

# A DataSpace Infrastructure for Astronomical Data

Robert Grossman

Magnify, Inc. and the University of Illinois at Chicago

Emory Creel and Marco Mazzucco

University of Illinois at Chicago

Roy Williams

California Institute of Technology

May 4, 2001

**This is a draft of a paper which later appeared R. L. Grossman, C. Kamath, W. Philip Kegelmeye, V. Kumar, and R. Namburu, *Data Mining for Scientific and Engineering Applications*, Kluwer Academic Publishers, 2001, pages 115-123.**

## Abstract

This article describes an internet infrastructure for working with data called DataSpace. A distributed DataSpace application containing data from the 2MASS and DPOSS astronomical data sets is also described. DataSpace is designed so that client applications supporting the remote analysis and distributed mining of data are easy to build.

## 1 Introduction

The web today provides an infrastructure which is better suited for working with multimedia documents than for working with data. In this article, we describe an internet infrastructure called DataSpace, which is designed to support the publishing, analysis, and mining of remote and distributed data. We also show how DataSpace can be used to work with astronomical data, and, in particular, to create virtual observatories.

The core technology is a network based transfer protocol called the data space transfer protocol (DSTP) [9], a platform independent way to share

<b>Description</b>	<b>Document Web</b>	<b>Data Web</b>
retrieve	retrieve a remote document	analyze & mine remote data
search & mine	search engines using keyword search	data mining engines analyzing distributed data
language	HTML, XML	XML, PMML [19], SOAP
protocol	HTTP	DSTP
server	web server	data server
hardware	workstation	workstation clusters
network	10K document can be retrieved < 1 sec with commodity internet	1 MB data set retrieved < 1 sec with OC-3 network

Table 1: Next generation networks with high bandwidth, quality of service and security will enable the emergence of a data web, broadly analogous to today’s document web.

data over a network. It is based on a data model which views data as a collection of distributed columns which can be linked for analysis through “global keys” called universal correlation keys (UCK’s). In this way, DataSpace provides a means for remote data analysis and distributed data mining which is independent of the type of data storage used, whether it be files, databases, or distributed data warehouses.

The DSTP protocol enables an application in one location to locate, access, and analyze data from several other locations. It also reduces the dependency on the data file, because it correlates data based on UCK’s in different data sets. The DSTP protocol provides the infrastructure for a global *data space* in the same way as the http provides the infrastructure for a global *document space*.

Section 2 describes related work and provides some background material. Section 3 contains a quick introduction to DataSpace. Section 4 describe the virtual observatory we created with DataSpace. Section 5 shows the DSTP protocol underlying a typical query to the virtual observation. Section 6 contains the conclusion and summary.

## 2 Related Work

Today there are several common formats for scientific and engineering data available on the web:

*Flat Files.* The most common format is probably files, either simple ASCII files with data arranged in rows and columns, self describing file formats such as HDF [13], or binary formats, accompanied by programs for extracting and working with the data. Often the relevant data file is selected using a cgi program which allows the user to query on the metadata.

*HTML Tables.* If the data set is small enough, HTML tables can be used so that the data is directly available for viewing.

*Databases.* Another common format is to put the data into a database and to provide access to the database using a form to produce an sql query, which is passed to the database via an API (application program interface) such as ODBC (open database connectivity) [16] or JDBC (java database connectivity) [14]. If the number of records returned is small, the records can be directly displayed in HTML; if not, the records can be output to a file.

*XML.* Some data is beginning to be stored in XML (extensible markup language), especially if the data sets are small.

Access to data is usually through a cgi program. If the data is in a database, access may involve ODBC or JDBC. Sometimes agents may be used. Access may also involve SOAP (simple object access protocol) [17] and HTTP.

Unless specialized programs are available for accessing file based data, file based data usually requires that the entire data be transferred before analysis can begin. Although databases are efficient for writing, selecting and reading records, they are not designed to provide a mechanism for the efficient statistical analysis of data. Data in HTML tables can easily be displayed, but not otherwise manipulated. Finally, XML provides a very flexible import and export format, but does not directly support the analysis of data. Also, it is very verbose which is not efficient for large data sets. DataSpace is designed to provide an alternative mechanism for storing, accessing, and analyzing data, which is complementary to the formats and access methods just described.

### 3 Data Space

In this section, we give a quick introduction to some of the key concepts underlying DataSpace.

*Data Sites.* A data site is simply a web site with a DSTP server listening on port 5040 for requests from DSTP clients.

*Catalogs and Metadata.* Adding data to a site is easy.

1. First, create a text file, say `dposs.dat` containing rows and columns of data.
2. Second, add the name of the text file to one of the catalog files at the site, say `catalog.dstp`, which contains the names of all the data files.
3. Third, add a file `dposs.ds` containing simple metadata about the data file, such as the names of the attribute, their data types, etc.

*Columns of Data.* Although data is physically loaded into a data site using files, DataSpace logically organizes the data into columns and provides efficient physical access to one or more columns of data.

*Universal Correlation Keys.* Each column of data is associated with one or more universal correlation keys or (UCK's). Each UCK is associated with a globally unique identifier. Two distributed columns of data with the same UCK may be meaningfully compared. An example of a UCK the number of seconds since January 1, 1970. Another example is right ascension and declination measured in degrees.

### 4 A DataSpace Based Virtual Observatory

Using the dataspace infrastructure, we developed an application which simultaneously works with two geographically distributed astronomical source catalogs: 2MASS (Two Micron All Sky Survey) survey [1] data from Cal-Tech and DPOSS (Digital Palomar Observatory Sky Survey) survey [2] data. The 2MASS data are in the optical wavelengths (0.4 - 0.7 micron), while the DPOSS data are in the infrared (1.2 - 2.2 micron) range.

Our goal was to create a virtual observatory supporting the statistical analysis of many millions of stars and galaxies with data coming from both surveys. In order to support this type of analysis, we must identify the same star in both surveys. It is a simple matter: if two star images are in the same position, they must be the same star. However, each star has a finite angular size in the sky, so there is a slight “fuzziness” in position, and this is exacerbated for galaxies. A typical query is of the form: “Find all pairs, one from the DPOSS catalog and one from 2MASS, whose angular separation is less than a given tolerance”.

In the DSTP application we built, this query is visualized by coloring the data as follows:

- red (appears only in 2MASS), or
- blue(appears only in DPOSS), or
- magenta (successful join: appears in both surveys).

At CalTech, we populated one DSTP server with 2MASS data. For the studies described, we populated another DSTP server with DPOSS data. This was easier than putting a second DSTP server at the DPOSS home site. The DPOSS DSTP server contains a list of stars within a given region of the sky, with magnitudes in three different colors. The 2MASS DSTP servers contains a list of stars in a given region, accompanied by extensive data about the stars.

The client DSTP application formulates and sends requests to a compute server, which interprets the client’s request, gets the relevant data from the DPOSS and 2MASS servers, performs the fuzzy join, and sends the resulting stars and galaxies back to the client.

An architecture like this, which uses separate servers for data and computational tasks, can scale to large numbers of distributed services which may be communicating very large amounts of data at very high bandwidth.

This virtual observatory we developed was introduced at Supercomputing 2000 in Dallas in early November, 2000. Figure 1.1 is a sample output of our client. The position of the light sources in both data sets were given in polar coordinates, right ascension and declination. For this experiment a tolerance of 2 seconds was chosen. This tolerance was suggested to us physicists at California Institute of Technology familiar with the two data sets. The data was raw data, with out any attempt to eliminate noise. Also, in order to achieve a reasonable graph, we restricted the data in both catalogs to the portion of the sky in the right ascension range of 183 to 270 degrees

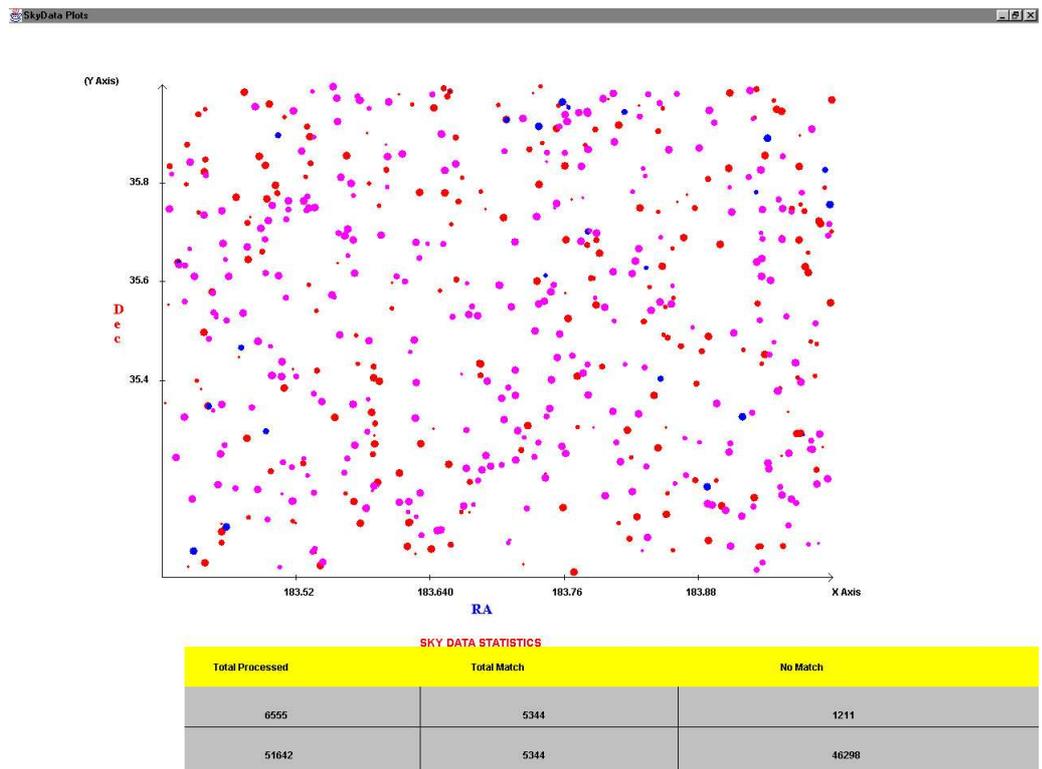


Figure 1: An example of a DSTP client.

and declination in the range of 17 to 47 degrees (somewhere over the north pole).

As can be seen, a match rate of about 82 percent was observed between these two data sets. Something that can not be observed from a screen shot is the continuous change of both the graph and the statistics below the graph during the runtime of this application. This is one of the many benefits of the DSTP protocol over a simple file transfer, its ability to stream data and allow a client to process the data in real time.

## 5 A Sample DSTP Session

The Data Space Transfer Protocol (DSTP) is a protocol for retrieving data from remote nodes in the internet. It is broadly based upon the Network News Transfer Protocol (NNTP) [15], a protocol for retrieval of news articles. The DSTP server uses a stream connection and NNTP-like commands and responses. It is designed to accept connections from hosts and to provide a simple interface to the data columns on the server. A DSTP server functions as an interface between DSTP applications and remote data.

Rather than describe the protocol and commands formally, this section contains an annotated session between a DSTP Server named S and a DSTP client named C (not their real names). First we see the client contacting the server on tcp port 5040.

```
C: Trying 131.193.181.121...
C: Connected to ncdm121.lac.uic.edu.
C: Escape character is '^]'.
S: 200 ncdm121.lac.uic.edu DSTP server v0.90 ready
S: .
```

The client can now request which universal correlation keys (UCK's) the server is prepared to serve. The right ascension is assigned the UCK 10 and the declination is assigned the UCK 11, which are the polar coordinates for the sky data.

```
C: list uck
S: 214 List of UCKs follows
S: <UCK ID="10" NAME="ra" NUMATTR="42"/>
S: <UCK ID="11" NAME="dec" NUMATTR="42"/>
S: <UCK ID="91" NAME="designation" NUMATTR="42"/>
S: <UCK ID="20" NAME="j_m" NUMATTR="42"/>
S: <UCK ID="21" NAME="h_m" NUMATTR="42"/>
S: <UCK ID="22" NAME="k_m" NUMATTR="42"/>
S: <UCK ID="23" NAME="id_opt" NUMATTR="42"/>
S: <UCK ID="30" NAME="b_m_opt" NUMATTR="42"/>
S: <etc> ...
S:.
```

We assume now, for this example, the client has access to another DSTP server which also can serve the UCK's 10 and 11. So the client requests that the server set these to indexes for correlation.

```
C: set uck 10
S: 211 UCK 10 ra selected
S: .
C: set uck 11
S: 211 UCK 11 dec selected
S: .
```

Next the client can ask what data files the server has which contain these UCK's:

```

C: list uck datafile
S: 234 List of Datafiles follows
S: <DATAFILE NAME="2mass.dat" NUMRECORDS="2912952" DSFILENAME="2mass.ds"
    UCKNAME="ra" UCKID="10" NUMATTR="42" UCKNAME="dec" UCKID="11"
    NUMATTR="42" UCKNAME="designation" UCKID="91" NUMATTR="42"
    UCKNAME="j_m" UCKID="20" NUMATTR="42" UCKNAME="h_m" UCKID="21"
    NUMATTR="42" UCKNAME="k_m" UCKID="22" NUMATTR="42" UCKNAME="id_opt"
    UCKID="23" NUMATTR="42" UCKNAME="b_m_opt" UCKID="30" NUMATTR="42"
    UCKNAME="r_m_opt" UCKID="31" NUMATTR="42" UCKNAME="glon"
    UCKID="13" NUMATTR="42" UCKNAME="glat" UCKID="14" NUMATTR="42"
    UCKNAME="j_h" UCKID="40" NUMATTR="42" UCKNAME="h_k" UCKID="41"
    NUMATTR="42" UCKNAME="j_k" UCKID="42" NUMATTR="42"/>
S: <DATAFILE NAME="2mass_s.dat" NUMRECORDS="536" DSFILENAME="2mass_s.ds"
    UCKNAME="ra" UCKID="10" NUMATTR="5" UCKNAME="dec" UCKID="11"
    NUMATTR="5" UCKNAME="j_m" UCKID="20" NUMATTR="5" UCKNAME="h_m"
    UCKID="21" NUMATTR="5" UCKNAME="k_m" UCKID="22" NUMATTR="5"/>
S: .

```

Then the client selects the desired data file:

```

C: set datafile 2mass.dat
S: 230 Datafile 2mass.dat selected
S: .

```

After the data file is selected, the client can request the corresponding meta-  
data:

```

C: metadata
S: 250 Metadata follows
S: <DATAFILE NAME="2mass.dat" NUMRECORDS="2912952" DSFILENAME="2mass.ds"/>
S: <UCKNAME="ra" UCKID="10" NUMATTR="42" />
S: <UCKNAME="dec" UCKID="11" NUMATTR="42" />
S: <UCKNAME="designation" UCKID="91" NUMATTR="42" />
S: <etc>
S: <ATTRIBUTE-DESCRIPTOR NUMBER="1" NAME="ra" DATA-TYPE="real"

```

```

    USE-AS="continuous" UNIT="deg" />
S: <ATTRIBUTE-DESCRIPTOR NUMBER="2" NAME="dec" DATA-TYPE="real"
    USE-AS="continuous" UNIT="deg" />
S: <ATTRIBUTE-DESCRIPTOR NUMBER="3" NAME="err_maj" DATA-TYPE="real"
    USE-AS="continuous" UNIT="arcsec" />
S: <ATTRIBUTE-DESCRIPTOR NUMBER="4" NAME="err_min" DATA-TYPE="real"
    USE-AS="continuous" UNIT="arcsec" />
S: <ATTRIBUTE-DESCRIPTOR NUMBER="5" NAME="err_ang" DATA-TYPE="integer"
    USE-AS="continuous" UNIT="deg" />
S: <etc>
.

```

Finally the client requests the data itself. First the client sets the range for the rows of interest, 10 through 20, in this example. Then the client issues the data command along with a list of attributes which are desired.

```

C: set line 10 20
S: 270 Lines from 10 to 20 selected
S: .
C: data 1 10 34
S: 242 [1, 10, 34] retrieved - data follows
S: 183.400288,27.646980,183.400288,14.797,0.90
S: 183.400334,39.956638,183.400334,10.999,1.84
S: 183.400358,41.809029,183.400358,12.889,0.86
S: <etc>
S: .

```

Finally the client closes the connection with the server:

```

C: quit
S: 205 GoodBye
S: .
C: Connection closed by foreign host.

```

## 6 Discussion and Summary

In this article, we have described how DataSpace can be used to provide a simple means to analyze remote astronomical data, such as the DPOSS and 2MASS surveys. We have also illustrated how DataSpace provides a simple way to support distributed data mining queries, such as the fuzzy match identifying objects appearing in both catalogs. With DataSpace, virtual astronomical observatories can easily be built and be used to provide access to a wide range of astronomical data. Moreover, with the universal correlation keys used by DataSpace, survey data can also be used by other applications.

## 7 Acknowledgements

The DSTP based virtual observatory was created by a project team co-led by Robert Grosman and Roy Williams. The team consisted of Shirley Connelly, Emory Creel, Sivakumar Harinath, Marco Mazzucco, Andriy Turinskiy, and Scott Wahlstrom.

## References

- [1] See the web site: <http://www.ipac.caltech.edu/2mass/>
- [2] See the web site: <http://astro.caltech.edu/rrg/science/dposs.html>
- [3] I. Foster and C. Kesselman, Globus: A Metacomputing Infrastructure Toolkit, International Journal Supercomputing Applications, Volume 11, pages 115-128, 1997.
- [4] I. Foster and C. Kesselman, editors, The Grid: Blueprint for a New Computing Infrastructure, Morgan Kaufmann, San Francisco, 1998.
- [5] A. Grimshaw and W. Wulf, The Legion Vision of a Worldwide Virtual Computer, Communications of the ACM, Volume 40, pages 39-45, 1997.
- [6] A. Grimshaw, A. Ferrari, G. Lindahl, and K. Holcomb, Metasystems, Communications of the ACM, Volume 41, pages 46-55, 1998.
- [7] R. L. Grossman, S. Bailey, A. Ramu, B. Malhi and A. Turinsky, The Preliminary Design of Papyrus: A System for High Performance, Distributed Data Mining over Clusters, in Advances in Distributed and

Parallel Knowledge Discovery, H. Kargupta and P. Chan, editors, AAAI Press/The MIT Press, Menlo Park, California, 2000, pages 259-275.

- [8] R. L. Grossman, S. Bailey, A. Ramu, B. Malhi and H. Sivakumar, A. Turinsky, Papyrus: A System for Data Mining over Local and Wide Area Clusters and Super-Clusters, Proceedings of Supercomputing 1999, IEEE.
- [9] S. Bailey, E. Creel, R. Grossman, S. Gutti, and H. Sivakumar, A High Performance Implementation of the Data Space Transfer Protocol (DSTP), Large-Scale Parallel Data Mining, M. J. Zaki and C.-T. Ho, editors, Springer-Verlag, Berlin, 2000, pages 55-64.
- [10] R. L. Grossman, An Open, Web-Based Infrastructure for Working with Data, submitted for publication.
- [11] H. Sivakumar, S. Bailey, R. L. Grossman, Pockets: The Case for Application-level Network Striping for Data Intensive Applications using High Speed Wide Area Networks, Proceedings of Supercomputing 2000, IEEE.
- [12] E. Creel, R. L. Grossman, M. Mazzuco, Data Space: Protocols and Services for Distributed Data Mining and Remote Data Analysis, submitted for publication.
- [13] The Hierarchical Data Format web site is <http://hdf.ncsa.uiuc.edu>.
- [14] George Reese, Database Programming with JDBC and Java, 2nd Edition, O'Reilly, 2000.
- [15] Phil Lapsley, Network News Transfer Protocol, February 1986.
- [16] Roger Sanders, Hands on ODBC 3.5 Developer's Guide, Osborne McGraw-Hill, 1998.
- [17] Retrieved from <http://www.w3.org/SOAP/> on March 3, 2001.
- [18] R. Wolski, N. Spring, J. Hayes, The Network Weather Service: A Distributed Resource Performance Forecasting Service for Metacomputing, UCSD Technical Report TR-CS98-599, September, 1998
- [19] The Predictive Model Markup Language (PMML) web site is <http://www.dmg.org>.