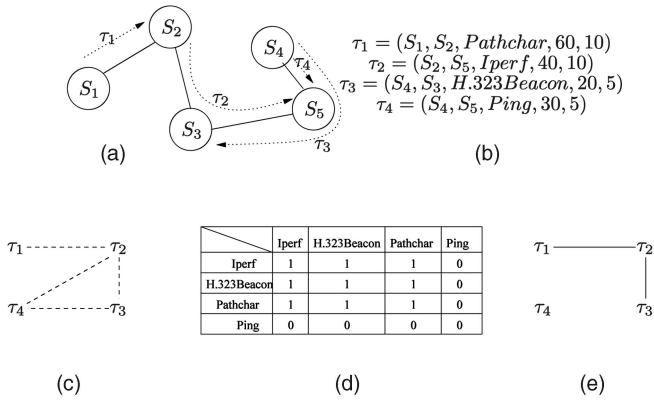


Fig. 3. Task conflict graph. (a) Measurement topology. (b) Task set. (c) Potential task conflict graph. (d) Tool conflict matrix. (e) Task conflict graph.

**Problem.** Given measurement topology  $G = (N, E)$  and offline specified measurement task set  $\Gamma = \{\tau_1, \tau_2, \dots, \tau_n\}$ , find the schedule of measurement jobs such that all deadlines (equal to periods) can be met while preventing conflicts and adhering to the MLA constraint  $\psi$ . For an on-demand measurement request  $J_{k,r}$  schedule it as early as possible without violating the deadlines of offline tasks in  $\Gamma$ , conflict constraint, and MLA constraint.

## 4 MEASUREMENT SCHEDULING ALGORITHMS

In this section, we first present an offline scheduling algorithm to construct a schedule table for a given set of periodic measurement tasks  $\Gamma = \{\tau_1, \tau_2, \dots, \tau_n\}$ . Then, we present an online algorithm to schedule an on-demand measurement request  $J_{k,r}$  without missing deadlines of periodic tasks. We first assume the existence of a central regulator that governs the global schedule and later relax this assumption.

### 4.1 Offline Scheduling Algorithm

In our measurement scheduling framework, a central regulator collects all specifications of periodic measurement tasks and builds a schedule table that determines the times when measurement jobs can start and stop at each server. To build such a table, the first step is to make a *task conflict graph* by combining the measurement topology  $G$  and the task set  $\Gamma$ . Fig. 3 shows an example problem. For the given measurement topology and the task set in Figs. 3a and 3b, we examine each pair of tasks  $\tau_i$  and  $\tau_j$  to see if they share the same source server, destination server, or part of the paths between the source and destination servers. If so, the two tasks may “potentially” conflict if scheduled concurrently. In Figs. 3a and 3b,  $\tau_1$  and  $\tau_2$  share  $S_2$  and, thus, we add a potential dependency edge between them in the potential task conflict graph, as in Fig. 3c.  $\tau_2$  and  $\tau_3$  share the path and, thus, a dependency edge is added. On the other hand,  $\tau_1$  does not share any network resource with  $\tau_3$  and, thus, no edge is added. Even if two tasks share network resources, they may not actually conflict, depending on the active measurement tools used. Based on our empirical studies in [15], we could determine which tools conflict if they run concurrently. The result is summarized by the tool conflict matrix in Fig. 3d. For example, Iperf and Pathchar

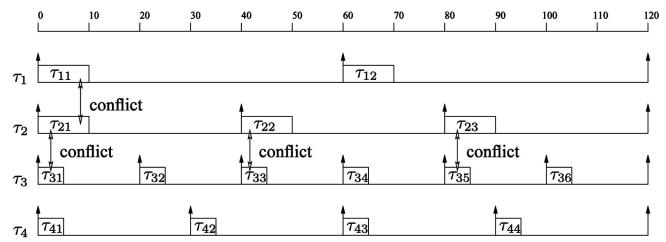


Fig. 4. No orchestrated schedule.

conflict if they run concurrently on the same server since both intensively use the server and channel resources for active measurement. On the other hand, Ping just injects small probing packets and, hence, does not conflict with any other tools. Considering the tool conflict matrix, the potential task conflict graph in Fig. 3c can be converted to the final task conflict graph in Fig. 3e. The edge between two tasks in the task conflict graph means that they should be scheduled in a mutually exclusive manner; otherwise, a conflict happens, resulting in misleading reports.

Now, we can consider only the final task conflict graph to compute the offline schedule. One obvious solution is to start a measurement job at the source server at its release time, without considering the measurement conflict and MLA constraints. Fig. 4 shows such a schedule for the problem given in Fig. 3. In the figure, the upward arrows indicate the release times of the periodic measurement tasks. The schedule, however, causes a number of conflicts that result in a misleading report of the actual network performance. Another approach is to run only a single measurement job at any time instant by using a non-preemptive EDF scheduling algorithm. Fig. 5 shows such a schedule for the same problem. It can completely prevent conflicts. However, it does not allow concurrent execution of multiple jobs, even if they do not conflict, which degrades the schedulability.

We aim to find a schedule in between these two extremes such that conflicts are completely prevented while maximizing the concurrent execution whenever possible. For this, we propose the EDF with Concurrent Execution (EDF-CE) algorithm that schedules measurement jobs in the EDF order while allowing concurrent execution if jobs do not conflict. The algorithm is formally described as follows:

**EDF-CE:** for the given task conflict graph, find the measurement schedule during a hyperperiod

**Input:** task set  $\Gamma$  and task conflict graph

**Output:** start time  $st_{ij}$  and finish time  $ft_{ij}$  for each job  $\tau_{ij}$  in a hyperperiod

**begin procedure**

1. Initialize  $rt\_list$  with the ordered list of all release times in a hyperperiod
2. Initialize  $ft\_list = \{\}$  /\* ordered list of finish times\*/
3. Initialize  $pending\_job\_queue = \{\}$
4. **do**
5.  $time =$  get the next scheduling time point from  $rt\_list$  and  $ft\_list$
6. add all newly released jobs at  $time$  to  $pending\_job\_queue$  in EDF order
7. **for** each job  $\tau_{ij}$  in  $pending\_job\_queue$  in EDF order

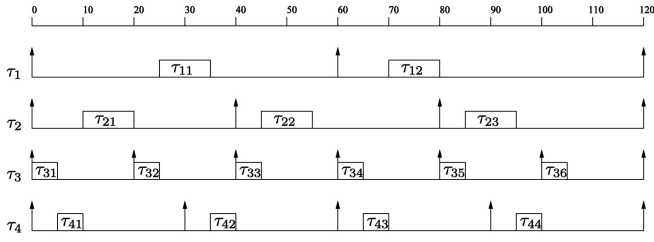
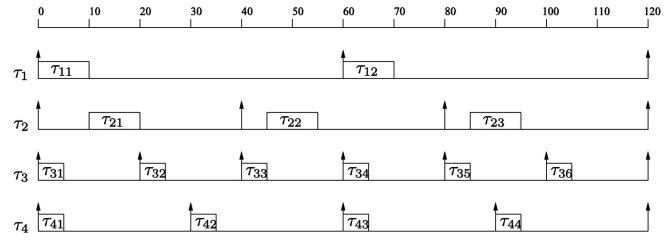


Fig. 5. Orchestration based on a single processor nonpreemptive EDF schedule.

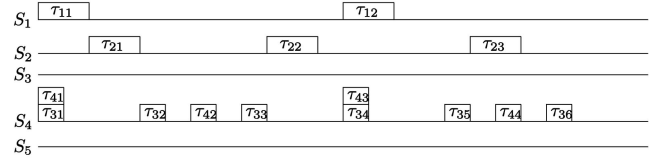
```

8.   if  $\tau_{ij}$  does not conflict with any of already
scheduled jobs at time and
9.   scheduling  $\tau_{ij}$  at time does not violate MLA
constraint  $\psi$ 
10.   $st_{ij} = \text{time}$  and  $ft_{ij} = \text{time} + e_i$ 
11.  if  $ft_{ij}$  is later than the deadline of  $\tau_{ij}$ 
12.  return error /* infeasible task set */
13.  end if
14.  remove  $\tau_{ij}$  from pending_job_queue
15.  add  $ft_{ij}$  to ft_list in order
16.  else
17.  do nothing /*  $\tau_{ij}$  will be considered again at the next
scheduling time point in the outer loop */
18.  end if
19.  end for
20. until time == hyperperiod
end procedure
    
```

The EDF-CE algorithm maintains the ordered list of release times *rt\_list* and the ordered list of finish times *ft\_list*. Line 1 initializes *rt\_list* with all release times in a hyperperiod. In Fig. 6, the release times are 0, 20, 30, 40, 60, 80, 90, 100, and 120. Line 2 initializes *ft\_list* as empty since no job is scheduled yet. Note that the only time points when we need to make a scheduling decision are either when a new job is released or a current executing job is finished. Thus, we call times in *rt\_list* and *ft\_list* “scheduling time points.” In addition, the algorithm maintains a *pending\_job\_queue* that holds all jobs released but not scheduled in the EDF order. Line 3 initializes it as empty. The **do-until** loop from lines 4 to 20 progresses the virtual time variable *time* up to a hyperperiod while determining the schedule at all scheduling time points. Line 5 moves *time* to the next scheduling time point. Then, line 6 adds all newly released jobs to the *pending\_job\_queue*. The **for** loop from lines 7 to 19 examines the pending jobs in the EDF order and determines whether they can start at *time* without causing any conflict and without violating MLA  $\psi$  (see lines 8 and 9). If so, job  $\tau_{ij}$ ’s start time  $st_{ij}$  is determined as *time* and its finish time  $ft_{ij}$  is determined as  $\text{time} + e_i$  in line 10. If the finish time  $ft_{ij}$  is later than the deadline of job  $\tau_{ij}$  in line 11, we cannot construct a feasible schedule that meets all deadlines and, hence, return an error in line 12. If we can meet the deadline of  $\tau_{ij}$ , we can continue. In line 14,  $\tau_{ij}$  is removed from the *pending\_job\_queue*. Also, its finish time  $ft_{ij}$  is added to *ft\_list* so that  $ft_{ij}$  can be considered as a new scheduling time point in the outer **do-until** loop. If  $\tau_{ij}$  cannot be scheduled at *time* (line 16), it is kept in the *pending\_job\_queue* and can be considered again at the



(a)



(b)

Fig. 6. EDF-CE schedule. (a) EDF-CE schedule. (b) Schedule table for each server.

next scheduling time point by the outer loop. Note that the algorithm tries to concurrently start as many jobs as possible in the EDF order at *time* as long as they neither conflict with nor violate the MLA. Fig. 6a shows such an EDF-CE schedule for the same problem of Fig. 3. At time 0 of the EDF-CE schedule, note that  $\tau_{11}$ ,  $\tau_{31}$ , and  $\tau_{41}$  are executed concurrently but not  $\tau_{41}$ , which maximizes the concurrent execution, guaranteeing no conflict.

Once we find the EDF-CE schedule, we can convert it to the measurement schedule table of each server, considering the source server of each job. Fig. 6b shows the schedule tables of all five servers. Such constructed schedule tables are transferred to corresponding servers so that they can start and stop the planned measurement jobs.

## 4.2 Online Scheduling of On-Demand Measurement Requests

At runtime, while each server executes periodic measurement tasks according to the precomputed schedule table, a network engineer can request an on-demand measurement  $J_k$ . For now, we assume that such a request is received by the central regulator.

Upon the arrival of an on-demand request  $J_k = (\text{src}_k, \text{dst}_k, \text{tool}_k, e_k)$ , our goal is to serve it as early as possible, without missing any deadlines of periodic measurement tasks. For this, we propose a *recursive push* algorithm that recursively pushes offline scheduled periodic jobs within their deadlines. This push can create a leftover time, called a *slack*, as early as possible and this slack time can be used to schedule  $J_k$ . The basic idea of recursive push can be best illustrated in Fig. 7, which shows the same EDF-CE schedule as above. Suppose that an on-demand request  $J_k = (S_2, S_3, \text{Iperf}, 10)$  arrives at time 50. We assume that  $J_k$  conflicts with  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$ , as shown by the modified task graph. The central regulator cannot allow  $J_k$  to start at its arrival time 50 since it conflicts with  $\tau_{22}$ . Thus, the central regulator calculates the maximum slack from starting at 55. For this, the central regulator calls

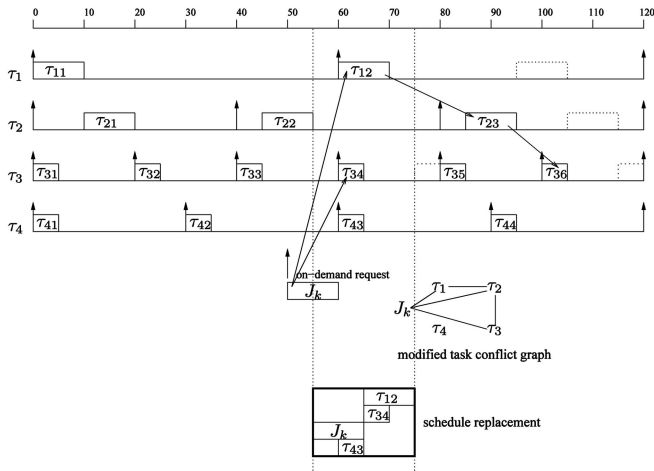


Fig. 7. Recursive pushing for maximum slack calculation.

**push**( $\tau_{12}$ ) and **push**( $\tau_{34}$ ) to determine how much  $\tau_{12}$  and  $\tau_{34}$  can be pushed to make the maximum slack for  $J_k$ . The **push** operation is recursive. To determine the maximum **push** of  $\tau_{12}$ , we first have to know the maximum **push** of the dependent job  $\tau_{23}$ . Thus, **push**( $\tau_{12}$ ) recursively calls **push**( $\tau_{23}$ ) and, in turn, **push**( $\tau_{23}$ ) calls **push**( $\tau_{36}$ ). On the other hand,  $\tau_{36}$  does not conflict with any other offline scheduled jobs while being pushed up to its deadline  $d_{36} = 120$ . Such a job with which the recursion can terminate is called a *terminal job*. Similarly,  $\tau_{34}$  is also a terminal job. For a terminal job  $\tau_{ij}$ , the *push* procedure can determine its new pushed finish time  $new\_ft_{ij}$  and new pushed start time  $new\_st_{ij} = new\_ft_{ij} - e_i$ , without any further recursive calls. The **push** operation is formally defined as follows:

**push**: return the new start time of input jobs after maximum push

**Input**:  $\tau_{ij}$

**Output**: new start time after maximum push  $new\_st_{ij}$

**begin procedure**

1. if  $\tau_{ij}$  has no conflicting jobs scheduled up to  $d_{ij}$  /\* terminal job \*/
2. slide  $\tau_{ij}$  from  $st_{ij}$  to  $d_{ij} - e_i$  until MLA violation is observed at  $t_{MLA}$  ( $t_{MLA} < d_{ij}$ ).
3. if  $t_{MLA}$  is found, the new finish time  $new\_ft_{ij} = t_{MLA}$ . otherwise,  $new\_ft_{ij} = d_{ij}$ .
4. new start time  $new\_st_{ij} = new\_ft_{ij} - e_i$ .
5. **else** /\* not a terminal job \*/
6. new finish time  $new\_ft_{ij} = d_{ij}$ .
7. **for** each conflicting task  $\tau_{i'j'}$  up to  $d_{ij}$
8.  $new\_ft_{ij} = \min(new\_ft_{ij}, \mathbf{push}(\tau_{i'j'}))$ .
9. **end for**
10. slide  $\tau_{ij}$  from  $st_{ij}$  to  $new\_ft_{ij} - e_i$  until MLA violation is observed at  $t_{MLA}$  ( $t_{MLA} < new\_ft_{ij}$ ).
11. if  $t_{MLA}$  is found, the new finish time  $new\_ft_{ij} = t_{MLA}$ . otherwise keep  $new\_ft_{ij}$ .
12. new start time  $new\_st_{ij} = new\_ft_{ij} - e_i$ .
13. **end if**
14. return  $new\_st_{ij}$ .

**end procedure**

This algorithm returns the new start time  $new\_st_{ij}$  after maximally pushing  $\tau_{ij}$ . If  $\tau_{ij}$  is a terminal job, its new finish time can be pushed up to its deadline  $d_{ij}$  if we could ignore the MLA constraint. In order to consider the MLA constraint, in line 2, we slide  $\tau_{ij}$ 's execution interval up to  $d_{ij}$  to find the earliest time point  $t_{MLA}$  when the MLA constraint can be violated, if any. If such time point  $t_{MLA}$  is found,  $t_{MLA}$  is the latest possible pushed finish time of  $\tau_{ij}$  without violating the MLA constraint. Thus,  $new\_ft_{ij}$  is set to  $t_{MLA}$  in line 3. Otherwise, the new finish time can be pushed up to  $d_{ij}$ , that is,  $new\_ft_{ij} = d_{ij}$  in line 3. Once the new finish time is determined, line 4 can simply calculate the new start time, that is,  $new\_st_{ij} = new\_ft_{ij} - e_i$ .

If  $\tau_{ij}$  is not a terminal job, lines 7, 8, and 9 recursively call **push** for all dependent jobs to figure out the minimum new start time of all dependent jobs. If we ignore the MLA constraint, the minimum of the deadline  $d_{ij}$  and the new pushed start times of all dependent jobs is the latest possible new finish time  $new\_ft_{ij}$  for  $\tau_{ij}$ . Lines 10 and 11 can advance the new finish time  $new\_ft_{ij}$ , considering the MLA constraint in the same way as in the terminal job case. With  $new\_ft_{ij}$ , line 12 calculates the new start time as  $new\_st_{ij} = new\_ft_{ij} - e_i$ . Finally, line 14 returns  $new\_st_{ij}$ .

Considering the  $new\_st_{ij}$  of all dependent jobs of  $J_k$ , we can calculate the maximum slack that can be used for the on-demand request  $J_k$  starting from the current scheduling time point  $t$ . If the maximum slack is larger than the required execution time  $e_k$  and also if executing  $J_k$  from  $t$  to  $t + e_k$  does not violate the MLA constraint, the central regulator sets time  $t$  as the start time of  $J_k$  and pushes dependent periodic jobs as needed. The piece of schedule affected by  $J_k$  (see "schedule replacement" in Fig. 7) is transferred to the corresponding servers so that they can temporarily use the updated schedule piece instead of the original schedule to accommodate  $J_k$ . If the above condition does not hold, the central regulator examines the next scheduling time point to recalculate the maximum slack and so on until it finds enough slack time during which  $J_k$  can be executed without violating the MLA constraint.

### 4.3 Distributed Implementation of Scheduling Algorithms

The aforementioned scheduling algorithms assume a central regulator that collects all offline/online measurement requests and builds/updates the global schedule. A centralized regulator is popular in NMIs because it is convenient to initiate and collect measurements into a central database. However, such a centralized mechanism could incapacitate an NMI when there is a failure of the central regulator. Also, some applications require distributed measurement scheduling to gain greater flexibility to dynamically determine the locations of measurement data collection and subsequent analysis. To address these issues, this section presents a mechanism to implement the above scheduling algorithms in a decentralized way.

In a distributed setting, measurement requests (for example, add/remove periodic measurement tasks and on-demand measurement jobs) arrive at their local servers, possibly concurrently. If each server concurrently updates the schedule upon the arrival of requests, it breaks the

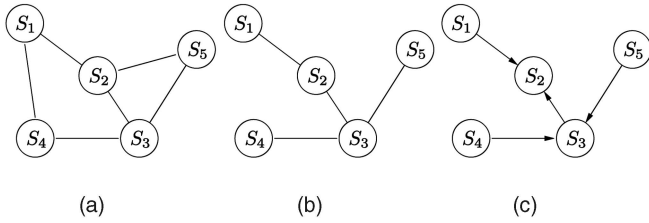


Fig. 8. Minimal spanning tree for the measurement server topology. (a) Measurement topology. (b) Minimal spanning tree. (c) Initial lock placement.

consistency of the schedule and, in turn, creates measurement conflicts. Therefore, the issue is to serialize the distributed concurrent requests such that the schedule can be updated in a consistent way. For this, we propose using Raymond’s algorithm [16], developed for distributed synchronization. This section describes how Raymond’s algorithm works with our scheduling algorithms, maintaining the schedule consistency in a distributed way.

For the measurement topology given in Fig. 8a as an example, we first create the minimal spanning tree as in Fig. 8b. This tree is used to maintain a tree-wide single lock, with minimal exchange of messages [16]. The basic idea is to allow only the lock holder to commit the arrival of a request at a time, which assures the global serialization of concurrent requests. In the initialization phase, we place the lock at any server, say,  $S_2$  in the example in Fig. 8b, and make each server set its *dir* variable to the neighbor toward the lock holder, as shown in Fig. 8c.

Upon the arrival of a new request at a server, the server exchanges messages with others along the spanning tree and eventually gets the lock. Then, it commits the arrival of the request by sending this commitment information to all of the affected servers. All of the servers that receive this commitment run the same EDF-CE (for an add/remove request of a periodic task) or the recursive push algorithm (for an on-demand job) to update its schedule table. This procedure can be best illustrated by the example in Fig. 9. Suppose that the initial lock holder is  $S_2$ , as shown in the leftmost tree. Also, assume that an on-demand job  $J_1(S_5, S_4, Iperf, 10)$  arrives at  $S_5$  at time  $t_1$ . Since  $S_5$  is not the lock holder, it enqueues its ID  $S_5$  in  $S_5$ ’s queue and sends a LOCK-REQUEST message to the neighbor  $S_3$  pointed to by its *dir* variable.  $S_3$  is not the lock holder, either and, thus, it enqueues the requester’s ID  $S_5$  and sends a LOCK-REQUEST message to the neighbor  $S_2$  pointed to by its *dir* variable. In the meantime, suppose that another request  $J_2(S_1, S_3, Iperf, 10)$  arrives at  $S_1$  at time  $t_2$ . Since  $S_1$  is not the lock holder, it enqueues its ID  $S_1$  and sends a LOCK-REQUEST message to  $S_2$  pointed to by its *dir* variable. When  $S_1$ ’s LOCK-REQUEST reaches the lock holder  $S_2$ ,  $S_2$ ’s queue is empty and it is not updating the schedule (not in the critical section) and, thus, it can immediately yield its lock to the requester  $S_2$  by sending a LOCK-APPROVAL message to the requester  $S_1$ . It is no longer the lock holder and sets its *dir* variable to  $S_1$  (see the second tree). When  $S_1$  receives the LOCK APPROVAL at  $t_3$ , it notices that the head of the queue is its own ID and thus can enter the critical section to commit  $J_2$ ’s arrival. Since the  $J_2$ ’s arrival needs to be viewed by all affected servers in a

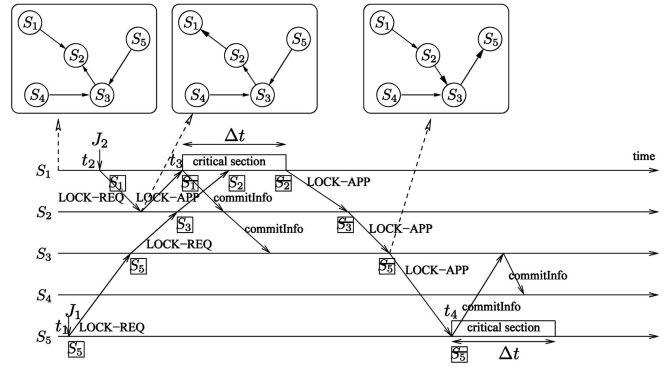


Fig. 9. Distributed schedule update.

consistent way,  $S_1$  adds the sufficient delay  $\Delta t$  of commitment transmission to  $t_3$  and considers  $t_3 + \Delta t$  as the committed arrival time of  $J_2$ . Then,  $S_1$  sends the commitment information ( $J_2$  and  $t_3 + \Delta t$ ) to all of the affected servers  $S_2$  and  $S_3$ . Now,  $S_1$ ,  $S_2$ , and  $S_3$  can run the same recursive push algorithm for inserting  $J_2$  with the same committed arrival time of  $t_3 + \Delta t$ . When the LOCK-REQUEST message from  $S_3$  arrives at  $S_2$ ,  $S_2$  is not the lock holder and its *dir* is pointing to  $S_1$ . Thus, the LOCK REQUEST is forwarded to  $S_1$ . When the LOCK REQUEST reaches  $S_1$ , it is the lock holder, but it is already in the critical section to commit  $J_2$ . Thus,  $S_1$  enqueues the requester’s ID  $S_2$ . After that,  $S_1$  leaves the critical section at  $t_3 + \Delta t$ . At this time,  $S_1$  notices that the head of its local queue is  $S_2$  and thus sends a LOCK-APPROVAL message to  $S_2$  and sets its *dir* toward  $S_2$ .  $S_2$  and  $S_3$ , in turn, forward the LOCK APPROVAL and update their *dir* variables according to the head of their queues until the LOCK APPROVAL reaches  $S_5$ . When  $S_5$  receives the LOCK APPROVAL at time  $t_4$ , it notices that the head of its queue is itself and thus can enter the critical section to commit the arrival of  $J_1$ . The commitment phase is the same as that of  $J_2$ . As a consequence, the concurrent arrivals of  $J_1$  and  $J_2$  are globally serialized in the order of  $J_2$  and  $J_1$ , with the consistent commitment times of  $t_3 + \Delta t$  and  $t_4 + \Delta t$ . Therefore, the schedule can be updated in a globally consistent way, assuring the conflict-free scheduling property.

For the complete and formal description of this distributed schedule update procedure, readers are referred to [16]. The procedure inherits the proven properties of Raymond’s algorithm, such as minimal message exchange for assuring serializability, deadlock freedom, no starvation, and fault tolerance.

#### 4.4 Measurement Federation Issues across ISP Borders

Collecting measurement data within a single ISP domain is not sufficient for distributed computing applications because they often span network paths across multiple ISP domains. For example, application service providers such as Vonage rely on multiple ISPs for delivering worldwide voice-over-IP (VoIP) and videoconferencing services. To serve their customers and meet the service-level agreements (SLAs), ISPs need to support interdomain measurements that could produce end-to-end Internet measurements. For facilitating such interdomain measurements, “NMI federations” [17], [5]

have emerged where multiple ISPs agree upon a common measurement policy to cooperate with each other.

This section discusses the interdomain NMI federation issues and explains how our scheduling framework can be incorporated into the federation. For building an NMI federation, all of the participating ISPs should agree on the following:

1. sharing each other's measurement server topology,
2. bounding the amount of measurement traffic (that is, the MLA constraint  $\psi$ ),
3. authenticated and secure access to measurement resources, and
4. sharing collected measurement data.

First, the measurement server topology of an ISP can be securely revealed only to other ISPs in the same federation by using the agreed upon authentication and encryption methods, as will be discussed later. Thus, every measurement server in the NMI federation can have the federation-wide view of the server topology and, thus, can determine the schedule of measurement tasks, even if they span across multiple ISPs. Second, the agreed measurement traffic bound, that is, the MLA constraint  $\psi$ , can be enforced in our scheduling algorithms, as explained in Sections 4.1 and 4.2, and, thus, it can be complied with across multiple ISPs. Third, for authenticated and secure access to measurement resources across ISP borders, all ISPs can use preagreed upon authentication and encryption techniques. For example, upon the arrival of a new measurement request, they can use a centralized Kerberos [18] authentication server with a Data Encryption Standard (DES) or triple DESs. This can verify that the requesting domain belongs to the same NMI federation and also prevents intruders from eavesdropping on the request for deciphering the authentication mechanism and impersonating a member of the NMI federation. Finally, the collected measurement data can be shared by multiple ISPs as needed by distributed computing applications by using "Request/Response" schemas being developed by the Global Grid Forum [17].

We envision that the growth of an NMI federation mostly involves political hurdles rather than technical ones. Since the application and ISP communities are realizing the importance of the NMI federation for interdomain distributed computing, we believe that all of the political hurdles will be overcome, resulting in a worldwide NMI federation. Note that such efforts have already been started by communities such as the Global Grid Forum, Internet2 in the US, and DANTE in Europe [17], [5].

## 5 CASE STUDY WITH NWS FOR DISTRIBUTED COMPUTING

### 5.1 NWS Network Prediction

We now apply our measurement scheduling algorithms to the widely used NWS [9] that can provide network performance forecasts.<sup>3</sup> In this application, one challenge is the gap between the original measurement time requirement of NWS and the actual temporal behavior of our

3. NWS uses periodically measured network status and forecasts the network performance. Due to its ability to forecast the network performance, NWS has been adopted by a number of networked job schedulers such as AppLes [19], Legion [20], and Globus/Nexus [21].

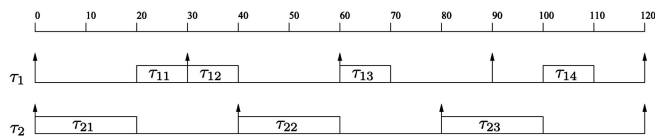


Fig. 10. Jitter of intersampling time points of network status.

scheduling algorithms. More specifically, NWS periodically issues measurement requests, expecting a periodic sampling of network status. However, the scheduler cannot serve the requests at exactly the desired times due to resource conflicts with other measurement requests. As such, any scheduler that tries to avoid conflicts inevitably creates a jitter in the intersampling times of the network status. This section presents a simple method for compensating for the intersampling jitter.

NWS relies on a continuous and periodic sampling or *Pure Periodic Sampling* (PPS) of network status. It uses the periodically sampled network status data to maintain the history of network performance, which, in turn, is used to generate ongoing and dynamic network performance forecasts. The forecast time window is the same as the sampling period.

However, it is not always possible for the measurement scheduling algorithm to provide pure periodic network status data, especially when multiple measurement tasks are running. This can be explained by Fig. 10, which shows an example conflict-free schedule of two periodic measurement tasks that have a conflict relation. We can note that the task  $\tau_2$  is scheduled purely periodically with constant intersampling times. However, the intersampling times of task  $\tau_1$  vary for every instance. To avoid conflict of multiple concurrent tasks, the actual scheduling time points inevitably deviate from the periodic release time points by any conflict-free scheduling algorithm. Our EDF-CE also produces such intersampling time jitter since it is designed to guarantee the periodic deadlines and not pure periodic execution of jobs. In fact, with our EDF-CE, the intersampling jitter of task  $\tau_i$  can vary from  $e_i$  (when a job instance is scheduled just before its deadline and the next one is scheduled at the release time) to  $2p_i - e_i$  (when a job instance is scheduled at the release time and the next one is scheduled just before the deadline).

Although our EDF-CE causes intersampling jitter between two consecutive jobs, it bounds the jitter by meeting the end-of-period deadlines. Therefore, it can still be used for NWS with simple interpolations of the collected network status data. The interpolation transforms the actual measured data to pure periodic data by using piecewise linear interpolation. To explain this, let us consider Fig. 11. In the figure,  $(t_{i-1}, y_{i-1})$ ,  $(t_i, y_i)$ , and  $(t_{i+1}, y_{i+1})$  show the sequence of the  $(i-1)$ th,  $i$ th, and  $(i+1)$ th *actual periodic sampling* (APS), whose intersampling time is not always the same as the period  $p$ . To transform the APS sequence to a PPS sequence, which we call *Transformed Periodic Sampling* (TPS), we can draw piecewise lines between pairs of two APS points. For example, we can draw a line between  $(t_{i-1}, y_{i-1})$  and  $(t_i, y_i)$ , as shown in Fig. 11. With this line, we can estimate the measurement value  $\hat{y}_i$  at the pure periodic



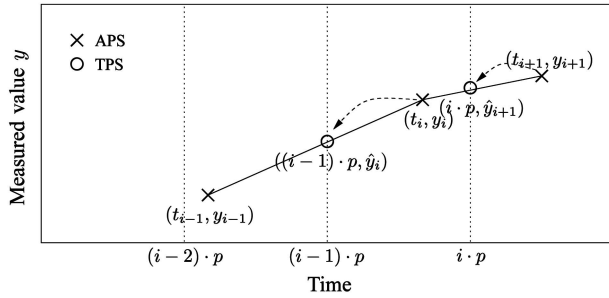


Fig. 11. APS transformation to TPS.

$i$ th sampling time, that is,  $(i-1) \cdot p$ . Specifically,  $\hat{y}_i$  is given as follows:

$$\hat{y}_i = y_i + \frac{y_{i+1} - y_i}{t_{i+1} - t_i} ((i-1) \cdot p - t_i).$$

Thus, the APS data  $(t_i, y_i)$  with intersampling jitter can be transformed into the pure periodic TPS data  $((i-1) \cdot p, \hat{y}_i)$ , as shown in Fig. 11. Similarly,  $(t_{i+1}, y_{i+1})$  can be transformed into  $(i \cdot p, \hat{y}_{i+1})$ . Now, the NWS can use the TPS data, rather than the original measured data, to provide the network performance forecast. With this simple interpolation method, we will show in Section 6.3 that our EDF-CE can work well with NWS to produce accurate forecasts.

## 5.2 Use of NWS for Distributed Computing

Due to its ability to forecast the network status, NWS can be used for a number of distributed computing applications that rely on networked computational resources. In this section, we sketch the scenarios where the NWS based on our conflict-free measurements can help two typical examples of distributed computing, that is, Grid computing and videoconferencing.

In Grid computing, users often transfer computational jobs involving large data sets (sometimes even on the scale of terabytes) to remote computing sites [22]. If proper network paths are not selected, then jobs that traverse problematic paths could hold up the speedy completion of other queued jobs that traverse problem-free paths and thus significantly degrade the overall efficiency of grid computing. This problem can be avoided by using the NWS. It can accurately monitor and predict the network status by using the conflict-free (and, so, not misleading) measurement data. The job scheduler of grid computing can use such network status data to select the computing sites and network paths in such a way that ensures the optimal efficiency of the overall computing. In addition, the ability of network status prediction can allow the job scheduler to dynamically change the network path selections before a severe performance degradation happens.

In videoconferencing, interactive sessions involving three or more participants are established using call-admission controllers that manage Multipoint Control Units (MCUs). MCUs combine the admitted voice and video streams from participants and generate a single conference stream that is multicast to all of the participants. If a call-admission controller selects problematic network paths between the participants and MCUs, the perceptual quality of the conference stream could be seriously affected by

impairments such as video frame freezing, audio dropouts, and even call-disconnects. Using the NWS, such a problem can be avoided. The call-admission controllers can consult the NWS to figure out the network paths that can satisfy the application QoS requirements. In addition, the network status forecasts from NWS can also be used to monitor whether the current selection will experience problems due to the paths that may soon degrade the application QoS severely. In such cases, the call-admission controllers can dynamically change to alternate network paths that are identified to satisfy QoS requirements for the next forecasting period.

The selected paths in cases of both grid computing and videoconferencing can be enforced in the Internet by using MPLS explicit routing or by exploiting path diversity based on multihoming or overlay networks [23], [24].

## 6 EXPERIMENTAL RESULTS

In this section, we evaluate the performance of our measurement scheduling algorithms. We first perform simulations with synthetic measurement tasks to show the maximum schedulability by the EDF-CE algorithm and the average response times of on-demand requests by the recursive push algorithm. Then, we present performance evaluation results on an actual Internet2 testbed. Finally, we present the case study results of applying our EDF-CE to NWS.

### 6.1 Performance Evaluation Results Using Synthetic Tasks

Our synthetic task set is comprised of four periodic active measurement tasks:  $\tau_1$ ,  $\tau_2$ ,  $\tau_3$ , and  $\tau_4$ . The period  $p_i$  of each task  $\tau_i$  is randomly generated from [1,000 sec, 10,000 sec]. The execution time  $e_i$  of each task  $\tau_i$  is randomly generated from [100 sec, 999 sec]. Since the measurement topology and intertask conflict relations can be represented by a task conflict graph, we conduct this experiment as changing only the task conflict graph. The task conflict graph of the four tasks is randomly created using a parameter called a *conflict factor*. The conflict factor represents the probability that there is a conflict edge between any two tasks. Therefore, when the conflict factor is 1, the task conflict graph is fully connected. If the conflict factor is 0, there is no edge between tasks.

For each sample task set and task conflict graph, we use the “maximum schedulable utilization”  $\sum_{i=1}^4 e_i/p_i$  as the performance metric. We determine the maximum schedulable utilization by gradually increasing execution times  $e_i$  until the scheduling algorithms fail to construct a feasible schedule.

We compare three scheduling algorithms:

- **No Orchestration.** This schedules measurement jobs at their release times, without considering measurement conflicts.
- **EDF.** This schedules only one measurement job at a time by using the nonpreemptive EDF algorithm, just like a single processor EDF scheduling.
- **EDF-CE.** This is proposed in this paper.

Fig. 12 shows the maximum schedulable utilization while increasing the conflict factor. Here, we assume a large MLA  $\psi$ , say, 50 Mbps, and thus avoid any MLA

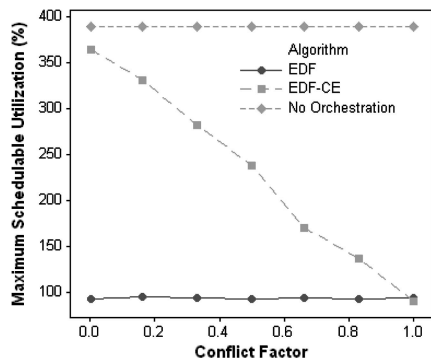


Fig. 12. Maximum schedulable utilization by three scheduling algorithms.

bottlenecks when finding the schedule. Each plotted point in the figure is the average of 1,000 random sample task sets. EDF's maximum schedulable utilization is constantly bounded under 100 percent, regardless of the conflict factor, since it does not allow concurrent execution, even when possible. On the other hand, our EDF-CE algorithm can maximally utilize the concurrent execution whenever possible. When the conflict factor is zero, EDF-CE allows concurrent execution of all four tasks. This is similar to scheduling the four tasks on four independent processors. Thus, the maximum schedulable utilization reaches up to 400 percent. As the conflict factor increases, the maximum schedulable utilization gradually decreases. When the conflict factor is 1, that is, when all four tasks conflict with each other, EDF-CE automatically degenerates to the single processor EDF and, hence, gives the maximum schedulable utilization of 100 percent. The result shows that EDF-CE is leveraging the "maximal but only possible" concurrent execution by explicitly considering the conflict dependency among tasks. The No-Orchestration approach always gives the maximum schedulable utilization of 400 percent since all four tasks can be concurrently executed, ignoring the conflict dependency. This, however, causes many conflicts, as will be shown in Section 6.2, resulting in many misleading reports of the actual network performance.

Fig. 13 illustrates how the maximum schedulable utilization of EDF-CE is bounded by the MLA constraint  $\psi$  and conflict factor. As expected, a higher value of  $\psi$  accommodates a larger number of concurrent jobs and, hence, produces a higher maximum schedulable utilization. For a given  $\psi$  value, the maximum schedulable utilization is constant up to a certain point of the conflict factor and then starts decreasing. Such a trend explains that  $\psi$  is the bottleneck when the conflict factor is small, whereas the conflict dependency becomes the bottleneck when the conflict factor is large.

To study the performance of the "recursive push" algorithm for handling on-demand measurement requests, we simulate random arrivals of on-demand jobs and schedule them over the offline EDF-CE schedule. The offline specified task set consists of four periodic tasks as before and their execution times and periods are randomly generated from [1 minute, 10 minutes] and [20 minutes, 200 minutes], respectively. The execution times and interarrival times of on-demand jobs are also randomly generated from

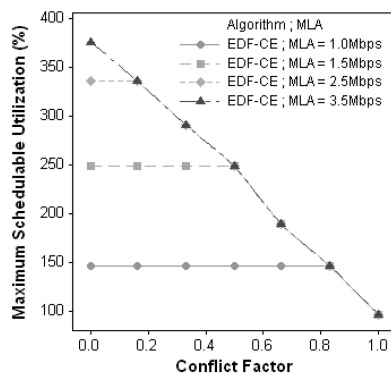


Fig. 13. Effect of MLA  $\psi$  and conflict factor to EDF-CE.

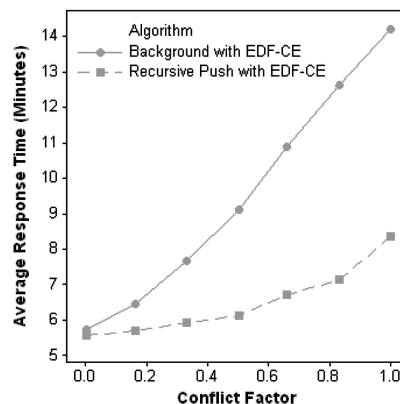


Fig. 14. Average response time of on-demand jobs.

[1 minute, 10 minutes] and [20 minutes, 200 minutes], respectively. The performance metric is the average of the response times for 1,000 on-demand jobs. We compare our recursive push algorithm with a background approach that schedules an on-demand job in the earliest gap present in the offline EDF-CE schedule, within which the on-demand job can execute to completion. Fig. 14 shows that our recursive push algorithm can significantly improve the responsiveness for on-demand measurement requests. Note that the average response time in both the background and recursive push cases increases as the conflict factor increases. This is because a higher conflict dependency among tasks reduces the concurrent execution of jobs and thus reduces the gaps available to schedule the on-demand jobs.

To estimate the overhead of online scheduling, we measure the algorithm's runtime for each on-demand job on a 2.4 GHz Pentium 4 Linux PC. Fig. 15 shows the average times as increasing the number of periodic tasks while fixing the conflict factor as 0.8. Even for a large number of periodic tasks with a high conflict factor, our recursive push algorithm can find the slack and calculate the updated schedule within tens of milliseconds. This is a negligible delay compared with typical measurement task execution times on the order of minutes.

In order to study the overhead of the distributed implementation of the scheduling algorithms, we simulate both centralized and distributed implementations of the recursive push algorithm in a large-scale network. For the network topology, we use the Waxman topology with 1,000 nodes produced by the BRITe tool [25]. From the

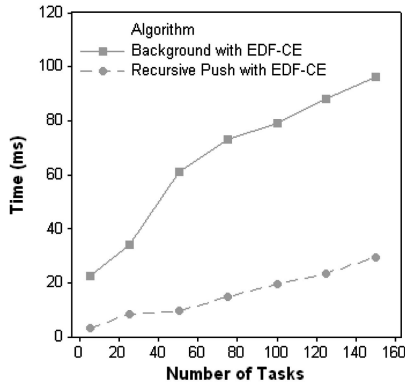


Fig. 15. Online schedule overhead for on-demand jobs.

topology of 1,000 nodes, we randomly select  $N$  nodes as the measurement servers creating a “measurement topology” with  $N$  measurement servers over the network topology with 1,000 nodes. We consider three different  $N$ s: 100, 200, and 300. These choices represent NMIs with a reasonably large number of servers, noting that the largest NMI deployment today, that is, the NLANR AMP project [3], has around 150 measurement servers deployed all over the world. On top of the measurement topology, we use a synthetic task set with 100 offline periodic measurement tasks. The period  $p_i$  and the execution time  $e_i$  of each task  $\tau_i$  are randomly generated from [20 minutes, 200 minutes] and [1 minute, 10 minutes], respectively. Then, the execution times of all 100 tasks are scaled such that the total utilization  $\sum_{i=1}^{100} e_i/p_i$  becomes 50 percent. Each task is assigned with randomly selected *src* and *dst* servers. With this offline periodic task set, we generate the offline schedule by using the EDF-CE algorithm. Given the offline schedule, we simulate the random arrival of 1,000 on-demand jobs with random execution times following the exponential distribution with the average 5 minutes. We conduct the simulation as we increase the average arrival rate from 10 jobs/hour to 150 jobs/hour following the Poisson distribution. Each on-demand job is assigned with randomly selected *src* and *dst* servers. In the following figures, we report the average of 100 simulation runs.

Fig. 16 compares the average response times of on-demand jobs by the centralized and distributed implementations of the recursive push algorithm. The centralized and distributed implementations show almost the same response times. This is because the total message passing delay to transfer the lock in the distributed implementation is at most 2 seconds, even in a large-scale measurement topology with 300 servers, as shown in Fig. 17. Also, such a delay does not increase with the increase of the arrival rate. This is due to the message minimization capability of Raymond’s algorithm as the number of requests increases [16]. Another interesting observation in Fig. 16 is that the response time is smaller when the number of servers is larger. This is because the on-demand job workload is scattered over a larger number of servers and, hence, the per-server workload is smaller.

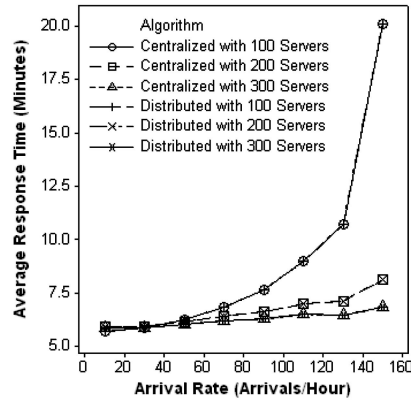


Fig. 16. Response time comparison of the centralized and distributed implementations.

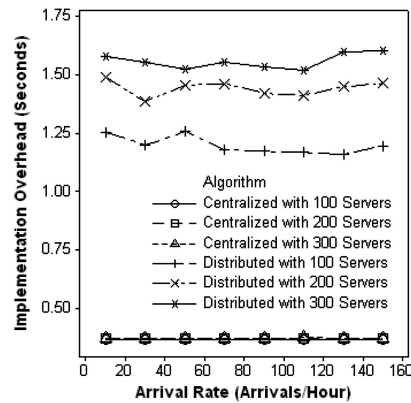


Fig. 17. Implementation overhead comparison of the centralized and distributed implementations.

### 6.2 Performance Evaluation Results on an Internet2 Testbed

We have actually implemented and deployed our scheduling algorithms in an NMI that is being used to monitor network paths on the Internet2 Abilene network backbone. The scheduling framework consists of a “Scripting Language Interface” and a central regulator, as shown in Fig. 18. The scripting language interface provides a generic and automated way to input measurement specifications such as measurement server topology, periodic measurement tasks, and MLAs. These specifications are interpreted by the central regulator to construct schedule timetables for the measurement servers. The constructed schedule timetables are transferred to the corresponding servers to initiate the measurement jobs at the planned times.

Our Internet2 testbed has five sites, each of which is equipped with a measurement server, as shown in Fig. 19a. To collect the actual measurement data, we run five periodic measurement tasks, as shown Fig. 19b. The resulting task conflict graph is shown in Fig. 19c.

Fig. 20 shows the H.323 Beacon mean opinion score (MOS) reports<sup>4</sup> measured between sites 3 and 4 by task  $\tau_2$ . To compare EDF-CE and No Orchestration, we pick the

4. MOS measurements reported by the H.323 Beacon are based on the E-Model [2], which is a computational model standardized by ITU-T to estimate the perceptual user quality for VoIP. The MOS values are reported on a quality scale of 1 to 5, with the [1, 3] range being poor, the [3, 4] range being acceptable, and the [4, 5] range being good. MOS values close to 4.41 are desirable for high-quality VoIP.

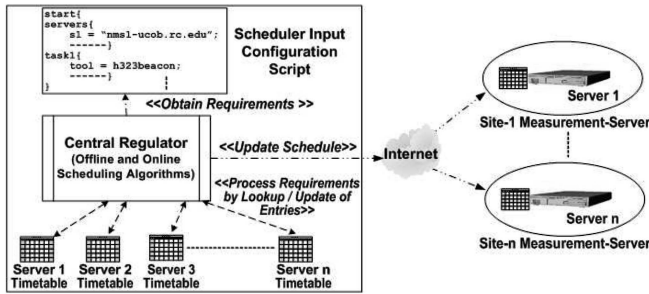


Fig. 18. Structure of the measurement scheduling framework.



Fig. 19. Internet2 testbed setup.

same 12-hour time frames in two consecutive days. For the 12-hour time frame on the first day, we use the No-Orchestration method to run all five measurement tasks, as shown in Fig. 19b, and collected the MOS reports from  $\tau_2$ . For the 12-hour time frame on the second day, we use EDF-CE and collect the same reports. From these two experiments, we can observe that the proposed EDF-CE guarantees zero conflict, whereas No Orchestration causes 50 percent instances of  $\tau_2$  to overlap with other tasks. All of the overlaps in the No-Orchestration schedule are indeed conflicts since all of the tools used in Fig. 19b are CPU intensive and channel intensive. In terms of MOS accuracy, however, we are not sure which curve better reflects the reality of the network status since we do not know the “true real” network status. In order to have a good representation of the reality of the network status between sites 3 and 4, we run only  $\tau_2$  over a week-long period. The results are shown in Fig. 21. Based on the figure, we can affirm that MOS fluctuation between 4.31 and 4.42 is natural in reality between sites 3 and 4. The MOS values in Fig. 20, which were collected using EDF-CE, match the representation of the reality in Fig. 21 well. In contrast, the MOS reports by the No-Orchestration method in Fig. 20 show a much larger fluctuation, which seems abnormal compared with Fig. 21. We can conclude that these abnormal fluctuations are due to the 50 percent instances of  $\tau_2$  conflicting with other tasks.

Although we do not present the data for the Iperf of  $\tau_1$ ,  $\tau_3$ , and  $\tau_5$  and the Pathload of  $\tau_4$  due to the page limit, we observed the similar measurement anomalies in No Orchestration but not in EDF-CE. From these observations, we can justify the importance of measurement orchestration for the correct estimation of the network status.

### 6.3 Case Study Results with NWS

In this section, we show how well our EDF-CE can work in combination with NWS. For this purpose, we use actual trace data obtained from NWS measurements for a path traversing a T1 connection with a total bandwidth 1.5 Mbps [26]. The trace data corresponds to “hourly” samples of available bandwidth on the T1 line over a two-day period. We assume that the trace data reflects the “actual” network

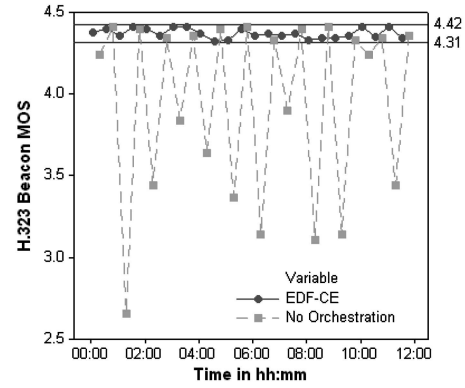


Fig. 20. H.323 Beacon MOS measurements between sites 3 and 4.

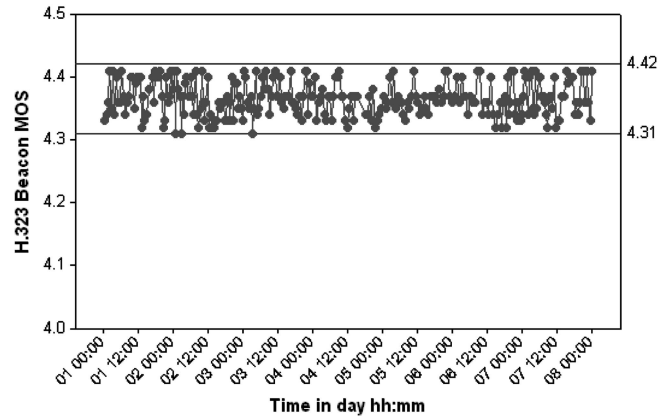


Fig. 21. H.323 Beacon MOS measurements over a week's period between sites 3 and 4.

performance trend on the path. Using this trace data, we generate two sample sequences: One representing the ideal PPS with a period of 2 hours and another one that corresponds to the APS by EDF-CE. The 2-hour-based PPS is obtained by considering every other sample in the trace data, assuming that no other monitoring tasks are present. This 2-hour-based PPS is the ideal measurement sequence expected by NWS for 2-hour look-ahead forecasts. The actual measurement sequence APS, however, inevitably deviates from the PPS due to scheduling of conflicting tasks. To model the APS, we simulate the EDF-CE with four measurement tasks in Fig. 22. Note that task  $\tau_1$  serves the NWS by providing 2-hour-based sampling of available bandwidth. This simulation provides the intersampling time distribution of  $\tau_1$ , which is used to select the samples from the trace data. The selected samples approximately represent the actual sequence of samples obtained by  $\tau_1$  as scheduled by EDF-CE in the presence of three other tasks.<sup>5</sup> Fig. 23 shows the reality of the available bandwidth (that is, 1-hour-based trace data), the NWS forecasted bandwidth using PPS, and the same using APS. The NWS forecasting using the ideal PPS closely matches the actual trend. However, the NWS forecasting using APS has nonnegligible differences from the reality. This is because of the

5. Due to the 1-hour-based granularity of the original trace data, the sequence of selected samples is only an approximation, with quantization errors of up to 1 hour. However, it is acceptable in terms of showing how the aforementioned simple interpolation can resolve the intersampling time jitter caused by EDF-CE, which ranges from 0 to 4 hours.

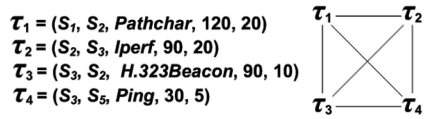


Fig. 22. Four task example for determining intersampling times of APS data.

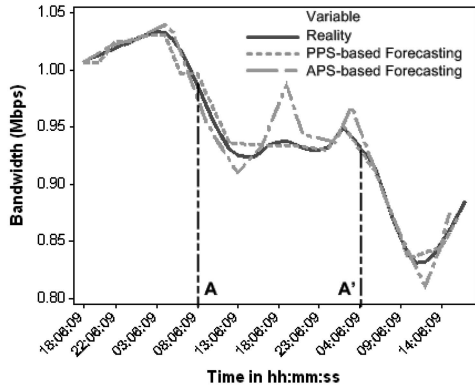


Fig. 23. Comparison of forecasts of PPS and APS.

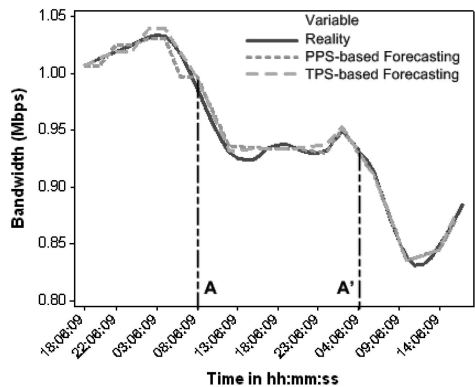


Fig. 24. Comparison of forecasts of PPS and TPS.

intersampling time jitter caused by EDF-CE. This problem can be fixed by a simple transformation of sampled data, as described in Section 5. Fig. 24 shows that the NWS forecasting using the transformed samples, denoted by *TPS*, can be very close to the ideal forecasting by *PPS*.

### 7 CONCLUSION AND FUTURE WORK

In this paper, we identify the measurement conflict problem, which results in misleading measurements of the network status when multiple conflicting measurement tools are executing at the same time on the same server or path. From the observation, we formulate the measurement scheduling problem as a real-time scheduling problem.

For the optimal schedulability of periodic measurement tasks, we use the EDF principle, which has been proven to be optimal in single processor preemptive scheduling and performs well in general settings. Our significant enhancement is to leverage the concurrent execution, which clearly differentiates the measurement scheduling problem from the classical real-time scheduling problems. Our enhanced EDF algorithm, called EDF-CE, allows concurrent execution of multiple measurement jobs not only on the isolated servers and paths but also on the same server and path as long as they do not conflict, that is, there are no misleading

reports. This significantly improves the schedulability and thus allows us to run measurements more frequently or to save significant time for on-demand requests.

We also propose an online scheduling algorithm to serve on-demand measurement requests as early as possible. The online algorithm can steal the maximum slack without violating any periodic deadlines and thus can almost immediately schedule the on-demand requests. Therefore, the response times of on-demand requests can be significantly reduced compared to their background processing.

Our proposed scheduling algorithms have actually been implemented and deployed on the Internet2 Abilene network. The actual experimental results demonstrate the pertinence and trustworthiness of our proposed scheduling algorithms.

### ACKNOWLEDGMENTS

This work has been supported in part by the Research Settlement Fund for the new faculty of Seoul National University (SNU) and by the Ohio Board of Regents. A preliminary version of this paper appeared in the *Proceedings of the IEEE Real-Time Systems Symposium (RTSS)*, 2005 [1]. Chang-Gun Lee is the corresponding authors for this paper.

### REFERENCES

- [1] P. Calyam, C.-G. Lee, P.K. Arava, and D. Krymskiy, "Enhanced EDF Scheduling Algorithms for Orchestrating Network-Wide Active Measurements," *Proc. 26th IEEE Int'l Real-Time Systems Symp. (RTSS '05)*, 2005.
- [2] P. Calyam, W. Mandrawa, M. Sridharan, A. Khan, and P. Schopis, "H.323 Beacon: An H.323 Application Related End-to-End Performance Troubleshooting Tool," *Proc. ACM SIGCOMM Workshop Network Troubleshooting (NetTs '04)*, 2004.
- [3] T. McGregor, H.-W. Braun, and J. Brown, "The NLNR Network Analysis Infrastructure," *IEEE Comm. Magazine*, 2000.
- [4] P. Calyam, D. Krymskiy, M. Sridharan, and P. Schopis, "TBI: End-to-End Network Performance Measurement Testbed for Empirical Bottleneck Detection," *Proc. First IEEE Int'l Conf. Testbeds and Research Infrastructures for the Development of Networks and Communities (TridentCom '05)*, 2005.
- [5] E. Boyd, J. Boote, S. Shalunov, and M. Zekauskas, "The Internet2 E2E piPES Project: An Interoperable Federation of Measurement Domains for Performance Debugging," Internet2 Technical Report, 2004.
- [6] A. Tirumala, L. Cottrell, and T. Dunigan, "Measuring End-To-End Bandwidth with Iperf Using Web100," *Proc. Fifth Passive and Active Measurement Workshop (PAM '03)*, 2003.
- [7] A. Downey, "Using Pathchar to Estimate Internet Link Characteristics," *Proc. ACM Ann. Conf. Applications, Technologies, Architectures, and Protocols for Computer Comm. (SIGCOMM)*, 1999.
- [8] C. Dovrolis, P. Ramanathan, and D. Moore, "Packet Dispersion Techniques and Capacity Estimation," *IEEE/ACM Trans. Networking*, 2004.
- [9] R. Wolski, N. Spring, and J. Hayes, "The Network Weather Service: A Distributed Resource Performance Forecasting Service for Metacomputing," *Future Generation Computer Systems*, 1999.
- [10] J. Liu, *Real-Time Systems*. Prentice Hall, 2000.
- [11] S. Kalidindi and M. Zekauskas, "Surveyor: An Infrastructure for Internet Performance Measurements," *Proc. Ninth INET Conf.*, 1999.
- [12] M. Alves, L. Corsello, D. Karrenberg, C. Ogut, M. Santcross, R. Sojka, H. Uijterwaak, and R. Wilhelm, "New Measurements with the RIPE NCC Test Traffic Measurements Setup," *Proc. Fourth Passive and Active Measurements Workshop (PAM '02)*, 2002.
- [13] B. Gaidioz, R. Wolski, and B. Tourancheau, "Synchronizing Network Probes to Avoid Measurement Intrusiveness with the Network Weather Service," *Proc. Ninth IEEE Int'l Symp. High-Performance Distributed Computing (HPDC '00)*, 2000.

- [14] H. Wang, D. Zhang, and K. Shin, "Change-Point Monitoring for Detection of DoS Attacks," *IEEE Trans. Dependable and Secure Computing*, vol. 1, no. 4, pp. 193-208, Oct.-Dec. 2004.
- [15] P. Callyam, C.-G. Lee, P.K. Arava, D. Krymskiy, and D. Lee, "OnTimeMeasure: A Scalable Framework for Scheduling Active Measurements," *Proc. Third IEEE/IFIP Workshop End-to-End Monitoring Techniques and Services (E2EMON '05)*, 2005.
- [16] K. Raymond, "A Tree-Based Algorithm for Distributed Mutual Exclusion," *ACM Trans. Computer Systems*, 1989.
- [17] GGF NMWG Request/Response Schema, nmwg.internet2.edu, 2006.
- [18] J. Steiner, C. Neuman, and J. Schiller, "Kerberos: An Authentication Service for Open Network Systems," *Proc. Usenix Ann. Technical Conf.*, 1998.
- [19] F. Berman and R. Wolski, "Scheduling from the Perspective of the Application," *Proc. Fifth Int'l Symp. High-Performance Distributed Computing (HPDC '96)*, 1996.
- [20] A. Grimshaw, W. Wulf, J. French, A. Weaver, and P. Reynolds, "Legion: The Next Logical Step towards a Nationwide Virtual Computer," Technical Report CS-94-21, Univ. of Virginia, 1994.
- [21] T. Defanti, I. Foster, M. Papka, R. Stevens, and T. Kuhfuss, "Overview of the I-WAY: Wide Area Visual Supercomputing," *Int'l J. Supercomputer Applications*, 1996.
- [22] D. Reed and C. Mendes, "Intelligent Monitoring for Adaptation in Grid Applications," *Proc. IEEE*, 2005.
- [23] R. Prasad, M. Jain, and C. Dovrolis, "Effects of Interrupt Coalescence on Network Measurements," *Proc. Sixth Passive and Active Measurement Workshop (PAM '04)*, 2004.
- [24] J. Han and F. Jahanian, "Impact of Path Diversity on Multi-Homed and Overlay Networks," *Proc. 38th IEEE Ann. Conf. Dependable Systems and Networks (DSN '04)*, 2004.
- [25] Boston Univ., BRITE: Representative Internet Topology Generator, <http://www.cs.bu.edu/brite>, 2006.
- [26] Middleware Initiative, *NWS User's Guide*, [http://archive.nsf-middleware.org/documentation/NMI-R5/0/gridscenter/NWS/users\\_guide.htm](http://archive.nsf-middleware.org/documentation/NMI-R5/0/gridscenter/NWS/users_guide.htm), 2006.



**Prasad Callyam** received the BS degree in electrical and electronics engineering from Bangalore University, India, in 1999 and the MS degree in electrical and computer engineering from the Ohio State University in 2002. He is currently working toward the PhD degree in electrical and computer engineering at Ohio State University. He is also currently a senior systems developer/engineer at OARnet, a division of the Ohio Supercomputer Center. His current research interests include network management, active/passive network measurements, voice and video over IP, and network security. He is a student member of the IEEE.



**Chang-Gun Lee** received the BS, MS, and PhD degrees in computer engineering from Seoul National University, Seoul, in 1991, 1993, and 1998, respectively. He is currently an assistant professor in the School of Computer Science and Engineering, Seoul National University. Previously, he was an assistant professor in the Department of Electrical and Computer Engineering at the Ohio State University from 2002 to 2006, a research scientist in the Department of Computer Science at the University of Illinois, Urbana-Champaign, from 2000 to 2002, and a research engineer in the Advanced Telecommunications Research Laboratory at LG Information and Communications, from 1998 to 2000. His current research interests include real-time systems, complex embedded systems, ubiquitous systems, QoS management, wireless ad hoc networks, and flash memory systems. He is a member of the IEEE and the IEEE Computer Society.



**Eylem Ekici** received the BS and MS degrees in computer engineering from Bogazici University, Istanbul, in 1997 and 1998, respectively, and the PhD degree in electrical and computer engineering from the Georgia Institute of Technology, Atlanta, in 2002. He is currently an assistant professor in the Department of Electrical and Computer Engineering at the Ohio State University. He is an associate editor of the *Computer Networks Journal* and *ACM Mobile Computing and Communications Review*. He has also served as a cochair of the technical program committee (TPC) of the 2007 IFIP/TC6 Networking Conference. His current research interests include wireless sensor networks, vehicular communication systems, and next-generation wireless systems, particularly routing and medium access control protocols, resource management, and analysis of network architectures and protocols. He is a member of the IEEE.



**Mark Haffner** received the BS degree in electrical engineering from the University of Cincinnati in 2006. He is currently working toward the MS degree in electrical and computer engineering at the Ohio State University. His current research interests include active/passive network measurements, radio frequency (RF) circuit design, and software-defined radios. He is a student member of the IEEE.



**Nathan Howes** is currently working toward the BS degree in computer science and engineering at the Ohio State University. His current research interests include active/passive network measurements and network security. He is a student member of the IEEE.

▷ For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).