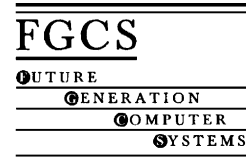




ELSEVIER

Available at
www.ComputerScienceWeb.com
POWERED BY SCIENCE @ DIRECT®

Future Generation Computer Systems 19 (2003) 945–955



www.elsevier.com/locate/future

Experimental studies using photonic data services at IGrid 2002

Robert L. Grossman^{a,*}, Yunhong Gu^a, Don Hamelburg^a, Dave Hanley^a,
Xinwei Hong^a, Jorge Levera^a, Dave Lillethun^b, Marco Mazzucco^a,
Joe Mambretti^b, Jeremy Weinberger^b

^a University of Illinois at Chicago, Chicago, IL, USA

^b Northwestern University, Evanston, IL, USA

Abstract

We describe an architecture for remote and distributed data intensive applications that integrates optical path services, network protocol services for high performance data transport, and data services for remote data analysis and distributed data mining. We also present experimental evidence using geoscience data that this architecture scales to long haul, high performance networks.

© 2003 Elsevier Science B.V. All rights reserved.

Keywords: High performance networks; Distributed data mining; Data webs; TCP; UDP

1. Introduction

A fundamental research challenge is to develop systems for remote data analysis and distributed data mining which scale to very large data sets. The data may be at rest in the sense that it resides on remote disks and tapes or it may be in motion in the sense that it is collected and streamed from a remote instrument.

The analysis and mining of this type of data is difficult for several reasons:

- (1) For fixed window size, the bandwidth of reliable network protocols such as TCP, which requires acknowledgments for each packet, is roughly approximated by $1/RTT$, where RTT is the round trip time for a packet to travel between the client and server. For an application across the Atlantic this gives a rough limit of approximately 5 Mb/s, even

if the path has a capacity many times greater, such as 1 Gb/s link.

- (2) In the last few years, we have gained a better understanding of the primitives required to integrate data mining with databases. On the other hand, we do not yet have a good understanding of the primitives required for remote data analysis and distributed data mining.

In this paper, we introduce an architecture for remote data analysis and distributed data mining that integrates services to set up optical paths, network protocols designed for high performance networks, and data services providing simple primitives supporting the remote analysis and distributed mining of large data sets. We also describe experiments showing the speedup gained with this approach for some typical data mining algorithms such as computing simple correlations for streaming geoscience data.

In Section 2, we describe related work. In Section 3, we describe the architecture we introduce called photonic data services (PDS). In Sections 4–6, we

* Corresponding author.

E-mail addresses: grossman@uic.edu (R.L. Grossman), marco@dmg.org (M. Mazzucco).

describe the three main layers of PDS. In Section 7, we describe two experimental studies involving PDS applications. Section 8 is the summary and conclusion.

A preliminary version of this paper describing a prior version of the architecture and experiments with an OC-12 network appeared in the conference proceedings [1].

2. Background and related work

In this paper, we are concerned with supporting remote data analysis and distributed data mining applications with high performance data transport services. In addition, many applications will also require high performance compute services, which we do not address. Today, these would typically be provided by local compute clusters or by virtual compute clusters accessed via a computational grid [2] or computational middleware architecture [3].

It has long been recognized that TCP is not appropriate for high performance applications on long haul networks. The reason is simple: the TCP protocol requires acknowledgment of each packet. This limits the bandwidth to be a function of $1/RTT$, where RTT is the time required to send a packet and receive an acknowledgment [4, p. 186].

One approach to improving TCP performance for data intensive applications is to adjust the TCP window size to be the product of the bandwidth and the RTT delay of the network [4]. This approach requires modifying and tuning the kernel of each of the operating systems transporting the packets and ensuring that the networking hardware can support these large or jumbo packets.

Another approach to overcoming the limitations of TCP is to stripe TCP over several standard TCP network connections. In contrast to the first approach, this can be done at the data middleware or application level. This approach has been implemented several times, including Pockets [5] and GridFTP [6]. It has been observed that the performance of striped TCP begins to level off after about 12–20 sockets, effectively limiting its usefulness to OC-3 or OC-12 networks.

There have been three main architectural approaches to date for distributed data mining: agent based systems, data grid-based systems, and data web based systems. We consider each in turn.

The first approach is to use agents over commodity networks to move data, remotely control the data mining algorithms at the different sites, and collect the intermediate results and models. Systems with this architecture include the JAM system developed by Stolfo et al. [7], the BODHI system developed by Kargupta et al. [8], the Kensington system developed by Guo and coworkers [9], and the Papyrus system developed by Grossman et al. [10].

The second approach is to use grid-based middleware. Systems with this architecture include those developed by Parthasarathy and Subramonian [11], Moore et al. [12], and Grossman et al. [10]. More recently, Globus has emerged as the dominant middleware for working with distributed clusters [2]. The Globus infrastructure for data intensive computing is called the data grid, and includes services for parallel TCP striping (GridFTP), and data replication services (Globus Replica Catalog and Globus Replica Management) [13].

Other grid middleware services that have been used for data mining include the DataCutter developed by Saltz and coworkers [14] and Discovery Net developed by Patrick and Guo [15]. For example, Du and Agrawal [16] recently used the DataCutter for some distributed data mining experiments.

The third approach and the one used in this paper is to use data webs, which are web based infrastructures for data [17]. Unlike grid middleware which is built over authentication, authorization and access (AAA) control mechanisms for rationing and scheduling presumably scarce high performance computing resources [2], data webs are built using W3C standards and emerging standards for web services and packaging (SOAP and XML). Data webs in contrast to data grids are designed to encourage the open sharing of data resources without AAA controls, in the same way that the web today encourages the sharing of document resources without AAA controls [18].

For small data sets, data webs use W3C standards and emerging standards to manage both the data and meta-data. These include HTTP, DSTP (DataSpace transfer protocol) and other emerging standards for transport [19], and SOAP and XML for packaging [18]. For large data sets, this infrastructure is used just for the *meta-data*, while specialized network protocols and data services (the PDS described below) are used to manage the *data* itself. Providing separate

mechanisms for control paths and data paths is an old idea in high performance computing going back to at least the IBM High Performance Storage System (HPSS). Developing the appropriate data web services and protocols to work with large remote and distributed data is a fundamental research challenge.

3. Photonic data services

Today data intensive applications working with remote and distributed data are generally based upon standard networking (IP) and transport (TCP) protocols. For data mining applications running on commodity networks analyzing small data sets these protocols work extremely well. Data mining applications involving large, distributed data sets have generally used specialized networks such as NSF's vBNS network or the Internet 2's Abilene network. In practice, very large bandwidth applications have to be scheduled on these networks and require the use of specialized transport protocols [5].

This paper describes applications based upon the following specialized services.

Intelligent signaling for paths by applications. As optical networking architectures become more common, a new possibility is emerging. A bandwidth demanding application can request an optical connection between the data sources and the data sinks for a specific application. More specifically, the application can request the setup, status and tear down of the required optical paths. Clearly there is a cross-over point; for short transfers of small data, TCP is clearly preferable, while for long transfers of very large data, a dedicated optical path is clearly preferable.

With today's MEMS technology, reconfiguration of a light path in an optical add/drop multiplexer (OADMs) or optical cross-connects (OXC) can be done in milliseconds. New cross-connect technologies such as semiconductor optical amplifiers (SOAs) would shorten this to nanoseconds [20]. Today, as we will see below, applications requiring moving >0.5 GB for remote data analysis and distributed data mining can benefit from requesting specialized lambda paths. In the future, the cross-over point will be much smaller.

Specialized network protocols. It is clear that TCP does not currently perform well for moving large data

sets over long distances. Recently there has been a flurry of activity investigating alternative protocols [5,6,13,21,22]. Our assumption is that one or more of these alternative protocols will be commonly used by data intensive applications, supplementing TCP when required.

Specialized data services. Moving data is different than moving bits. In addition, we assume that data intensive applications will also make use of specialized protocols and services for working with data, services that support meta-data operations, record and attribute selection, missing data, sampling, and related data services.

In this paper, we introduce the idea of integrating: (1) specialized photonic path services; (2) high performance network protocols; (3) high performance data services providing data mining primitives for remote data analysis and distributed data mining. We call these integrated services PDS.

As an example, in Section 7.3, we describe a distributed data mining application in which 1.8 GB of vegetation data over a region specified by latitude and longitude coordinates is correlated with 1.8 GB of climate data over the same region. Both data sets are in the US in different locations, but the client doing the correlation is in Amsterdam.

Assume that both data sources are connected to the client by 10 Gb/s links. Today, the data would be moved to a common location using a standard network protocol such as TCP, merged, and then correlated. Merely moving the data across the Atlantic takes over 3000 s, as we will see in Section 7.3.

Our assumption is that in the future dedicated links will not be common, but the ability to set up specialized photonic paths on a per application basis will be. Using the PDS described below, a photonic path can be set up in less than 1 min and the two 1.8 GB streams transported and merged in less than 70 s, as we will see in Section 7.3. As the path services software matures, we expect the setup time to be reduced substantially, so that a data mining computation that today requires approximately 1 h could be done in approximately 1 min.

For the purposes here, the layered network model we use is a simple extension of the standard 5-layer model in which we split the top layer into two. One of these provides specialized data services for remote data analysis and distributed data mining and the other

is the top application layer. Here is the layered model we are using:

- (1) *Physical links.* We assume that the physical links are provided by multichannel wavelength-division multiplexed (WDM) communications, as well as by Ethernet, and other technologies.
- (2) *Path services layer.* We assume that there are services allowing us to set up paths between devices, tear down paths, check the status of paths, setup routing, etc.
- (3) *Internet layer.* This layer provides a common network addressing and routing across multiple networks. For our applications, we use the Internet Protocol (IP) in this layer.
- (4) *Network protocol services layer.* We assume that there are transport services including TCP, UDP, and other more specialized protocols providing high performance over the paths. Our applications use specialized high performance protocols in this layer.
- (5) *Data services layer.* We assume that there are standard services for moving data such as SOAP-based web services, as well as more specialized data services designed for performance networks.
- (6) *Application layer.* We assume that the remote data analysis and distributed data mining applications can request standard and specialized network services depending upon the applications requirements.

In the next three sections, we describe the three service layers we have implemented and integrated to create PDS to support data analysis and data mining. Our implementation of the path services (layer 2) is called ODIN [23]; our implementation of the network protocol services (layer 4) is called Simple Available Bandwidth Utilization Library (SABUL) [22]; our implementation of data services (layer 5) uses high performance implementations of the DSTP we have developed [24]. Using these three services, we developed an application for merging two high bandwidth data streams of earth science data employing a join algorithm called the continuously generated merge (CGM) [25].

The work described in this paper is the first time we have integrated these three service layers and performed experimental studies using them.

Integrating PDS with web services is straightforward and will be described in a forthcoming publication. To do this, we simply added two additional layers between PDS layers 5 and 6: one for the description of data services (for example, WSDL) and the other for the discovery of data services (for example, UDDI) [18].

4. Path services—ODIN

The path services used in PDS are called the Optical Dynamic Intelligent Network Service Layer or ODIN [23]. ODIN receives requests for circuits by applications, which for PDS can be from layer 4 or 5 services, and contacts the required network switches, including both optical-domain DWDM switches and traditional Ethernet switches and IP routers, to set up the circuits.

ODIN consists of two sub-systems: one, called the TeraScale High Performance Optical Resource Regulator or THOR, interfaces to the optical fabric; while the other, called the Dynamic Ethernet Intelligent Transit Interface or DEITI, interfaces to the traditional Ethernet/IP fabric. We now describe these systems following [23].

ODIN is designed to dynamically provision and control global light paths. The ODIN subsystem THOR is based on new signaling methods for dynamically provisioning light paths. These light paths can be used to create Optical Virtual Private Networks (OVPNs), as well as to extend these light paths to edge resources through other types of dynamically provisioned paths, such as vLANs.

ODIN can be used to establish “services on demand” and, as noted, not only dynamically allocated light paths, but also dynamically allocated transient (or permanent) OVPNs. In part, ODIN accomplishes its functions through interactions with lower layer optical services such as THOR.

ODIN is designed to work within a single administrative domain and provide a predefined set of path services. Since the only allocation of paths is through ODIN, it has complete knowledge of the topology and current resource allocations within the administrative domain. ODIN accepts requests for path services from services and applications over the network. When resources are allocated to fulfill requests, ODIN communicates with the requisite network switches to

configure them as required. These switches can be optical-domain DWDM switches, Ethernet switches and/or IP routers.

With the current implementation of PDS, layer 4 or 5 services determine whether a photonic path is required. If so, a request is sent to the ODIN server. The ODIN server: (1) finds the closest optical endpoints for the specific source and destination; (2) establishes the optical connection; (3) activates other related network components to provide the required upper layer services, such as MPLS and Tagged VLAN switching.

5. Network protocol services—SABUL

In this section, we follow [22] and describe a network protocol we have developed called SABUL. SABUL is designed for high performance data transfer and is especially useful on long haul networks.

The idea behind SABUL is simple. SABUL combines the UDP protocol in order to send data at a high rate with the TCP protocol in order to do this in a reliable fashion. UDP has no flow control, rate control, or reliable transmission mechanisms. SABUL implements these control functions in a separate TCP control channel. This approach is in contrast to the approach of other high performance protocols such as NETBLT [26], which combine the data and control channels.

In SABUL, the packets on the UDP channel consist of the usual UDP header plus a 32 bit field for a sequence number. On the TCP channel, each packet consists of a list of lost data packets, a field stating the requested data rate, and a field reserved to report the state of the receiver's available buffer size. We define the *communication state* information to be the information contained in these TCP packets.

The flow is assumed to be unidirectional. Data is sent to the receiver over the UDP channel, while current communication state information is sent over the TCP channel, from the receiver to the sender. Since the communication state information is passed over TCP, its arrival is ensured; since the amount of this information is relatively small, it has a negligible effect on the overall performance of SABUL.

One of the advantages of SABUL is its continuous updating of state information. In contrast, NETBLT uses a mechanism that sends buffers of data at a

fixed rate. At the end of transmission of each buffer, the receiving side of NETBLT sends the sender a list of packets that were lost in the transmission of this buffer. The sender then resends these packets; the process continues until all packets in the buffer are accounted for. Then the next buffer can be transmitted by NETBLT. NETBLT needs to block until all packets are accounted for on the sending side before sending another buffer. This process can be further delayed since packet loss information is transmitted unreliably by the receiver to the sender (since this information is sent over UDP). Another deficit of NETBLT is that it needs to wait for at least one round trip time to get each update of packets lost.

In SABUL, however, each time the receiver notices at least one missing packet, it uses the TCP channel to transmit to the sender a list of packets that were lost. It does not have to block the sending of packets over the UDP channel to wait for an incoming packet containing the communication state information. This allows for changing the rate and flow of data, and retransmission of any missing packets during the transmission of the data. The list of missing packets is updated every time a missing packet is received. If during a predefined amount of time no packet was lost, and thus no transmission sent to the sender on the TCP channel, the receiver sends a notification of this fact to the sender with communication state information. This allows the sender to empty its buffer of packets that have successfully been received and adjust the rate and flow if necessary.

6. Data services—DSTP

In this section we follow [19] and describe data services designed to be component services or primitives for distributed data mining applications. The services are part of a protocol called DSTP, which we have developed. Access to DSTP data can also be gained through proxy servers, such as SOAP-based servers. The experiments below used high performance DSTP servers, which we have developed [24].

Distributed columns of numerical data. The data model for DSTP Version 2 is simple. Data is divided into rows (data records) and columns (data fields or data attributes). Both may be distributed over the web. Access to the data itself is through a DSTP

server, or through a proxy service, such as through a SOAP-based server. Physically, the data itself may be stored as files, in databases, or using other specialized storage mechanisms. Logically, data is just a distributed collection of columns. Version 3 of DSTP may support a more complex data model.

Universal correlation keys. A Universal Correlation Key (UCK) is a globally unique id (GUID) and is used for relating columns of data on two different DSTP servers. Each column of data is associated with at least one column of UCKs.

Multi-dimensional UCKs. UCKs may be combined to provide multi-dimensional keys. This is essential for working with scientific and engineering data, such as the geoscience data used in the experiments below. For example, this data uses latitude and longitude as the UCKs.

Column based meta-data. Associated with each column of data is attribute meta-data and with each data set (a collection of columns) data set meta-data. DSTP applications may or may not use this meta-data. On the other hand, this meta-data is essential for building and deploying statistical models. DSTP servers provide a simple mechanism for associating meta-data to columns and collections of columns. DSTP applications often access meta-data using SOAP/XML based web services.

UCKs enable distributed columns to be correlated in the following fashion: pairs (k_i, x_i) , where k_i is a UCK value and x_i an attribute value, on DSTP server 1 can be combined with pairs (k_j, y_j) on DSTP server 2 to produce a table (x_k, y_k) in a DSTP client. The DSTP client can then, for example, find a function $y = f(x)$ relating x and y . This simple mechanism of distributed columns identified by UCKs (perhaps vector valued) is sufficient information for many data mining algorithms.

Depending upon the request, DSTP servers may return one or more columns, one or more rows, or entire tables. DSTP uses XML to describe the meta-data. DSTP applications can also access small data sets using XML so that DSTP is compatible with SOAP. On the other hand, for efficiency and scalability, the default is for data *itself* to be transmitted as records delimited by carriage returns, with fields delimited by commas. The DSTP client may also indicate that a specialized high performance protocol such as SABUL should be used for the data channel. To sum-

marize, the DSTP protocol uses XML for meta-data and small data, while data is typically streamed, with large amounts of data streamed using SABUL or other high performance network protocols.

The DSTP protocol includes commands for retrieving meta-data, retrieving UCKs, retrieving data and subsets of data, and mechanisms for sampling, working with missing data, and merging by UCKs.

We note that some of the IGrid demonstrations, such as browsing data from the Protein Data Bank, required 30 s or more when using SOAP/XML for the meta-data while less than 1 s when using specialized DSTP streaming protocols. This is due to the overhead to parse XML, which becomes more of an application burden as the number of attributes grows.

7. Experimental studies

7.1. IGrid 2002 testbed

We assume that our network consists of dense WDM optical devices together with standard Ethernet/IP devices. For our experiments we used the Chicago area OMNInet [27] and the global Terra Wide Data Mining (TWDM) testbed [28].

OMNInet is an optical networking testbed deployed in the Chicago metropolitan area. OMNInet currently provides 1 and 10 GE services between Northwestern, The University of Illinois at Chicago, and the StarLight facility in Chicago. OMNInet is operated by a research consortium consisting of iCAIR at Northwestern, the Electronic Visualization Laboratory at the University of Illinois at Chicago, Argonne National Laboratory, SBC, and Nortel.

The TWDM testbed is a testbed built on top of DataSpace for the remote analysis, distributed mining, and real time exploration of scientific, engineering, business, and other complex data. Currently, the TWDM testbed consists of five geographically distributed workstation clusters linked by optical networks with StarLight in Chicago as the optical interchange point. These sites include StarLight itself, the Laboratory for Advanced Computing at UIC, iCAIR at Northwestern University, SARA in Amsterdam, and CANARIE in Ottawa. SARA is connected to StarLight via the Netherlands' Surfnet network,

Ottawa is connected to Starlight via Canda's CANARIE network, and UIC and iCAIR are connected via OMNinet.

The setup for the experiments in the next two sections was as follows. Three-node DSTP clusters were located at the SARA research facility in Amsterdam, at the University of Illinois at Chicago, and at the StarLight Facility in Chicago. StarLight and the University of Illinois at Chicago are located several miles apart. The SARA cluster and the StarLight cluster were connected via a 10 Gb/s Surfnet link. The UIC cluster and StarLight cluster were connected via a 10 Gb/s OMNinet link.

The machine in Amsterdam was a dual P4, 1700 MHz, with 512 M RAM. The machines in Chicago were dual PIII, 1000 MHz, with 512 M RAM. The machines were all running Linux, with the 2.4.x kernels.

7.2. PDS application: lambda FTP

Our first series of experiments measured an application we developed called Lambda FTP, a high performance implementation of FTP using SABUL. The testing was done between two three-node clusters, one located at StarLight in Chicago and one located at SARA in Amsterdam. The results are reported in the table below. Standard TCP provided about 4.36 Mb/s of bandwidth, while each of the three SABUL streams provided about 900 Mb/s of bandwidth, so that the aggregate SABUL bandwidth was about 2.7 Gb/s between the two three-node clusters.

We note that DSTP servers can stream data using SABUL so that these performance numbers are also applicable to streaming DSTP data. In general, though, DSTP clients perform some type of computation so that actual bandwidth is dependent upon the particular data web application. An illustration of this is contained in the next section.

File transfer between Chicago and Amsterdam (Mb/s)

TCP stream	SABUL stream 1	SABUL stream 2	SABUL stream 3	Aggregate SABUL stream
4.36	902.8	902.9	907.1	2712.8

7.3. PDS application: lambda joins

Our second series of experiments involved testing a basic operation in distributed data mining, merging two data streams by a common key. For this series of experiments, we merged two NCAR data sets [29]. One of the data sets was located at UIC in Chicago and one was located at SARA in Amsterdam. The merging was done at StarLight in Chicago.

The two DSTP streams of data were merged by an algorithm called the CGM, which we developed for merging two high bandwidth data streams [25]. In this case, we merged the two streams using latitude, longitude, and time as the common vector valued key. Once distributed streaming data has been merged in this way, simple statistical counts using a finite buffer can be done in a variety of ways to support the interactive exploration of large remote data sets [30].

Merging distributed data streams by common keys is a basic operation in distributed data mining. Other examples include merging satellite imaging data of different modalities by latitude/longitude and merging network route dumps by source address.

We now briefly review the CGM algorithm following [25]. In the CGM algorithm we assume the data is partially presorted. Without loss of generality, assume there are two data streams, A and B, being drawn into a client in approximately ascending order and we are trying to merge on one UCK. The CGM algorithm depends upon two parameters: a parameter N determining the number of records in a window, which is used to buffer the streaming data, and N_h , the number of entries in two auxiliary hash tables. The algorithm has an even step and an odd step. The even steps of the algorithm are as follows:

- (1) The client grabs some fixed number of records N , from both stream A and stream B and places them in window A and window B, respectively (each has room for exactly N records).
- (2) A hash is done on the value of each UCK in window A and the record is placed in the appropriate location in hash table A, overwriting any previous record.
- (3) A hash is done on the value of each UCK in window B and if the value hashes to an occupied location in hash table A, both the records are merged. If the value does not hash to an occupied location in hash

table A, then the record is placed in the appropriate location in hash table B, overwriting any previous record.

In the odd steps of the algorithm, the above algorithm is executed, but reversing the roles of A and B.

The results in the table below are from the CGM algorithm running using TCP as the network protocol and DSTP as the data service protocol. Each data stream was 300 MB in size. The CGM algorithm used a hash table size of 50,000 and a window size of 10,000. The data was randomized by replacing every n th row (for example, for 10% every 10th row) with a random row that was within 50,000 lines of the current row.

As can be seen from the table, the average speed varied between 4 and 5 Mb/s.

Rand (%)	Match (%)	Time (s)	Data rate (Mb/s)
Merging two data streams without PDS			
2	96.6	513	4.68
10	89.9	540	4.44
20	81.5	531	4.52
33	73.1	563	4.26

The next two results below are from the CGM algorithm running SABUL as the network protocol and DSTP as the data service protocol. One DSTP server per cluster was used to send and receive data, so that the maximum bandwidth between the two clusters was 1 Gb/s. *The data size this time was 1.8 GB so that in total 3.6 GB of data were merged by the algorithm.* The average speed varies between 400 and 500 Mb/s. This means that CGM over SABUL was about $600\times$ faster on average, since the amount of data was $6\times$ greater and the elapsed time was about $100\times$ less.

When testing the algorithm we realized the largest single affect on the performance of the merge was the length of the record. The longer the record size the memory copying required, the greater the merge time. To illustrate this we ran two tests. In the first, both data files containing one UCK and one attribute; in the second, both data files containing one UCK and seven attributes.

Rand (%)	Match (%)	Time (s)	Data rate (Mb/s)
Merging two data streams with PDS (one attribute)			
2	99	53.3	550
10	91	52.4	550
20	83	56.2	512
33	78	54.6	527
Merging two data streams with PDS (seven attributes)			
2	99	66.3	434
10	92	65.7	438
20	82	64.2	449
33	79	65.1	442

8. Summary and conclusion

In this paper, we have introduced an architecture called PDS which integrates path services, network protocol services, and data services. With intelligent path services, distributed data mining applications can intelligently signal for a special photonic path, use this for distributed data mining, and then release it for use by other applications. With high performance network protocols, data mining applications can work effectively with remote Gigabyte size data sets over high performance networks. These types of protocols are sometimes several hundred times faster than traditional protocols over the same high performance networks. With specialized data services such as streaming merges, distributed data mining services can effectively correlate distributed Gigabyte size data sets.

In this paper, we have provided experimental evidence that this approach is practical and useful and that our implementations scale to remote Gigabyte size data sets that are distributed over thousands of miles. Compared to current implementations of data mining primitives for merging two data streams and computing counts, our PDS services are significantly faster. For example, to stream two 1.8 GB data streams of geoscience data across the Atlantic, merge the results using latitude, longitude and time as keys, and compute simple counts required over 1 h with conventional services and less than 1 min using the PDS described in this paper. We emphasize that both experiments used the same high performance network.

We believe that the work described here is novel for the following reasons:

- (1) This is the first description of an *architecture* for integrating path services for photonic networks, network protocols designed for high performance networks, and data services supporting primitives to facilitate remote data analysis and distributed data mining. Integrated services such as these can provide the foundation for scaling distributed data mining to large data sets. We call this architecture PDS.
- (2) This is the first integration that allows applications (in our case data services middleware) to signal to optical networks requests to set up, check the status, and tear down photonic paths. The ability to configure paths and related services on an as-needed basis is essential if these services will one day move from today's experimental networks to tomorrow's production networks. Think of this as intelligent application signaling.
- (3) This is the first experimental study demonstrating the feasibility of the *distributed* mining of Gigabyte size data sets that are separated by thousands of miles and over 100 ms in packet round trip time.

There are three main areas of work for the future:

- (1) First, we plan on scaling photonic path services to multiple administrative domains. The experiments reported here used photonic paths services within the OMNInet domain, but paths between OMNInet and Surfnet were set up by hand. This is a PDS layer 2 research problem.
- (2) Second, we plan on investigating the performance of SABUL on larger clusters. The experiments reported here were limited to three-node clusters. This is a PDS layer 4 research problem.
- (3) Third, we plan on developing additional remote data analysis and distributed data mining algorithms using DSTP primitives. This is a PDS layer 5 research problem.

Acknowledgements

This work was partially supported by the National Science Foundation.

References

- [1] R.L. Grossman, Y. Gu, D. Hanley, X. Hong, J. Levera, M. Mazzucco, D. Lillethun, J. Mambretti, J. Weinberger, Photonic data services: integrating path, network and data services to support next generation data mining applications, in: Proceedings of the Next Generation Data Mining Conference, Kluwer Academic Publishers, Dordrecht, 2003.
- [2] I. Foster, C. Kesselman, The Grid: Blueprint for a New Computing Infrastructure, Morgan Kaufmann, San Francisco, CA, 1999.
- [3] NSF middleware initiative, September 2, 2002. <http://www.nsf-middleware.org>.
- [4] J. Walrand, P. Varaiya, High Performance Communication Networks, Morgan Kaufmann, San Francisco, CA, 2000.
- [5] R.L. Grossman, H. Sivakumar, S. Bailey, Pockets: the case for application-level network striping for data intensive applications using high speed wide area networks, in: Supercomputing, IEEE and ACM, 2000.
- [6] A. Chervenak, I. Foster, C. Kesselman, S. Tuecke, Protocols and services for distributed data-intensive science, in: ACAT2000 Proceedings, 2000, pp. 161–163.
- [7] S. Stolfo, A.L. Prodromidis, P.K. Chan, Jam: Java agents for meta-learning over distributed databases, in: D. Heckerman, H. Mannila, D. Pregibon, R. Uthurusamy (Eds.), Proceedings of the Third International Conference on the Knowledge Discovery and Data Mining, AAAI Press, Menlo Park, CA, 1997.
- [8] H. Kargupta, I. Hamzaoglu, B. Stafford, Scalable, distributed data mining using an agent based architecture, in: D. Heckerman, H. Mannila, D. Pregibon, R. Uthurusamy (Eds.), Proceedings of the Third International Conference on the Knowledge Discovery and Data Mining, AAAI Press, Menlo Park, CA, 1997, pp. 211–214.
- [9] J. Darlington, Y. Guo, J. Sutiwaraphun, H.W. To, Parallel induction algorithms for data mining, Lecture Notes in Computer Science 1280.
- [10] R.L. Grossman, S. Bailey, A. Ramu, B. Malhi, H. Sivakumar, A. Turinsky, Papyrus: a system for data mining over local and wide area clusters and super-clusters, in: Proceedings of Supercomputing 1999, IEEE and ACM, 1999.
- [11] S. Parthasarathy, R. Subramonian, Facilitating data mining on a network of workstations, in: Advances in Distributed and Parallel Knowledge Discovery.
- [12] R.W. Moore, C. Baru, R. Marciano, A. Rajasekar, M. Wan, Data-intensive computing, in: I. Foster, C. Kesselman (Eds.), The Grid: Blueprint for a New Computing Infrastructure, Morgan Kaufmann, San Francisco, CA, 1999, pp. 105–129.
- [13] Globus data grid, September 2, 2002. <http://www.globus.org/datagrid/>.
- [14] M.D. Beynon, T. Kurc, U. Catalyurek, C. Chang, A. Sussman, J. Saltz, Distributed processing of very large datasets with datacutter, Parallel Comput. 27 (11) (2001) 1457–1478.
- [15] Patrick, Y. Guo, The design of a platform for distributed kdd components, in: S. Parthasarathy, H. Kargupta, V. Kumar, D.

- Skillicorn, M. Zaki (Eds.), High Performance Data Mining, SIAM, Philadelphia, PA, 2002, pp. 63–78.
- [16] W. Du, G. Agrawal, Using general grid tools and compiler technology for distributed data mining: preliminary report, in: S. Parthasarathy, H. Kargupta, V. Kumar, D. Skillicorn, M. Zaki (Eds.), High Performance Data Mining, SIAM, Philadelphia, PA, 2002, pp. 51–61.
- [17] R. Grossman, M. Hornick, G. Meyer, Data mining standards initiatives, *Commun. ACM* 45 (8) (2002) 59–61.
- [18] W3C Semantic Web, September 2, 2002. <http://www.w3.org/2001/sw/>.
- [19] R. Grossman, M. Mazzucco, Dataspace—a web infrastructure for the exploratory analysis and mining of data, *IEEE Computing in Science and Engineering*.
- [20] M. Veeraraghavan, R. Karri, T. Moors, M. Karol, R. Grobler, Architectures and protocols that enable new applications on optical networks, *IEEE Commun. Mag.* 39 (3) (2001) 118–127.
- [21] A. Chervenak, I. Foster, C. Kesselman, C. Salisbury, S. Tuecke, The data grid: towards an architecture for the distributed management and analysis of large scientific datasets, *J. Network Comput. Appl.* 23 (2001) 187–200.
- [22] R.L. Grossman, M. Mazzucco, H. Sivakumar, Y. Pan, Simple available bandwidth utilization library for high-speed wide area networks, *J. Supercomputing* (2003), to appear.
- [23] D. Lillethun, J. Mambretti, J. Weinberger, Odin: path services for optical networks, in preparation. <http://www.icair.org>.
- [24] S. Bailey, E. Creel, R.L. Grossman, S. Gutti, H. Sivakumar, A high performance implementation of the data space transfer protocol (dstp), in: *Large-scale Parallel Data Mining*, Springer, Berlin, 2000, pp. 55–64.
- [25] M. Mazzucco, A. Ananthanarayan, R.L. Grossman, J. Levera, G.B. Rao, Merging multiple data streams on common keys over high performance networks, in: *Proceedings of the Conference, Supercomputing'2002*, IEEE and ACM, 2002.
- [26] D. Clark, M. Lambert, L. Zhang, NETBLT: a high throughput transport protocol, in: *Proceedings of the Conference, ACM-SIGCOMM'87 on Frontiers in Computer Communications Technology*, 1987, pp. 353–359.
- [27] J. Mambretti, OMNInet. <http://www.icair.org/omninet>.
- [28] Terra wide data mining testbed, September 2, 2002. <http://www.ncdm.uic.edu/testbeds.htm>.
- [29] National Center for Atmospheric Research, Community Climate Model, April 10, 2002. <http://www.cgd.ucar.edu/cms/ccm3/>.
- [30] R.L. Grossman, J. Levera, M. Mazzucco, Aggregate queries on streams of data using a small buffer, UIC Laboratory for Advanced Computing Technical Report, 2002.



Robert L. Grossman is the Director of the Laboratory for Advanced Computing and the National Center for Data Mining at the University of Illinois at Chicago, where he has been a Faculty Member since 1988. He is also the spokesperson for the Data Mining Group (DMG), an industry consortium responsible for the Predictive Model Markup Language (PMML), an XML language for data mining and predictive modeling. He is the President of the Two Cultures Group, which provides consulting and outsourced services focused on data. He has published over 75 papers in refereed journals and proceedings on internet computing, data mining, high performance networking, business intelligence, and related areas, and lectured extensively at conferences and workshops.



Yunhong Gu is a PhD candidate in Computer Science at UIC and a Research Assistant at the Laboratory for Advanced Computing. His research interests are data mining and high-speed networks.



Don Hamelburg is currently a Postdoctoral Research Associate at the Laboratory for Advanced Computing at UIC. He received his BS (Hons) from Fourah Bay College (1993), and his MS and PhD from Georgia State University (1997 and 2001, respectively). His current research interests include structural bioinformatics, data mining and computational biology.



Dave Hanley has a bachelor's degree from UIC and is pursuing his master's degree in Computer Science. He is the co-author of two computer books, *C: Just the FAQ's* and *Visual J++ Unleashed*. He is co-author of numerous papers and has won a series of computer programming contests.



Xinwei Hong is a Postdoctoral Research Associate at the Laboratory for Advanced Computing at UIC. He received his PhD in Electronics and Information Engineering from Huazhong University of Science and Technology in 1998. His current research interests include developing high-speed data delivery over high-performance wide-area networks.



Marco Mazzucco is a Postdoctoral Fellow at the University of Wales, Swansea. He is currently working on an ESPRERC funded project in theoretic computer science under the direction of Dr. Martin Otto. He also does research and consulting for the National Center for Data Mining. He received his PhD in Mathematics at UIC in 2000.



Jorge Levera received his Bachelor's degree in Computer Science from the Universidad Catolica Nuestra Senora de la Asuncion. In 1997 he joined the Laboratorio de Electronica Digital. In 1999 he won the Fulbright/LASPAU grant that sponsored his Master in Computer Science Program at the University of Illinois at Chicago (UIC). He is a PhD candidate in Computer Science at UIC

and a Research Assistant at the Laboratory for Advanced Computing. His research interests are data mining and high-speed networks.



Joe Mambretti is the Director of the International Center for Advanced Internet Research (iCAIR) at Northwestern University, the Director of the Metropolitan Research and Education Network (MREN), member of the StarLight/STAR TAP partnership, a member of the I-WIRE consortium, and a principal researcher on the OMNInet project. The mission of iCAIR is to accelerate

leading-edge innovation and enhanced digital global communications through advanced Internet technologies, in partnership with the international community. iCAIR accomplishes that mission by undertaking large-scale (e.g., global, national, region, metro) projects in several key areas, high performance resource intensive applications, advanced middleware and metasystems, and large scale optical network infrastructure.



Dave Lillethun is a Research Associate at iCAIR and a principal architect and software engineer and protocol developer for ODIN, THOR, and DEITI, and investigator of methods for intelligent application signaling. Dave designed and developed a number of successful experiments related to innovative intelligent application signalling techniques on the OMNInet testbed. Dave also has an interest in developing new techniques for medical imaging. Previously, he was a designer and developer for a major software corporation.

and a Research Assistant at the Laboratory for Advanced Computing. His research interests are data mining and high-speed networks.



Jeremy Weinberger received his BA degree in computing and information systems from Northwestern University and is currently a research associate at iCAIR. He is manager of Operations for OMNInet and a principal architect and developer for ODIN, THOR and DEITI and related protocols, and investigator of new intelligent signaling methods. He has also been involved in experimentation and

development of IETF DiffServ QoS methods, and developed experiments for a regional Science Grid. He also developed several successful international DiffServ experiments with CERN and various Asia Pacific research centers.