

GRID RESOURCE MANAGEMENT AND SCHEDULING FOR DATA STREAMING APPLICATIONS

Wen ZHANG

Chongqing Military Delegate Bureau

General Armament Department of PLA, Chongqing 400060, P.R. China

This work was carried out when the author was with:

Department of Automation

Tsinghua University, Beijing 100084, P.R. China

Junwei CAO

Research Institute of Information Technology

Tsinghua National Laboratory for Information Science and Technology

Tsinghua University, Beijing 100084, P.R. China

e-mail: jcao@tsinghua.edu.cn

Yisheng ZHONG

Department of Automation

Tsinghua University, Beijing 100084, P.R. China

Lianchen LIU

National CIMS Engineering Research Center

Tsinghua University, Beijing 100084, P.R. China

Cheng Wu

National CIMS Engineering Research Center

Tsinghua University, Beijing 100084, P.R. China

Manuscript received 27 April 2009; revised 6 November 2009

Communicated by Liberios Vokorokos

Abstract. Data streaming applications bring new challenges to resource management and scheduling for grid computing. Since real-time data streaming is required as data processing is going on, integrated grid resource management becomes essential among processing, storage and networking resources. Traditional scheduling approaches may not be sufficient for such applications, since usually only one aspect of grid resource scheduling is focused. In this work, an integrated resource scheduling approach is proposed and coordinated resource allocation of CPU cycles, storage capability and network bandwidth is implemented. Resource allocation is performed periodically with updated information on resources and applications and heuristic search for optimal solutions is used to assign various resources for running applications simultaneously. Performance metrics considered in this work include data throughput and utilization of processors, storage, and bandwidth, which are actually tightly coupled with each other when applied for grid data streaming applications. Experimental results show dramatic improvement of performance and scalability using our implementation.

Keywords: Grid computing, data streaming, resource management, genetic algorithm

Mathematics Subject Classification 2000: 65K10, Secondary 90C59

1 INTRODUCTION

Data streaming applications are becoming more popular recently, such as astronomical observations, large-scale simulation and sensor networks, which brings new challenges to grid resource management. Characteristics of such applications include that

1. they are continuous and long running in nature;
2. they require efficient data transfers from distributed data sources;
3. it is often not feasible to store all the data in entirety for later processing because of limited storage and high volumes of data to be processed;
4. they need to make efficient use of high performance computing (HPC) resources to carry out computation-intensive tasks in a timely manner.

Grid computing [1] paves a new way to enable such applications with cross-domain resource sharing and service integration supports.

When there is a shortage of CPU processing capability located at data sources, there is a requirement that data have to be streamed to computational resources for processing. For example, LIGO (Laser Interferometer Gravitational-wave Observatory) [2, 3] is generating 1TB scientific data per day and trying to benefit from processing capabilities provided by the Open Science Grid (OSG) [4]. Since

most OSG sites are CPU-rich but storage-limited with no LIGO data available, data streaming supports are required in order to utilize OSG CPU resources. In such a data streaming scenario, data should be pulled rather than pushed to the computational system in the form of streams of tuples, and processing is continuously executed over these streams as if data were always available from local storage. What's more, data arrival rates must be controlled to match the processing speeds to avoid waste of computational capacity or data overflow. Meanwhile, processed data have to be cleaned up to save space for the subsequently coming data.

Grid enabled data streaming applications require management of various grid resources, e.g. CPU cycles, storage capability and network bandwidth. In this paper, an integrated resource management and scheduling system for grid data streaming applications is developed to improve data throughput, processor utilization, storage usage and bandwidth utilization in a coordinated way. When a new application arrives, admission control is invoked to decide whether to start or queue it. Accepted applications are allocated with appropriate resources at the end of each scheduling period, together with running ones. A heuristic approach, genetic algorithm (GA) [5] is applied to find satisfactory resource allocation scheme in the given scheduling period with updated information of resources and applications. Scheduling evolves periodically with updated status of resources and applications since the grid is a shared environment where resources are not dedicated. Based on the Globus toolkit [6], the system is able to discover and manage resources geographically distributed and belonging to different management domains in a transparent and secure way. Evaluation results show excellent performance and scalability of this system.

The rest of this paper is organized as follows: Section 2 provides a formal representation of the optimization issue with predefined performance metrics; corresponding algorithms are elaborated in Section 3; performance evaluation results are illustrated in Section 4; related work is discussed in Section 5; and Section 6 concludes the paper.

2 GRID DATA STREAMING – PROBLEM STATEMENT

As mentioned above, an integrated resource management and corresponding scheduling algorithms are required to make full resource utilization while keeping optimal performance of each data streaming application. The approach tries to accommodate as many applications as possible simultaneously to make the best use of resources in preconditions that requirements of each application can also be met. In this section, the schedule issue is described in a formal way and performance metrics are defined.

2.1 Formal Representation

A resource pool R described in this work includes processors (P), storage (S) and network bandwidth (B) that have to be allocated to data streaming applications in

an integrated manner. Suppose n is the total number of processors P in the resource pool and there are m applications (A) for data streaming and processing.

$$R = \{P, S, B\}$$

$$P = \{p_i | i = 1, 2, \dots, n\}$$

$$A = \{a_j | j = 1, 2, \dots, m\}$$

Let s and b be the total storage space S and network bandwidth B , respectively. In general, the sum of processors, storage and bandwidth allocated to each application cannot exceed the total available resources.

$$\sum_{j=1}^m n_j \leq n$$

$$\sum_{j=1}^m s_j \leq s$$

$$\sum_{j=1}^m b_j \leq b$$

For each application a_j , there is a corresponding minimum requirement of resources that has to be met for the application to be executed.

$$R_j^{req} = \{n_j^{req}, s_j^{req}, b_j^{req}\}$$

$$n_j^{req} \leq n_j \leq n, j = 1, 2, \dots, m$$

$$s_j^{req} \leq s_j \leq s, j = 1, 2, \dots, m$$

$$b_j^{req} \leq b_j \leq \min(b, b_j^{\max}), j = 1, 2, \dots, m$$

All of the above constraints have to be met during resource scheduling and allocation for data streaming applications. Note that the bandwidth allocated to the application a_j is constrained by both available bandwidths locally b and remotely at the data source end b_j^{\max} . The major challenge is that the three different types of resources are correlated to each other inherently in deciding the performance of a scheduling and allocation scheme. In the next section, major performance metrics considered in this work are described.

2.2 Performance Metrics

There are many aspects of performance criteria when evaluating resource allocation schemes. Specifically, for data streaming applications, data throughput, the amount of data streamed and processed during a given period of time, is most important. Other performance metrics that have to be considered simultaneously are resource utilization and scheduling scalability.

Suppose an evaluation period includes l time units. A time unit t is a predefined minimum time period, based on which all resource scheduling and allocation are carried out. Let b_{jk}^{usg} ($j = 1, 2, \dots, m; k = 1, 2, \dots, l$) be the bandwidth usage of the application a_j during the k^{th} time unit and s_{jk}^{usg} ($j = 1, 2, \dots, m; k = 1, 2, \dots, l$) be the storage usage at the beginning of the k^{th} time unit. Note that actual resource usage of an application is usually different from corresponding resources allocated to an application. We can calculate the total data throughput TP as follows:

$$\begin{aligned} TP_{jk} &= b_{jk}^{\text{usg}}t + s_{jk}^{\text{usg}} - s_{j(k+1)}^{\text{usg}} \quad (j = 1, 2, \dots, m; k = 1, 2, \dots, l) \\ TP_j &= \sum_{k=1}^l TP_{jk} = \sum_{k=1}^l b_{jk}^{\text{usg}}t + s_{j1}^{\text{usg}} - s_{j(l+1)}^{\text{usg}} \quad (j = 1, 2, \dots, m) \\ TP &= \sum_{j=1}^m TP_j = \sum_{j=1}^m \sum_{k=1}^l b_{jk}^{\text{usg}}t + \sum_{j=1}^m (s_{j1}^{\text{usg}} - s_{j(l+1)}^{\text{usg}}). \end{aligned}$$

The total data processing throughput is the difference of storage usage plus all data streamed into the system during the given period. This is based on the assumption that just-in-time data cleanup is enabled and all processed data is cleaned up from storage at the end of each time unit. If the evaluation period covers all the makespan of an application, it is obvious that s_{j1}^{usg} and $s_{j(l+1)}^{\text{usg}}$ are both zero and the total data processing throughputs for a given application can be represented purely via bandwidth usage. To simplify the problem, we assume that data throughputs of each application TP_j are *comparable* with each other, so that a sum up of all TP_j can be used to represent the overall data throughput. If this is not the case in a real environment, some normalization has to be performed to weight data throughputs of different applications in terms of data throughput.

Resource utilization is another concern when enabling data streaming applications. It is straightforward to calculate storage and bandwidth usage percents of the k^{th} time unit as follows:

$$\begin{aligned} US_{jk} &= \frac{s_{jk}^{\text{usg}}}{s} \quad (j = 1, 2, \dots, m; k = 1, 2, \dots, l) \\ US_k &= \frac{\sum_{j=1}^m s_{jk}^{\text{usg}}}{s} \leq 1 \quad (k = 1, 2, \dots, l) \\ US &= \frac{\sum_{k=1}^l \sum_{j=1}^m s_{jk}^{\text{usg}}}{ls} \leq 1 \\ UB_{jk} &= \frac{b_{jk}^{\text{usg}}}{b} \quad (j = 1, 2, \dots, m; k = 1, 2, \dots, l) \\ UB_k &= \frac{\sum_{j=1}^m b_{jk}^{\text{usg}}}{b} \leq 1; k = 1, 2, \dots, l \\ UB &= \frac{\sum_{k=1}^l \sum_{j=1}^m b_{jk}^{\text{usg}}}{lb} \leq 1. \end{aligned}$$

The utilization of CPU cycles can be calculated indirectly via storage usage, since for data streaming applications it can be assumed the allocated processor is busy when there is available data in local storage, and idle when no data is available locally. Suppose that P_{jk} is the set of processors that are allocated to the application a_j during the k^{th} time unit. Let M_{ik} be a 2D array to identify whether the processor p_i is busy or idle at the k^{th} time unit.

$$M_{ik} = \begin{cases} 1 & \text{if } \forall j \quad p_i \in P_{jk} \text{ and } s_{jk}^{usg} > 0 \\ 0 & \text{if } \exists j \quad p_i \notin P_{jk} \text{ or } s_{jk}^{usg} = 0 \end{cases} \quad (i = 1, 2, \dots, n; k = 1, 2, \dots, l)$$

The processor usage percentage is calculated as follows:

$$UP_i = \frac{\sum_{k=1}^l M_{ik}}{l} \leq 1 (i = 1, 2, \dots, n)$$

$$UP = \frac{\sum_{i=1}^n \sum_{k=1}^l M_{ik}}{nl} \leq 1.$$

The resource management and scheduling issue for grid data streaming applications can be transformed into an optimization problem:

$$\begin{aligned} & P \\ & \max TP \\ & \max UP \\ & \min UB \\ & \text{s.t.} \\ & UP_i \leq 1 \quad (i = 1, 2, \dots, n) \\ & US_k \leq 1 \quad (k = 1, 2, \dots, l) \\ & UB_k \leq 1 \quad (k = 1, 2, \dots, l) \end{aligned}$$

where first two goals are to process more data and match data processing and streaming capability as much as possible while the third one is to utilize bandwidth in an economic way to avoid congestion. These three goals conflict in nature, and some tradeoffs have to be made. Currently, we focus more on the overall data throughput. Algorithms provided in the next section are processing-, storage-, and congestion-aware, so the last two goals can be maintained in most cases. Note that storage usage is not included in optimization goals because storage usage does not affect the ultimate throughput, but adequate storage will indeed increase the robustness of data processing. Available processors, storage and bandwidth are considered as constraints.

3 RESOURCE SCHEDULING AND ALLOCATION – KEY ALGORITHMS

There are two steps for resource scheduling: admission control is performed to decide whether a new application is started, according to its resource requirement and current status of available resources in the computing pool; the GA is performed periodically to improve resource utilization and to meet the requirements of active applications in an evolving way. Resource allocation is performed iteratively together with periodical scheduling of key parameters.

In this section, the key components of this resource management and scheduling scheme are described in detail. Overall flow chart for such a scheduling process is illustrated in Figure 1.

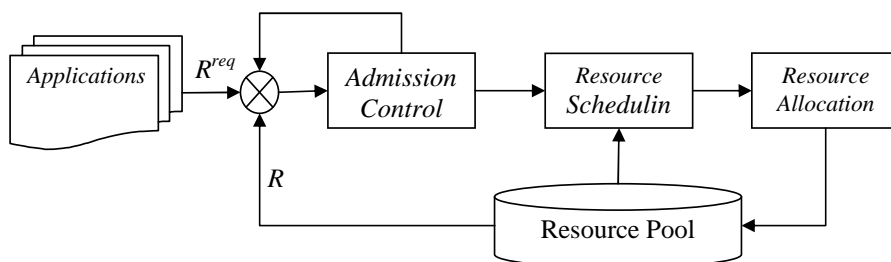


Fig. 1. The flow chart of resource scheduling and allocation for data streaming applications

A coming application has to specify explicitly its requirements of resources, R^{req} . Available resources R in the resource pool are monitored in real time. Both R^{req} and R are input to the admission control module to decide whether or not the coming application should be accepted and put to the active application set. Resource scheduling only works on active application set periodically to produce scheduling parameters. In this work, the GA is adopted as an evolving method to absorb dynamically changing resource and application status. Resource allocation takes scheduling parameters as inputs and generates final allocation schemes. Resource allocation occurs iteratively with a much higher frequency than resource scheduling to improve the overall system performance. These are described in details below.

3.1 Admission Control

It is obvious that the resource pool in general cannot support infinite applications simultaneously, and too many applications will lead to fierce resource competition which may decrease overall processing efficiency as a whole. The mechanism of admission control plays an essential role in our system for resource management.

When a new task is submitted, the admission control decides to accept it instantly or just to keep it in a waiting queue and resubmit it in future. This decision is made according to the usage status of resources and application requirements. Each

task can specify its minimum requirement of resources, e.g. it needs some processors of certain types, minimum bandwidth and storage. An XML schema is developed for application requirement specification.

Applications can specify CPU type, bandwidth and storage requirement. The system checks up current available resources. For example, an application compiled on X86_64 cannot run on I386 processors, so not every processor is suitable for the application. Suppose that the number of those available X86_64 processors is larger than n_j^{req} , and unallocated storage and bandwidth are both larger than s_j^{req} and b_j^{req} , respectively, the task can be immediately put to be active for execution. If any of resources is not available, the task would just be kept in the waiting queue.

Applications in the queue are called new eligible applications (NEAs). NEAs can have different priorities for resource scheduling and allocation. They are classified into different groups according to different priorities, and in each group, the first-come-first-serve (FCFS) strategy is applied.

3.2 Resource Scheduling

As CPU, bandwidth and storage are integrated as a whole in data streaming applications, there must be a so-called optimal operation point (OOP) to make balance in resource usage. The OOP defines the expected bandwidth, computing power, and granularity of data streaming (storage usage) which simultaneously maximizes the uses of bandwidth and CPU power. Our scheduler benefits greatly from using an optimal combination of resources, instead of making independent scheduling schemes for each type of resources.

Essentially, integrated scheduling of resources for applications is a NP-complete problem, and in this work the GA is adopted to find satisfactory, not necessarily optimal, solutions in a relatively short time. The GA is required to recalculate scheduling parameters in each scheduling period, with updated status information of resources and applications.

Using the GA, a single chromosome is composed of three parts of scheduling parameters, leading to an allocation scheme of processors, storage and bandwidth, respectively. The g^{th} generation of chromosome can be represented as follows:

$$CHROM_g = \{p_{jg}, s_{jg}, \alpha_g, \beta_g, \rho_g, \mu_{jg} \mid j = 1, 2, \dots, m\}.$$

Similar to definitions in Section 2.2, suppose that P_{jg} is the set of processors that are allocated to the application a_j during the g^{th} generation of GA evolving. To simplify the problem, each application just needs one processor, so p_{jg} is used above instead of P_{jg} . For an application, it is expected to run on a fixed processor till it is finished, or to put it another way, no migration from one processor to another occurs during task execution. Only parts of p_{jg} for new applications are involved in the GA evolving and previously assigned processors are fixed to be allocated to its current application. s_{jg} represents the maximum storage allocated to the application a_j . For a certain application, its lower limit of storage has to be larger than s_j^{req} and can be

set to be a proportion of s_{jg} . Details are included in the next section on resource allocation. Similarly, scheduling parameters above, α , β , ρ and μ_j ($j = 1, 2, \dots, m$) are corresponding to bandwidth allocation. How these parameters are associated with a bandwidth allocation scheme is also described in details in the next section. Three parts of a chromosome evolve independently with their own rules, decreasing computational complexity and avoiding meaningless heuristic searches.

The evaluation index, or fitness, of each chromosome is set to be data throughput, i.e., the amount of data processed in a scheduling interval. As mentioned before, we consider all the data for different applications equally. In calculating its throughput, information on applications and resources has to be updated and performance prediction is enabled using historical information. Given a chromosome with scheduling parameters, with historical performance information and some priority of resources and applications, data throughput in a scheduling interval can be calculated using formulations in Section 2.2. In a scheduling period, scheduling parameters are initiated from its direct foregoing period. During the evolution, two chromosomes are crossed to generate two new ones for the next generation, and genetic mutation happens in some chromosomes with a given probability. The chromosome that leads to the highest data throughput value can be achieved and corresponding scheduling parameters are used to generate a resource allocation scheme in a scheduling period.

Although it is hard to find the global OOP since applications and resources are varying constantly and it is not easy to define expected processors, bandwidth and storage precisely, evolving searching capability of the GA guarantees a satisfactory solution for each scheduling period. Allocation algorithms for processors, storage and bandwidth are given below.

3.3 Resource Allocation

Given a set of correlated scheduling parameters generated by the GA described above, resource allocation schemes can be achieved using the method given in this section. While scheduling parameters are fixed within each scheduling period, resource allocation schemes can still change; for instance, bandwidth allocation is an iterative process.

3.3.1 Processor Allocation

As mentioned above, processor assignment is a match making process. Both applications and resources can specify their own requirements. Processors can be classified into several groups according to their architectures. Similar processors in the same group may also have different frequencies that may lead to different data processing performance. NEAs can also be organized into several groups with different priorities. In each group, the selecting principle is FCFS.

Matchmaking is carried out to find appropriate processors for applications. The processors whose characteristics do not conflict with the application requirements

form a candidate set. Then applications with higher priorities find their matched processors first. In a *CHROM*, p_j is the number of the assigned processor for each application. In different generations of evolving in a given scheduling period, processor assignments of chromosomes in successive generations are independent to guarantee all possible assignments can be covered. As we have mentioned, no migration of applications exist, so in each scheduling period, the algorithm is only required to allocate processors for the NEAs.

3.3.2 Storage Allocation

The overall principle for storage allocation is to make full usage of storage to increase robustness while getting ready for new coming applications. If there are only few applications running in the resource pool, the storage allocated for each application can be set to a high value. While the number of applications increases, allocated storage for each application may be decreased. So quotas for each application can be scalable accordingly, and there must be some margin of storage for potentially new-coming applications.

Step 1: Initialization. As supposed, there are m applications in the pool, to generate m random numbers, $r_j \in (0, 1)$, $j = 1, 2, \dots, m$. Calculate each quota, q_j as follows:

$$q_j = \frac{r_j}{\sum r_j} (j = 1, 2, \dots, m).$$

Step 2: If $s_j = q_j s \geq s_j^{req}$, reserve these numbers for initially allocated storage for the application a_j . Otherwise, repeat step 1 until $s_j \geq s_j^{req}$ ($j = 1, 2, \dots, m$) becomes true, where s_j^{req} is the minimal required storage of application a_j as defined in Section 2.1.

Step 3: Repeat steps 1 and 2 until all the storage allocation schemes are set for the chromosomes in a population, and these would be initial values for the first generation.

Step 4: Chromosome crossing. Two chromosomes cross to generate new ones for the next generation as follows:

$$\begin{aligned} \varepsilon_{jg} s_{jg} + (1 - \varepsilon_{jg}) s_{j'g} &= s_{j(g+1)} \\ (1 - \varepsilon_{jg}) s_{jg} + \varepsilon_{jg} s_{j'g} &= s_{j'(g+1)} \end{aligned}$$

where $0 < \varepsilon_{jg} < 1$ is a random number evenly distributed in $(0, 1)$, and the other chromosome, j' , is selected at random.

Step 5: Repeat step 4 until all generations are covered.

A storage allocation scheme includes both the upper and lower limit values. As mentioned in Section 3.2, the lower limit of storage has to be larger than s_j^{req} and can be set to be a proportion of s_{jg} . Lower and upper limits are mainly used as thresholds to control start and end times of data streaming: when data amount

scratches the lower limit, more data should be transferred until the amount reaches the upper limit. Since there are also data cleanups involved, data amount keeps changing and varies between lower and upper limits.

The upper limit for each application is used to guarantee that the overall amount of data in local storage does not exceed available storage space. The lower limit is used to guarantee that data processing can survive network collapse when no data can be streamed from sources to local storage for a certain period of time, which improves system robustness and increases CPU resource utilization.

3.3.3 Bandwidth Allocation

Bandwidth allocation plays an important role in the whole resource allocation scheme, for appropriate bandwidth is indispensable to guarantee data streaming for applications to make them run constantly. Different from traditional bandwidth allocation, our scheme is storage aware, i.e., data streaming may be intermittent rather than continuous to avoid data overflow, for allocated storage for each application is limited. When the storage is full of data, streaming is halted for a while until some data has been processed and cleaned up so that storage is released for more data. At any moment, the amount of data in storage for each application is determined by both data streaming and cleanup processes.

A utility function $U_j(b_j)$ is introduced when its data is streamed with allocated bandwidth b_j . $U_j(b_j)$ should be concave, continuous, bounded and increasing in the interval $[b_j^{\text{req}}, b]$. Note that it is not necessary that identical utility functions are chosen for all applications. We try to maximize the sum of the utilities of all the applications, maintaining fairness among them.

Due to the usage status of storage and repertory policy with lower and upper limits of storage allocation, there are two possible states for each application a_j at any time, i.e., active and inactive, which indicates whether a data streaming is on or off. Let s_{jk} and λs_{jk} be the upper and lower limits of storage allocation. λ is a predefined fixed proportion of storage allocation between 0 and 1. Let A_{jk} be the state identity for the application a_j at the k^{th} step. A_{jk} can be initiated as active (1) and evolves as follows:

$$A_{j(k+1)} = \begin{cases} 1 & \text{if } A_{jk} = 1 & \text{and } s_{jk}^{\text{usg}} < s_{jk} \\ 0 & \text{if } s_{jk}^{\text{usg}} \geq s_{jk} \\ 0 & \text{if } A_{jk} = 0 & \text{and } s_{jk}^{\text{usg}} > \lambda s_{jk} \\ 1 & \text{if } s_{jk}^{\text{usg}} \leq \lambda s_{jk} \end{cases}$$

An iterative optimization algorithm is proposed in [7] and its convergence is analyzed, but it is required to be congestion-aware, which is hard to be satisfied in the wide-area Internet. According to our situation, we make some modification

upon it as follows:

$$b_{j(k+1)} = \begin{cases} [b_{jk} + \alpha_k U'(b_{jk})]_j & \text{if } A_{jk} = 1 \quad \text{and} \quad \sum_{A_{jk}=1} b_{jk} \leq \rho_k b \\ [\beta_k b_{jk}]_j & \text{if } A_{jk} = 1 \quad \text{and} \quad \sum_{A_{jk}=1} b_{jk} > \rho_k b \\ 0 & \text{if } A_{jk} = 0 \end{cases}$$

Here, b_{jk} is the bandwidth allocation for the application a_j at the k^{th} step of iterations. Note that there are many small iterations during a scheduling period. α_k , β_k and ρ_k are all positive sequences. ρ_k is the so-called safety coefficient to avoid bandwidth excess, where $\rho_k \in (0, 1)$, i.e., there is some margin from the full use of total bandwidth for flexibility and robustness. For the sake of convenience, they are usually substituted as a fixed value within multiple iterations of a scheduling period and only evolve when GA is applied to generate new scheduling parameters when a new scheduling period starts, as described in Section 3.2. $[\cdot]_j$ denotes a projection on the application a_j and can be calculated as

$$[x]_j = \min(b, b_j^{\max}, \max(b_j^{req}, x)).$$

$U'(\cdot)$ is the sub gradient of

$$U(\cdot) = \sum_{j=1}^M U_j(b_{jk})$$

and

$$U'(b_{jk}) = \frac{\partial U(\cdot)}{\partial b_{jk}}.$$

A popular utility function can be expressed as:

$$U_j(b_{jk}) = \mu_{jk} \ln(1 + b_{jk}), j = 1, 2, \dots, m$$

where μ_j stand for applications' coefficient. As we can see, given allocation schemes of processors and storage as mentioned above, with settled parameters such as α , β , ρ and μ_j , a bandwidth allocation scheme can be obtained, and performance metrics, e.g. data throughput and resource usage, can be calculated. For bandwidth allocation, it is transformed to finding the optimal set of such parameters, which is performed at the beginning of each scheduling period, as described in Section 3.2.

Essentially, the bandwidth allocation approach proposed in our work is a type of additive increase multiplicative decrease (AIMD) algorithm that is usually used in TCP congestion avoidance. Main feature of our approach is storage awareness, since data streaming is stopped if allocated storage is nearly full such that data overflow is avoided. The approach is also processing-aware, because the processing capacity can be reflected via storage usage. Our bandwidth allocation leads to on-demand data streaming, which is also congestion aware.

Bandwidth allocation is finally implemented in our system using GridFTP [8] through controlling start and end times of data transfers and tuning parallelism for each application. System implementation and experimental results are given below.

4 SYSTEM IMPLEMENTATION AND PERFORMANCE EVALUATION

In order to verify our approach proposed above, a grid environment has been established at Tsinghua University (Beijing, China). The Globus Toolkit 4.0.1 is deployed to provide common grid services and a simple Certificate Authority is set up to sign certificates for hosts and users which are used to establish a secure and transparent environment for data streaming applications. Detailed information on our system implementation and experimental results on our resource scheduling approach are included in this section.

4.1 System Implementation

Key components in the architecture of our resource management and scheduling system include a client tool for job submission, a management engine for admission control, a scheduler, a dispatcher and application wrappers.

Client Tool. This tool is an interface for users to submit their applications with requirement specification in XML format, including the executables, processor number and architecture, minimum bandwidth and storage requirements, data sources, etc. It is also capable of monitoring status of submitted applications and resources, halting and canceling submitted jobs.

Management Engine. The management engine accepts users' submissions and puts them into the queue, which can be accessed by the scheduler. Its main function is to provide grid supports for streaming applications, such as security, resource discovery and management. The components of Globus toolkit used here include GSI (Globus Security Infrastructure) and MDS (Monitoring and Discovery Service). The management engine is responsible for admission control.

Scheduler. This is the core component in the whole architecture and its key algorithms are discussed in detail in Section 3. It is responsible for generating the scheduling parameters. Its instruction will be executed by the dispatcher.

Dispatcher. The dispatcher is in charge of resource allocation, sending executables with their description files to appropriate processors and invoking. This component interacts with services provided by grid middleware, such as Globus GRAM.

Application Wrappers. This component parses the description file according to the XML schemas, initializes execution of executables, and starts data streaming to specified storage with allocated bandwidth. Also, results are sent back through the dispatcher.

Beside allocation of computational resources as most traditional grid resource management and scheduling systems do, our system also deals with allocation of bandwidth and storage to support real-time data transfer which is required by data streaming applications. Management and scheduling of processors, bandwidth and storage are carried out in an integrated rather than independent way.

4.2 Experiment Design

As shown in Figure 2, a case study for LIGO data analysis is taken as a typical data streaming application. Two streams of data from two observatories (H1 and L1), represented using two file lists, are inputs of the program, *rmon*. The program calculates the *r*-Statistic of two data sets to measure signal similarity. A signal has more possibility to be a gravitational wave candidate if it occurs at two observatories simultaneously. Once the program *rmon* is launched, an option file is used to define data stride and channels for each step of calculation. Data can be used in a streaming way since data is organized in time series and becomes obsolete after being processed. Note that this is a simplified example of LIGO data analysis since in a real world scenario there are also pre- and post- steps and a complete data analysis is usually carried out in a formal of pipelines.

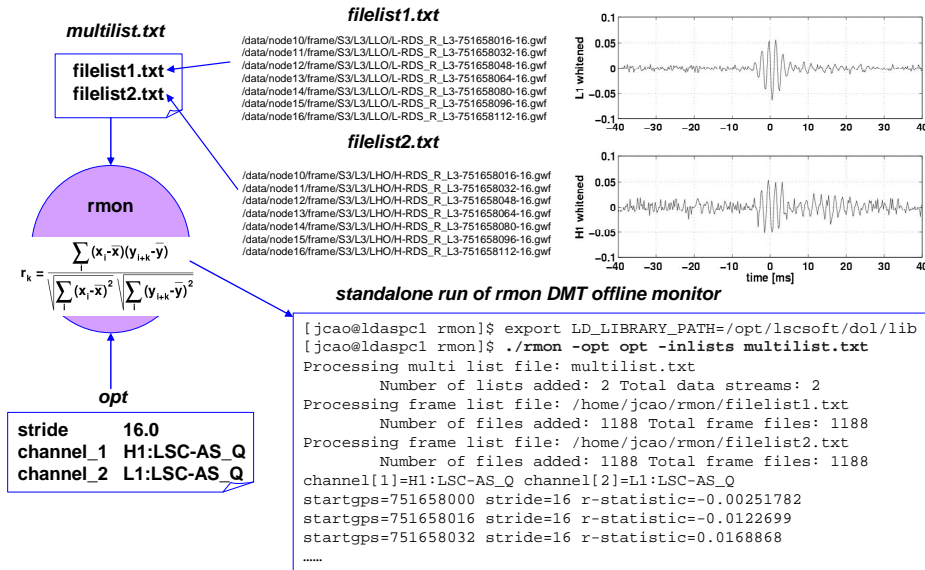


Fig. 2. Experiment design – LIGO data analysis case study

The grid environment established at the campus is connected to Internet with limited bandwidth, and a network file system (NFS) is configured to which all data streams are directed. Each application is given a unique number to identify its di-

rectory in the NFS where its data is streamed so that applications can access data as if data were always available locally. In a local area network of the same campus, NSF does not introduce heavy overhead for data access. Tasks are submitted at moments complying with negative exponential distribution law, and their requirements of resources are also explicitly expressed. Experiments are carried out using parameter values listed in Table 1.

Parameters	Values
m	30
n	16(5 + 4 + 7)
s	140 MB
b	30 Mbps
$n_j^{req}, j = 1, 2, \dots, m$	1
$s_j^{req}, j = 1, 2, \dots, m$	8 9 40 10 8 5 5 6 10 6 9 6 10 7 6 6 8 7 7 9 8 8 10 6 9 9 7 8 5 5
$b_j^{req}, j = 1, 2, \dots, m$	1 3 2 1 1 2 3 3 1 3 3 2 3 1 3 3 3 3 2 1 3 3 3 3 1 2 2 1 3 1
$b_j^{max}, j = 1, 2, \dots, m$	7 7 6 5 4 6 6 5 7 7 6 4 4 5 7 4 6 7 5 5 5 7 6 4 5 5 6 6 7 4
$\mu_j, j=1, 2, \dots, m$	15
l	10 000
t	1 s
λ	0.8

Table 1. Experiment design – parameter values

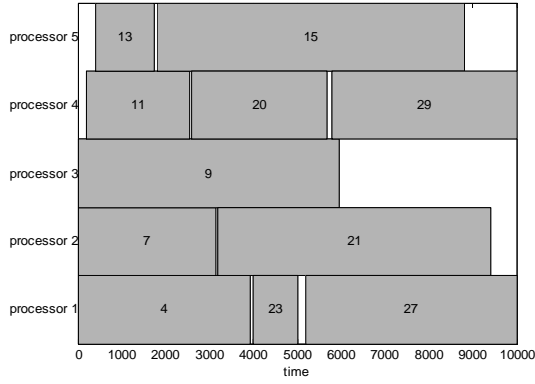
Totally 30 applications are submitted to a grid node with 16 processors, 140 MB disk space, and 30 Mbps bandwidth. Processors are divided into 3 groups, each with a different architecture. Each application requires one processor with different storage and bandwidth requirements. The experiment lasts for 10 000 seconds with 1 second per time unit. Resource scheduling is carried out once per 200 time units with updated resource and application information.

4.3 Experimental Results

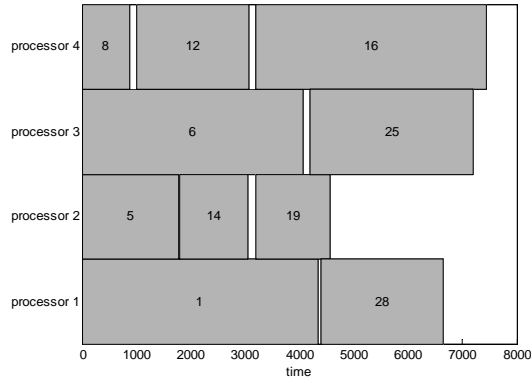
Figures 3–5 illustrate detailed resource usage information of the experiment. The processors are divided into 3 groups. Each group has a separate resource scheduling scenario. While utilization of processors cannot be 100 % since processors can become idle if no data is available for processing, these idle CPU times cannot be reused by other applications once allocated. Different from processor usage, storage and bandwidth usage can adapt to application requirements quickly.

4.4 Performance Evaluation

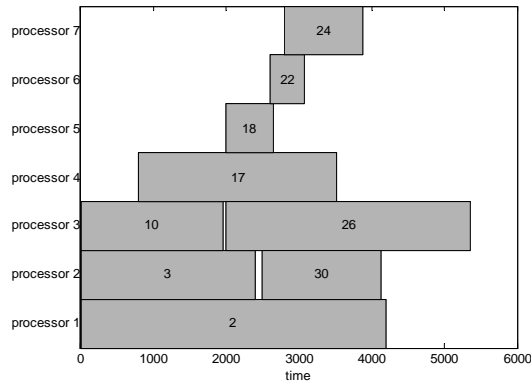
In the last section, raw data on resource usage of processors, storage and bandwidth are visualized. In this section, our approach described in Section 3 is evaluated using performance metrics described in Section 2.2.



(a) Processor Usage in Group 1



(b) Processor Usage in Group 2



(c) Processor Usage in Group 3

Fig. 3. Experiment results – processor usage

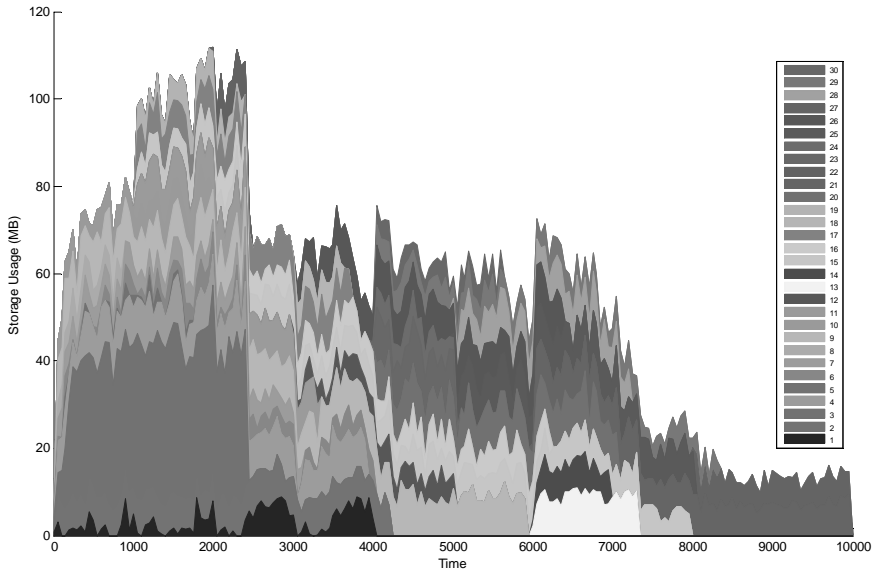


Fig. 4. Experiment results – storage usage

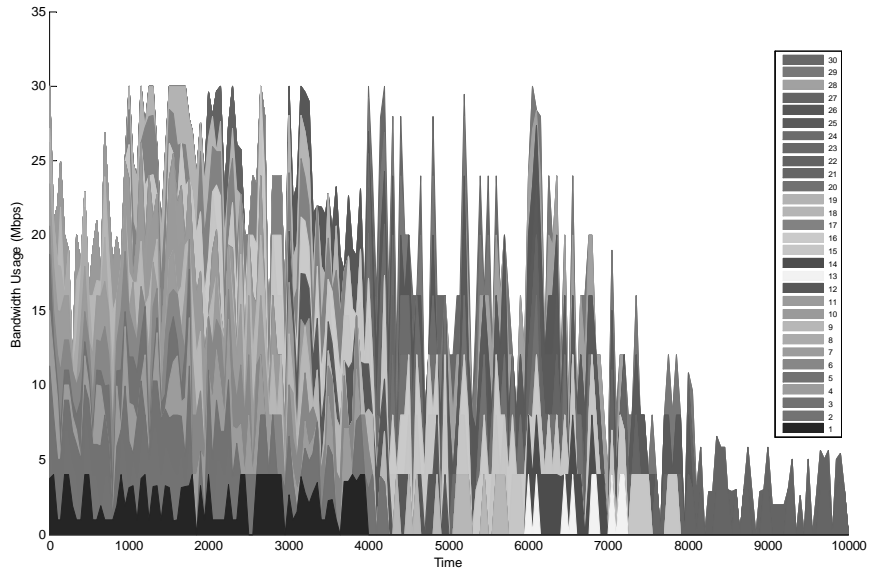


Fig. 5. Experiment results – bandwidth usage

4.4.1 Admission Control

In our experiment, as grid resources are limited compared with data streaming application requirements, it can be inferred that too many applications accepted by the grid without admission control would lead to low processing performance, measured using data throughput in our case. This is verified by further experimental results included in Figure 6.

As shown in Figure 6c), higher data throughput is achieved in the scenario with admission control than in that without admission control. Note that in the latter scenario, since no admission control mechanism is applied, all applications are started to run immediately once submitted and one processor may have to deal with more than one application at the same time. Inadequate data provision for each application and competition of computational power among applications lead to lower data processing performance as a whole.

Figures 6a) and 6b) also provide an overall picture of application status in both scenarios, where red bars with the character *P* stand for pending applications and pink bars with the character *R* mean the corresponding application is running. The makespan of 30 applications is different in the two scenarios. The numbers of completed applications are 29 and 25, respectively. It is obvious that the total makespan is longer in the situation when no admission control is applied.

4.4.2 Processor Utilization

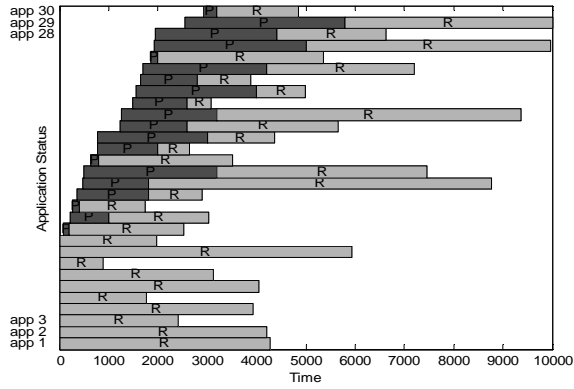
As shown in Figure 3, processors are allocated to applications, one for a single application exclusively. Groups 1 and 2 deal with more applications while they include less processors than Group 3, so average workload of processors is higher than that of Group 3. Because some applications may not be able to be executed on the processors in Group 3, they cannot be transferred to Group 3 for load balancing.

Once an application is started to run on a certain processor, it does not mean that the processor is busy all the time. Since all applications are dependent on data streaming supports, the processor of an application could be idle if no data is available at local storage. Higher processor utilization is always required for a better system performance.

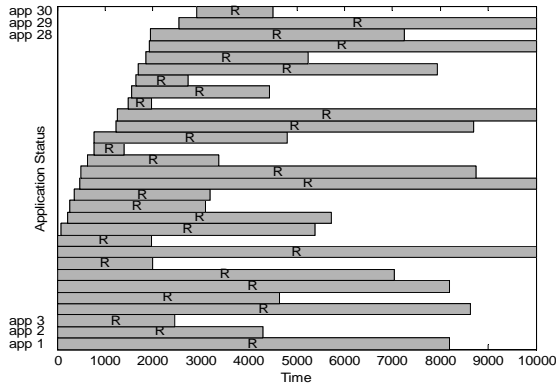
In this section, an additional experiment is carried out with a scheduling scheme in which allocation of computing, bandwidth and storage resources is made independently. As shown in Figure 7, compared with the integrated approach proposed in this work, the independent scheme results in lower processor utilization for most applications. Higher processor utilization is achieved using our integrated scheduling scheme.

4.4.3 Storage Awareness

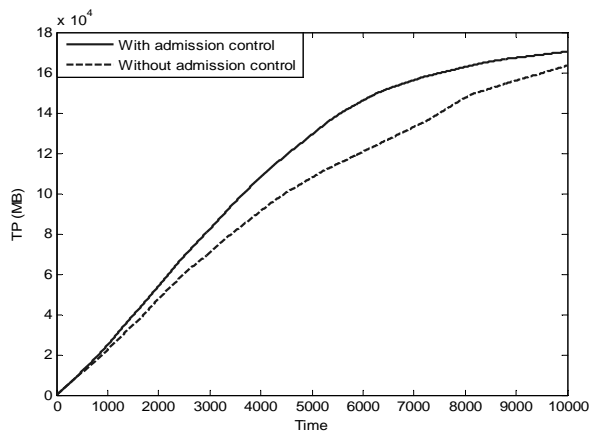
Another feature of our approach is that data streaming is storage-aware, i.e., data transfers are controlled by the usage of allocated storage, rather than spontaneously.



(a) With admission control



(b) Without admission control



(c) Comparison of data throughput

Fig. 6. Performance evaluation – admission control

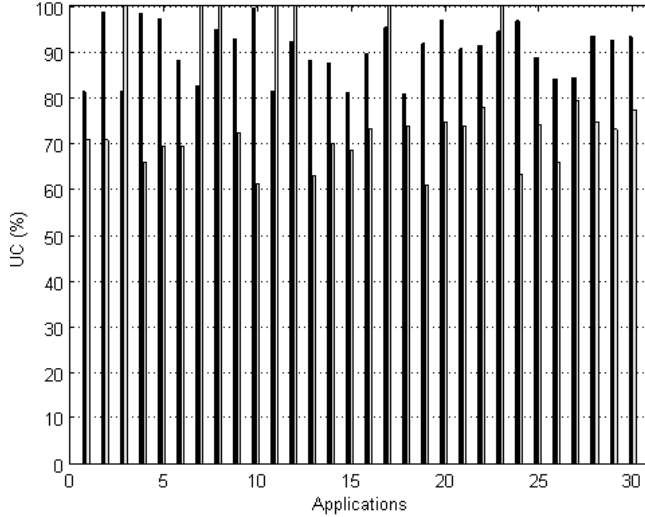


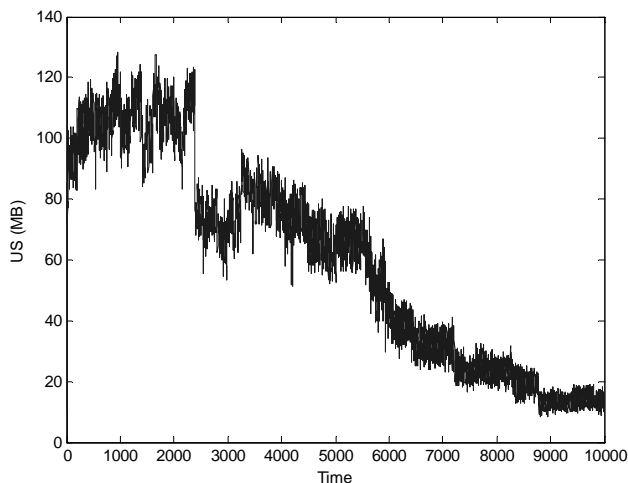
Fig. 7. Performance Evaluation – Processor Utilization

The principle here is that just-enough is OK, not the-more-the-better. Data streaming can be intermittent, not always continuous. In this way high volume of data can be processed with reasonable storage usage, as shown in Figure 8(a). Storage usage varies in a reasonable scope during the experiment. If data streaming is continuous and available storage space is large enough, storage usage can be of high volume, as illustrated in Figure 8(b). Small storage can achieve high throughput for data streaming applications with well-controlled data streaming and processing scheme.

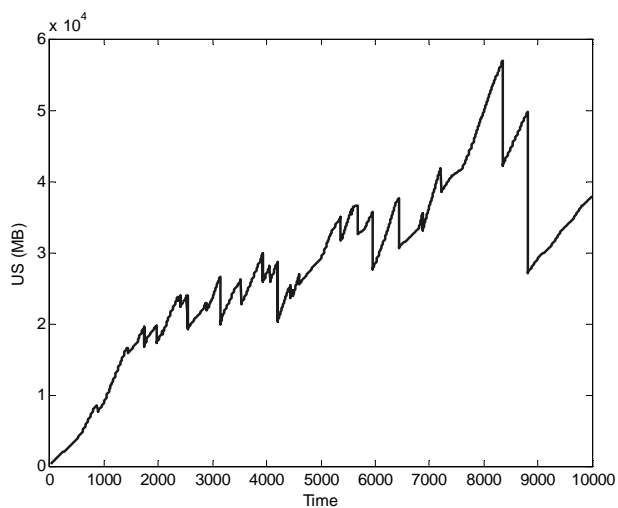
4.4.4 Iterative Bandwidth Allocation

Bandwidth is allocated to each running application to guarantee their data provision. Parameters for bandwidth allocation are obtained using the GA (as described in Section 3.2) and applied in each scheduling period. As described in Section 3.3.3, bandwidth allocation is an iterative process that is adaptive to the total available bandwidth and requirements of running applications.

To justify our iterative bandwidth allocation algorithm, additional experiments are carried out using an even bandwidth allocation method, where bandwidth is allocated to the running applications equally. As shown in Figure 9 a), since total available bandwidth is relatively low ($b = 30$ Mbps), with even bandwidth allocation, only 25 applications are finished in 10 000 time units. In Figure 9 b), available bandwidth is increased to 40 Mbps. In this case, each application can get enough data with the even allocation scheme. These results can be compared with Figure 6 a). Using our iterative approach, a better performance can be achieved with relatively low available bandwidth.



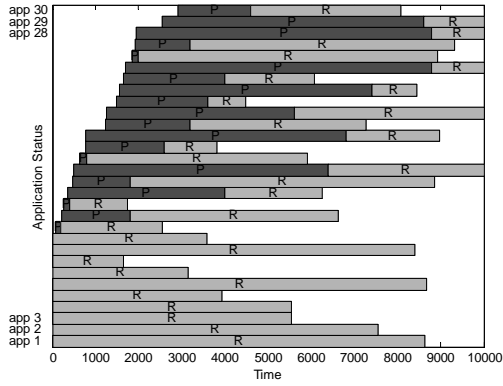
(a) Storage-Aware Data Streaming



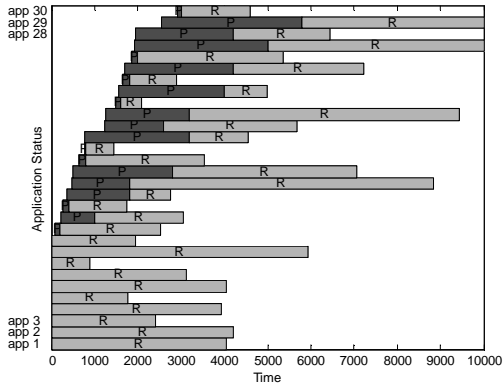
(b) Non-Storage-Aware Data Streaming

Fig. 8. Performance Evaluation – Storage Aware Data Streaming

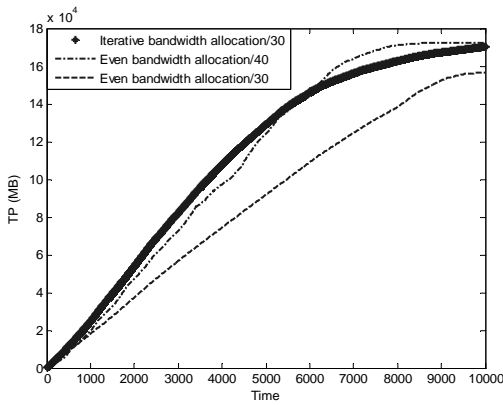
Comparison of data throughput using iterative and even bandwidth allocation is illustrated in Figure 9c). It is obvious that higher data throughput can be achieved using our iterative bandwidth allocation, though only relatively low available bandwidth is available. This is because our approach is storage- and processing-aware. Using the even allocation method, some applications may starve for data while others may be allocated redundant bandwidth, which reduces data processing efficiency and results in lower data throughput.



(a) With Even Bandwidth Allocation ($b = 30$)



(b) With Even Bandwidth Allocation ($b = 40$)



(c) Comparison of Data Throughput

Fig. 9. Performance Evaluation – Bandwidth Allocation

5 RELATED WORK

Stream processing used to be the focus of database research [9, 10, 11, 12] in recent years, and some tools and techniques have been developed to cope with efficient handling of continuous queries on data streams, while our work focuses on scheduling streaming applications on grid resources.

Grid scheduling has primarily focused on providing supports for batch-oriented jobs [13]. Most resource management infrastructures available for grid computing, such as Legion [14], Nimrod/G [15], Condor [16] and [17], are largely geared to support batch-oriented applications rather than streaming ones. Attention has also been paid to scheduling of interactive job in grid environments [18, 19, 20]. There are several existing efforts that are focused on grid resource management and scheduling for data streaming applications:

GATES. A middleware system, called GATES [21, 22, 23] (Grid-based AdapTive Execution on Streams), is developed for applications involving high-volume data streams and requiring distributed processing of data arising from a distributed set of sources. A resource allocation algorithm based on minimal spanning tree (MST) is developed, whose target is to create a deployment configuration, including

1. the number of data sources and their location;
2. the destination, i.e. the node where the final results are needed;
3. the number of stages in the application;
4. the number of instances of each stage;
5. how the instances connect to each other and
6. the node at which each stage is assigned.

Given a streaming application composed of several stages, the first three components of a deployment configuration are determined, and the emphasis of this resource allocation algorithm is to decide the last three components to give the applications the best chance to achieve best performance. Once a deployment configuration is finished, the launcher program can be called to automatically launch the application.

Streamline. Streamline [24, 25] as a scheduling heuristic is designed specially to adapt to the dynamic nature of grid environment and varying demands of a streaming application. It is expected to maximize throughput of the application by assigning best resources to the neediest stage in terms of computation and communication requirements, belonging to a general class of list scheduling algorithms. Stages are prioritized according to their estimated computation and communication costs, and the resources leading to minimal costs will be assigned to the stages. In this algorithm, precise estimation of computation and communication costs must be provided, which is not easy to be implemented.

Pegasus. Pegasus [26, 27] is aimed to address the problem of automatically generating job workflows for the grid, which helps map an abstract workflow defined in terms of application-level components to the set of available grid resources. It handles data transfers, job processing and data cleanups in a workflow manner, not in an integrated and cooperative way. In our approach for grid data streaming applications, data streaming, processing and cleanups have to be processed simultaneously instead of in a workflow manner.

The above projects mainly focus on management and scheduling of computational and data resources for grid applications, where network resources are not considered. The following two projects provide co-allocation of various grid resources but they are not targeting data streaming applications.

EnLIGHTened. As best-effort networks always introduce bottleneck for high-end applications, the EnLIGHTened computing [28] project is aimed to co-allocate any type of Grid resource: computers, storage, instruments, and deterministic, high-bandwidth network paths, including lightpaths, in advance or on-demand. Then external networks, especially dedicated optical networks, can be used as first class resources within grid environments, which will guarantee large data transfers associated with executions on remote resources. Its idea is similar with our integrated resource scheduling proposed in this paper, but we implement more fine-grained resource allocation to cope with specific requirements for data streaming applications, e.g. storage awareness.

G-lambda. G-lambda [29] project is to define a standard web service interface (GNS-WSI) between a grid resource manager and a network resource service from a network resource manager provided by commercial network operators. One grid scheduling system is developed to co-allocate computing and network resources with advance reservations through web service interfaces using the Grid Resource Scheduler (GRS), the Network Resource Management System (NRM), which is capable of GMPLS network resource management. This is a general framework for resource co-allocation, which does not pay enough attention to characteristics of data streaming applications, such as sustaining and controlled data provision.

6 CONCLUSIONS AND FUTURE WORK

Data streaming applications bring new challenges to grid resource management and scheduling, such as requiring real-time data provision and integrated resource allocation schemes. Different from existing resource management and scheduling schemes that only focus on computational resources, the system proposed in this paper takes computational, storage and network resources into account simultaneously and makes integrated management and scheduling schemes, which are proved to be feasible with improved performance and scalability.

The work described in this paper is mainly focused on resource utilization and overall data throughput of the system. Future work will address quality of service

(QoS) issues required by data streaming applications. Scheduling for data streaming pipelines will also be considered, which is more complicated with requirement of balancing among multiple stages and appropriate data provision. Ongoing work also includes scheduling data sharing among multiple data streaming applications to further improve system performance.

Acknowledgement

This work is supported by National Science Foundation of China (grant No. 60803017) and Ministry of Science and Technology of China under National 973 Basic Research Program (grants No. 2011CB302505 and No. 2011CB302805) and National 863 High-tech R & D Program (grants No. 2008AA01Z118 and No. 2008-BAH32B03).

Junwei Cao would like to express his gratitude to Professor Erik Katsavounidis of LIGO (Laser Interferometer Gravitational-wave Observatory) Laboratory at Massachusetts Institute of Technology for his support on LIGO Scientific Collaboration.

REFERENCES

- [1] FOSTER, I.—KESSELMAN, C.: *The Grid: Blueprint for a New Computing Infrastructure*. Morgan Kaufmann, San Francisco 1998.
- [2] DEELMAN, E.—KESSELMAN, C.—MEHTA, G.—MESHKAT, L.—PEARLMAN, L.—BLACKBURN, K.—EHRENS, P.—LAZZARINI, A.—WILLIAMS, R.—KORANDA, S.: GriPhyN and LIGO, Building a Virtual Data Grid for Gravitational Wave Scientists. Proc. 11th IEEE Int. Symp. on High Performance Distributed Computing, 2002, pp. 225–234.
- [3] CAO, J.—KATSAVOUNIDIS, E.—ZWEIZIG, J.: Grid Enabled LIGO Data Monitoring. Proc. IEEE/ACM Supercomputing Conf., Seattle, WA, USA, 2005.
- [4] PORDES, R.: For the Open Science Grid Consortium. The Open Science Grid, Proc. Computing in High Energy and Nuclear Physics Conf., Interlaken, Switzerland 2004.
- [5] HOLLAND, J. H.: *Adaptation in Natural and Artificial Systems*. University of Michigan Press 1975.
- [6] FOSTER, I.—KESSELMAN, C.: Globus: A Metacomputing Infrastructure Toolkit. Int. J. Supercomputer Applications, Vol. 11, 1997, No. 2, pp. 115–128.
- [7] KAR, K.—SARKAR, S.—TASSIULAS, L.: A Simple Rate Control Algorithm for Maximizing Total User Utility. Proc INFOCOM 2001.
- [8] ALLCOCK, B.—BESTER, J.—BRESNAHAN, J.—CHERVENAK, A. L.—FOSTER, I.—KESSELMAN, C.—MEDER, S.—NEFEDOVA, V.—QUESNAL, D.—TUECKE, S.: Data Management and Transfer in High Performance Computational Grid Environments. Parallel Computing, Vol. 28, 2002, No. 5, pp. 749–771.
- [9] ABADI, D.—CARNEY, D.—CETINTEMEL, U.—CHERNIACK, M.—CONVEY, C.—LEE, S.—STONEBRAKER, M.—TATBUL, N.—ZDONIK, S.: Aurora: A New Model

- and Architecture for Data Stream Management. *VLDB Journal*, Vol. 12, 2003, No. 2, pp. 120–139.
- [10] BALAZINSKA, M.—BALAKRISHNAN, H.—STONEBRAKER, M.: Contract-based Load Management in Federated Distributed Systems. Proc. 1st Symp. on Networked Systems Design and Implementation (NSDI), San Francisco, CA, March 2004.
 - [11] CHANDRASEKARAN, S.—COOPER, O.—DESHPANDE, A.—FRANKLIN, M. J.—HELLERSTEIN, J. M.—HONG, W.—KRISHNAMURTHY, S.—MADDEN, S. R.—REISS, F.—SHAH, M. A.: TelegraphCQ: Continuous Dataflow Processing. Proc. ACM SIGMOD Int. Conf. on Management of Data (SIGMOD '03), 2003.
 - [12] CHERNIACK, M.—BALAKRISHNAN, H.—BALAZINSKA, M.—CARNEY, D.—CETINTEMEL, U.—XING, Y.—ZDONIK, S.: Scalable Distributed Stream Processing. Proc. 1st Biennial Conf. on Innovative Data Systems Research (CIDR '03), Asilomar, CA, January 2003.
 - [13] NABRZYSKI, J.—SCHOPF, J. M.—WEGLARZ, J.: Grid Resource Management: State of the Art and Future Trends. Kluwer Academic Publishers, Sep. 2003.
 - [14] CHAPIN, S. J.—KATRAMATOS, D.—KARPOVICH, J.—GRIMSHAW, A. S.: The Legion Resource Management System. Job Scheduling Strategies for Parallel Processing, Springer Verlag 1999, pp. 162–178.
 - [15] BUYYA, R.—ABRAMSON, D.—GIDDY, J.: Nimrod/G: An Architecture for a Resource Management and Scheduling System in a Global Computational Grid. Proc. High Performance Computing ASIA, 2000.
 - [16] LITZKOW, M.—LIVNY, M.—MUTKA, M.: Condor – A Hunter of Idle Workstations. Proc. 8th Int. Conf. on Distributed Computing Systems, 1988, pp. 104–111.
 - [17] CAO, J.—JARVIS, S. A.—SAINI, S.—KERBYSON, D. J.—NUDD, G. R.: ARMS: An Agent-based Resource Management System for Grid Computing. Scientific Programming, Special Issue on Grid Computing, Vol. 10, 2002, No. 2, pp. 135–148.
 - [18] BASU, S.—TALWAR, V.—AGARWALLA, B.—KUMAR, R.: I-GASP: Interactive Grid Environment Provided by Application Service Providers. Proc. 1st Int. Conf. on Web Services (ICWS '03), Las Vegas, USA 2003.
 - [19] TALWAR, V.—AGARWALLA, B.—BASU, S.—KUMAR, R.: Architecture for Resource Allocation Services Supporting Remote Desktop Sessions in Utility Grids. Proc. 2nd Int. Workshop on Middleware for Grid Computing (MGC 2004), Toronto, Canada, Oct. 2004.
 - [20] CAO, J.—ZIMMERMANN, F.: Queue Scheduling and Advance Reservations with COSY. Proc. 18th IEEE Int. Parallel & Distributed Processing Symp., Santa Fe, NM, USA 2004, p. 63.
 - [21] CHEN, L.—REDDY, K.—AGRAWAL, G.: GATES: A Grid Based Middleware for Processing Distributed Data Streams. Proc. 13th IEEE Int'l. Sym. on High Performance Distributed Computing (HPDC-13), Honolulu, Hawaii (USA), June 4–6, 2004.
 - [22] CHEN, L.—AGRAWAL, G.: Resource Allocation in a Middleware for Streaming Data. Proc. 2nd Workshop on Middleware for Grid Computing (MGC '04), Toronto, Canada, October 18, 2004.

- [23] CHEN, L.—AGRAWAL, G.: A Static Resource Allocation Framework for Grid-based Streaming Applications. *Concurrency and Computation: Practice and Experience*, Vol. 18, 2006, pp. 653–666.
- [24] AGARWALLA, B.—AHMED, N.—HILLEY, D.—RAMACHANDRAN, U.: Streamline: A Scheduling Heuristic for Streaming Applications on the Grid. *Proc. SPIE Multimedia Computing and Networking*, Vol. 6071, 2006.
- [25] AGARWALLA, B.—AHMED, N.—HILLEY, D.—RAMACHANDRAN, U.: Streamline: Scheduling Streaming Applications in a Wide Area Environment. *Multimedia Systems*, Vol. 13, 2007, No. 1, pp. 69–85.
- [26] DEELMAN, E.—BLYTHE, J.—GIL, Y.—KESSELMAN, C.—MEHTA, G.—VAHI, K. et al.: Mapping Abstract Complex Workflows onto Grid Environments. *J. Grid Computing*, Vol. 1, 2003, No. 1, pp. 25–39.
- [27] RAMAKRISHNAN, A.—SINGH, G.—ZHAO, G.—DEELMAN, E.—SAKELLARIOU, R.—VAHI, K.—BLACKBURN, K.—MEYERS, D.—SAMIDI, M.: Scheduling Data-Intensive Workflow onto Storage-Constrained Distributed Resources. *Proc. 7th IEEE Int. Symp. on Cluster Computing and the Grid*, Rio de Janeiro, Brazil 2007, pp. 401–409.
- [28] BATTISTILLI, L. et al.: An Architecture for Co-allocating Network, Compute, and other Grid Resources for High-End Applications. *Proc. Int. Symp. on High Capacity Optical Networks and Enabling Technologies*, 2007, pp. 1–8.
- [29] TAKEFUSA, A. et al.: Coordination of a Grid Scheduler and Lambda Path Service over GMPLS. *Future Generation Computer Systems*, Vol. 22, 2006, No. 8, pp. 868–875.



Wen ZHANG works in Chongqing Military Delegate Bureau, General Armament Department of PLA. He received his Ph.D. in control engineering and applications from Department of Automation, Tsinghua University, in 2010. His research is focused on grid data streaming.



Junwei CAO is currently a Professor and Assistant Dean, Research Institute of Information Technology, Tsinghua University, China. He was a research scientist at MIT LIGO Laboratory and NEC Laboratories Europe. He received the Ph.D. in computer science from University of Warwick, UK, in 2001. He is a senior member of the IEEE Computer Society and a member of the ACM and CCF.



Yisheng ZHONG is currently a Professor at Department of Automation, Tsinghua University, China. He received the Ph.D. in electrical engineering from Hokkaido University (Japan) in 1988. His research interests include control theory, complexity, etc.



Lianchen LIU is currently an Associate Professor at Department of Automation, Tsinghua University, China. He received the Ph.D. from NanKai University (China). His research interests include large scale scientific resource sharing, distributed computing, etc.



Cheng WU is a Professor at Department of Automation, Tsinghua University (China), Director of National CIMS Engineering Research Center, and member of Chinese Academy of Engineering. He received his B.Sc. and M.Sc. from Department of Automation, Tsinghua University in 1962 and 1966, respectively. His research interests include complex manufacturing system scheduling, grid/cloud applications, etc.