

## 9A.5

### Deployable Suite of Data Mining Web Services for Online Science Data Repositories

Sara Graves, Rahul Ramachandran, Ken Keiser, and Manil Maskey  
*University of Alabama in Huntsville*

Christopher Lynnes and Long Pham  
*Goddard Earth Sciences Data and Information Services Center*

#### 1 Abstract

A project is currently underway to create a suite of specialized deployable data mining web services designed specifically for science data. The project leverages the Algorithm Development and Mining (ADaM) toolkit as the basis. The ADaM toolkit is a robust, mature and freely available science data mining toolkit that is being used by different research organizations and educational institutions worldwide. These deployable services will give the scientific community a powerful and versatile data mining capability that can be used to create higher order products such as thematic maps from current and future NASA satellite data records with methods that are not currently available. Many of the specialized data mining, pattern recognition, image processing and data preparation algorithms in ADaM are specifically geared towards satellite imagery, making these tools a perfect fit for NASA satellite data. In addition to providing specialized deployable data mining services, the suite will use the Earth Science Markup Language (ESML) to handle a variety of heterogeneous data formats seamlessly. ESML is another proven technology that, like ADaM, is in use by organizations worldwide. The deployable package of mining and related services are being developed using web services standards so that community based measurement processing systems can access and interoperate with them. The maturation of web services standards and technology sets the stage for a distributed "Service-Oriented Architecture" (SOA) for

---

Sara Graves  
Information Technology and Systems Center,  
University of Alabama in Huntsville,  
Huntsville, AL 35899  
Tel: 256-824-6066  
email: sgraves@itsc.uah.edu

NASA's next generation science data processing. This architecture will allow members of the scientific community to create and combine persistent distributed data processing services and make them available to other users over the internet. The Goddard Earth Sciences Data and Information Services Center (GES DISC) will serve as an operational site for demonstration.

#### 2 Introduction

The Mining Web Services (MWS) technology is funded by NASA to develop ways of orchestrating the use of web services to process scientific data at online data repositories. The Information Technology and Systems Center (ITSC) at the University of Alabama in Huntsville has teamed with the Goddard Earth Sciences Data and Information Services Center (GES DISC) for this purpose. This collaboration is practical since ITSC has long had data mining capabilities that are useful in the analysis of the scientific data, and the GES DISC has a large online repository of many NASA scientific data sets, together with an "engine" for applying user-supplied processing algorithms to the data. Previous approaches to mining this data with the ADaM toolkit would typically involve a researcher installing the ADaM software locally, requesting and transporting the data of interest to the same location and then programming scripts to execute mining workflows to process the data. While effective, this approach is not scalable when the data volumes become a significant obstacle to timely data transport operations. Also, any solution (script) developed typically remained an isolated capability at the researcher's location and not reusable by external collaborators or other researchers interested in similar investigations. The MWS team is implementing an architecture based on web services that will allow researchers to develop and test analysis workflows locally and then ultimately deploy the workflow against

large amounts of data at a remote data center. This removes the need to transport and provide redundant storage for large amounts of data and can potentially allow for reuse of analyses if a researcher so chooses.

### **3 Background**

The MWS project was structured to take advantage of several existing resources that provided capabilities that are important to the effort. These included web services technology that provides a standard way to remotely interface with programmatic components and to orchestrate the chaining of services through standardized service descriptions. The ADaM Toolkit has been developed and expanded over several years and provides a rich set of analysis tools. ESML provides the ability to programmatically interpret heterogeneous data formats, thus support data interoperability. The GES DISC has a large online repository of publicly available scientific data sets as well as existing internal process automation capabilities.

#### **3.1 ADaM (Algorithm Development and Mining) Toolkit**

The ADaM system (*Hinke, 1997a; Hinke, 1997b; Hinke, 2000; Ramachandran, 2000*) was originally developed in 1994 with the goal of mining large scientific data sets for geophysical phenomena detection and feature extraction, and has continued to be expanded and improved. ADaM provides knowledge discovery and data mining capabilities for data values, as well as for metadata. Thus, unlike most data mining software, ADaM has been designed for use with scientific and image data from the outset. ADaM includes not only traditional data mining capabilities such as pattern recognition, but also image processing and optimization capabilities, and many supporting *data preparation* algorithms that are useful in the mining process. Recently, ADaM was redesigned as a toolkit of discrete, independent components for the evolving computing landscape which can be used together in different combinations to perform many complex tasks. This redesign also allows the algorithms in ADaM to be easily packaged as grid or web services (*Rushing, 2005*) and is being extensively used by different research groups.

#### **3.2 ESML (Earth Science Markup Language)**

Data providers typically offer their data in different data formats. Dealing with heterogeneous formats can potentially become a major obstacle in the utility of distributed web services. However, there already exists a technology that can be leveraged to overcome this problem. The Earth Science Markup Language (ESML) is a specialized markup language for Earth Science metadata based on XML which provides a solution to this problem (*Ramachandran, 2004a; Ramachandran, 2004b*). ESML is being used by a number of research groups applications to seamlessly utilize datasets in heterogeneous formats. There are three components of ESML: the ESML Schema, the ESML Description Files and the ESML Library. The ESML Schema defines the grammar for writing ESML Description Files which contain external metadata with content and structural information for the corresponding data file format. These machine-readable and interpretable markups allow applications to parse ESML Descriptions and read the data files, eliminating the need for additional data conversion software. ESML Description Files can be generated at any time to allow data/application interoperability. Because ESML Description Files are external files, they do not modify the application or the data file itself. The ESML Library is used by applications to parse the relevant ESML Description File for the structural information and read the data. For ADaM's two internal data models, there are plans to write a Data Translation service using the ESML Library which will be able to read the ESML Description Files and data in different formats and translate them into one of the ADaM internal data models.

#### **3.3 GES DISC (Goddard Earth Sciences Data and Information Services Center)**

Since 2001, the GES DISC has been providing users with the ability to mine data at the archive where the data are stored (*Lynnes, 2001*). This avoids the need to deliver large volumes of data to the user for large-scale mining activities. The Near Archive Data Mining (NADM) system was first implemented for the Tropical Rainfall Measuring Mission. It allows users to run data mining or other data reduction algorithms of their own design either on newly acquired data, or by pulling data out of the archive in a data

mining campaign. The user-provided algorithms are uploaded and built through a set of web forms. This system was later implemented for data from the Terra and Aqua satellites. This data mining capability has been successful at enabling several data mining and reduction activities that would not otherwise be possible. However, the requirement for the user to supply the algorithm has limited its use to a select few “high-end” users. It has long been a goal of the GES DISC to make data mining algorithms available to a broader community. In fact, *Lynnes and Mack, 2001* proposed a scenario by which third-party data service providers (such as UAH/ITSC) would provide algorithms to run on the user’s behalf at the data archive. With the proposed mining web services it will finally be possible to do this in a standardized, cost-effective way.

Since GES DISC’s NADM is an operational data center, it is not the ideal site for the highly iterative process of developing a data mining “recipe”. Instead, ITSC will develop an external sandbox that allows the user to experiment with a small subset of the data and a variety of ADaM algorithms and of parameter settings. Once the user is satisfied with a recipe of modules and settings, a simple press of a button will cause this recipe to be invoked on large quantities of data within GES DISC’s Near Archive Data Mining system. The user will be able to mine interactively, or to request a mining subscription, which will run on newly produced or acquired data when it is ingested at the GES DISC. The interfaces and security mechanisms between the “sandbox” and data provider defined and implemented for this deployment can be used in other similar situations, requiring the integration of ADaM services into an existing processing system.

## **4 Approach**

For some time, the University of Alabama in Huntsville (UAH) Information Technology and Systems Center (ITSC) has been investigating the use of distributed services for use with data mining, subsetting, image processing, thematic map generation and other spatially oriented data applications. The mining web services technology (see Figure 1) builds on this expertise to provide a deployable web service suite, hardened for widespread use.

## **4.1 Web Services**

A web service is a software component designed to be used between distributed machines following standard internet protocols. Web services are the basic components of an SOA.

The two main styles of web services are Simple Object Access Protocol (SOAP) and Representational State Transfer (REST). SOAP provides a standard message protocol for communication based on XML. SOAP web services have two main conventions: any non-binary attachment messages must be carried by SOAP and the service must be described using Web Service Description Language (WSDL). SOAP was selected primarily for this WSDL feature, which has particularly useful for driving the composition of complex workflows. SOAP has become so popular that the terms “SOAP” and “web service” are often used interchangeably.

### **4.1.1 Mining Web Services**

ADaM components are implemented using standard C++ for portability across platforms. These discrete components can already be scripted together in a variety of ways to solve complex problems. This project will augment this existing technology by repackaging ADaM components as SOAP-based web services for easier and more dynamic integration in emerging distributed service oriented architectures. The suite of mining and related services will be packaged into deployable bundles to allow easy deployment in the web server environments at online data provider sites. An incremental development approach will be utilized, with different ADaM algorithms selected for three major staggered releases based on their utility and the requirements of the target demonstration deployments. The “deploy early and deploy often” philosophy will be used in releasing these packages to prospective users with an additional set of algorithms added at every new release. Also, a registration mechanism for advertising the different service capabilities of a package via ECHO (Earth Observing System Clearinghouse) (*ECHO, 2006*) will be created.

### **4.1.2 Translation Services**

The first release package will contain the ESML translation service to allow the use of the mining suite on several different science data formats. UAH/ITSC has already built similar grid

services for other projects such as the NSF large Information Technology Research effort, Linked Environments for Atmospheric Discovery (LEAD). The second package will contain translator services to support visualization of the results of ADaM image processing services. These translators will support integration of the image processing services with Open Geospatial Consortium (OGC) compatible web services, such as web mapping and web coverage services (WMS, WCS). WMS output is an image of the requested data, where WCS output is the actual data. This project will create specialized services that provide WMS/WCS interfaces to the final results of the analytical or image processing. A composition of ADaM image processing and WMS/WCS services will provide for visualization of results and merging with other data in a WMS/WCS client. This approach will provide a standardized and well-accepted approach for applications and users to receive the results of the image processing services. NASA supports the OGC standardization efforts and is a strategic OGC member. UAH/ITSC is a university associate member of OGC and actively supports the advancement of the organization's efforts.

#### **4.2 SOA**

A Service Oriented Architecture (SOA) packages together independent services that have well defined interfaces. These interfaces are designed such that they are interchangeably invoked by applications or other services implemented on heterogeneous platforms and languages. Some benefits or developing applications in within an SOA include re-use of existing software, reduced deployment time, minimal code changes reduces chance of introducing errors, and orchestration of data and computational resources at distributed locations.

#### **4.3 Service Workflows (chaining)**

A variety of solutions for orchestrating service workflows are being explored, but MWS is initially concentrating on using BPEL (Business Process Execution Language). The use of BPEL was chosen primarily because it has been well received by many communities and applications so it approaches the level of a standard, and as a result of this acceptance there are a number of evolving tools available for this technology. Acceptance of a common execution language allows for the portability of workflow definitions across more systems and hopefully easier

acceptance at other data centers as the benefits are demonstrated.

As part of this project, the suite of mining and related web services will be deployed at existing NASA data provider sites which will demonstrate the utility of this package in a variety of environments. The utility of this system to users will be the ability to define and execute service workflows, comprised in part of the services deployed at NASA data centers. The proximity of the services, such as data mining, image process, etc., at the data centers will allow access to large amounts of NASA data for user-defined analyses. These analyses workflows themselves can then be exposed as services for other researchers to utilize, assuming permission by the originators, resulting in potential for a rich set of functionality that researchers can deploy against large amounts of data without the overhead of local storage or processing resources.

#### **4.4 Implementation**

The initial implementation of web services for existing ADaM tools has been completed by ITSC. The implementation included placing SOAP wrapper around each ADaM operation. ITSC considered many options of the SOAP implementation before deciding on SOAP::Lite (SOAP, 2006).

SOAP::Lite is the Perl implementation of SOAP. Although Java's implementation of SOAP, Axis (Apache, 2006), provides wide variety of tools that assist in implementing services and clients, Java was not appropriate in this case because the ADaM tool was implemented in C++. ITSC had encountered numerous issues while interfacing Java Native Interface (JNI) to C++ for other ITSC projects. The Axis C++ offers a C++ implementation of SOAP, however, it has not been well-accepted by the software industry so far. The Python implementation of SOAP, pySOAP, has not been up-to-date. Therefore, due to ITSC's experience on other projects and the familiarity of Perl at Goddard Space Flight Center (GSFC), ITSC decided on SOAP::Lite for Perl as the web service implementation language.

Once UAH finished implementing the web services, the Web Service Definition Language (WSDL) for each of the services were published. Then each of the web services was tested for interoperability using Perl and Java clients.

In order to allow the data center full control over enforcing local policies for file creation, the SOAP wrappers specify the full path of the input file only, plus additional parameters as necessary. The data center constructs the output directory and filename and returns it the calling program in the SOAP response. This may then be used as the input file to the next step in the workflow. This interface mechanism allows the data center to control where files are created on its system and to implement safeguards against inappropriate usage.

#### 4.4.1 Web Service Composition

Science problems are usually too complicated for a single operation to undertake. A set of operations is often necessary to provide an appropriate solution to a complex problem. When using SOA to solve science problems, solution to such complex problems can be designed using a workflow that employs various web services. A workflow describes how tasks are orchestrated, what component performs them, what their relative order is, how they are synchronized, how information flows to support the tasks and how tasks are being tracked.

Currently, the industry standard for service orchestration is the Business Process Execution Language (BPEL). BPEL provides a standard XML schema for workflow composition of web services that are based on SOAP. There are other workflow composition tools that create workflow descriptions for a set of web services execution, however, the tools are not standardized yet. This standardized composition description is eventually deployed on a BPEL engines. A BPEL engine is required to process the instructions described in BPEL. Most of the BPEL engines are open source software, the such as the ActiveBPEL engine (*ActiveBPEL, 2006*). Various visual workflow composers are also freely available. The problem with having different BPEL engines is that even though the BPEL is a standard language, these engines require an additional deployment descriptor that describes engine specific additional dependencies. Thus, the deployment descriptor varies from engine to engine. Here, deploying a BPEL translates to taking a BPEL description and asking the engine to expose the workflow as a web service. Being able to expose a composition of a set of web services makes the solution itself a web service and thus reusable. It

is important to note that the BPEL engine can be totally isolated from the services being utilized.

Ideally, for science problems involving large data, data centers should be capable of hosting web services so that the data can be accessed easily, keeping the network bandwidth to minimum. The only required component at a data center to host web services functionality is an appropriate web server. Dropping the set of web services to the web server container should make the web services available for use.

ITSC has investigated various freely available workflow composers and their associated engines. The two BPEL engines that stand out the most are: ActiveBPEL and Sun's BPEL engine. ActiveBPEL Designer and NetBeans IDE (v5.5) provide the composer for the corresponding BPEL engines respectively. These composers also provide mechanism to create deployment descriptors and user friendly deployment interfaces. ITSC also investigated the BPEL Execution Engine (BEXEE). However, BEXEE is no longer being supported. FiveSight PXE BPEL Engine (*ActiveBPEL, 2006*) has also been investigated.

During the investigation, it was discovered that NetBeans IDE was the best in terms of the usability in creation of workflows. However, ActiveBPEL, provided an easier mechanism to deploy the workflows. For BEXEE and PXE, there were no specified workflow composers. However, deployment descriptors were manually created for the BEXEE and PXE engines. Manual creation of deployment descriptors was error prone. We noticed that the BPEL were transportable between all of these engines. As the result of our investigation, we chose ActiveBPEL as our BPEL engine. Furthermore, a plan was made to develop a mining-specific application for an easier and more user friendly interface to the creation of BPEL and its deployment descriptor. This application will alleviate the necessity of BPEL knowledge to the users.

#### 4.4.2 Demonstration Architecture

The architecture is shown in Figure 2. A workflow is generated within the Workflow Composer, based on experimentation in the Sand Box. This workflow is then deployed to a BPEL

engine (Flow #1), which returns a URL pointing to the WSDL for that workflow (Flow #2). This URL is then transmitted to the GES DISC via a Web Services request along with a specification of the data to be mined, such as the dataset to be mined and temporal or spatial constraints(3). This request is provided to the GES DISC's processing engine (4), also known as the Simple, Scalable Script-based Science Processor for Measurements (S4PM) (*Lynnes, 2006*) for data mining at the data archive. The S4PM engine is responsible for acquiring the data from the archive, then executing the requested workflow on each input file. However, note that the workflow itself is not provided to S4PM. Instead, S4PM uses the supplied WSDL URL to fetch the WSDL document and then invoke the corresponding Web Service in the BPEL engine (5), supplying the full path of the input file. The BPEL engine turns around to invoke the atomic Web Services in the proper order at the GES DISC(6). Finally, the output is transmitted to the end user from the data center (7).

#### 4.4.3 Demonstration

ITSC and GES DISC have worked together to produce a demonstration web service workflow involving the orchestration of mining web services at a client site, using services and data available at the repository.

The demonstration included a solution to a real science problem of creating a cloud mask in satellite images from a Geostationary Operational Environmental Satellite. The solution involves the use of k-means clustering algorithm to label each data point as either part of a cloud or not part of a cloud. There were preprocessing services to convert the generalized satellite format (GSF) cloud data to Attribute-Relation File Format (ARFF) format, which is the allowed data format for k-means operations. Furthermore, the result was post-processed using another web service to convert into an image file for visualization.

ActiveBPEL Designer tool was used as the workflow composer. The ActiveBPEL engine resided in the workflow designer's computer. This demonstration began with the construction of the workflow in the Sandbox (Fig. 2) using a three-step process:

1. Importing the WSDL for required web services using ActiveBPEL Designer.

2. Creating the workflow with the imported services using ActiveBPEL Designer.

3. Deploying the workflow to the BPEL Engine using ActiveBPEL Designer.

The result of this workflow, a URL to the WSDL, was transmitted along with the dataset, start time, and stop time via a SOAP request to the Data Mining Service at the GES DISC. The Data Mining Service deposited the request in the S4PM engine, which proceeded to acquire the data files from a local archive. As each file was acquired, it was processed according to the workflow by invoking the corresponding Web Service presented by the BPEL engine. The results were staged to an FTP directory, from which the user ITSC retrieved them and verified the output.

#### 4.4.4 Proposed Architecture

The architecture is shown in Figure 2. A workflow is generated within the Workflow Composer, based on experimentation in the Sand Box. This workflow is then deployed to a BPEL engine, which returns a URL pointing to the WSDL for that workflow. This URL is then transmitted to the GES DISC via a Web Services request along with a specification of the data to be mined, such as the dataset to be mined and temporal or spatial constraints. This request is provided to the GES DISC's processing engine, also known as the Simple, Scalable Script-based Science Processor for Measurements (S4PM) (*Lynnes, 2006*) for data mining at the data archive. The S4PM engine is responsible for acquiring the data from the archive, then executing the requested workflow on each input file. However, note that the workflow itself is not provided to S4PM. Instead, S4PM uses the supplied WSDL URL to fetch the WSDL document and then invoke the corresponding Web Service in the BPEL engine, supplying the full path of the input file. The BPEL engine turns around to invoke the atomic Web Services in the proper order at the GES DISC.

#### 4.4.5 Current Status

A deadline of September 30, 2006 was established to implement a prototype that followed the proposed architecture (Figure 2). An initial prototype addressed the scientific problem (cloud mask) described in the demonstration section above.

The GES DISC has successfully hosted the mining web services developed at ITSC by deploying them at the GSFC's web service container. WSDLs for the web services hosted at the GSFC have been made available. ITSC has successfully tested those services using simple clients generated using the WSDLs. Externally generated workflow descriptions have been created and successfully executed through service interactions with the GES DISC.

## 5 Conclusions/Results

Initial prototyping efforts have been successful and encouraging for the use of web service workflows in providing a mechanism where remote users will be able to deploy custom data analysis solutions at data repositories. This approach will allow users to run these analyses on large amounts of data available at the repositories, while still maintaining control on the details of what services are used and the parameters employed. The solution workflows themselves can then be exposed as web services for other researchers to further incorporate in other studies.

Additional investigation is underway to prototype how effective this approach can be for more complicated data analysis workflows.

## 6 Acknowledgments

We gratefully acknowledge funding from NASA's ACCESS program for this project. In addition, ADaM was originally developed under a NASA Research Announcement, and later refactored as part of a NASA AISRP grant.

## 7 References

ActiveBPEL, The Open Source BPEL Engine, <http://www.activebpel.org/>, accessed Jun, 2006

AXIS, Apache AXIS SOAP Implementation: <http://ws.apache.org/axis/>, accessed October 2006.

ECHO, NASA's Earth Observing System Clearinghouse: <http://www.echo.nasa.gov>, accessed October 2006.

Hinke, T., J. Rushing, S. Kansal, S. J. Graves, H. S. Ranganath, and E. Criswell, "Eureka Phenomena Discovery and Phenomena Mining System," presented at AMS 13th International Conference on Interactive Information and Processing Systems (IIPS) for Meteorology, Oceanography and Hydrology, 1997a.

Hinke, T., J. Rushing, H. S. Ranganath, and S. J. Graves, "Target-Independent Mining for Scientific Data: Capturing Transients and Trends for Phenomena Mining," presented at Proceedings Third International Conference on Data Mining (KDD-97), Newport Beach, California, 1997b.

Hinke, T., J. Rushing, H. S. Ranganath, and S. J. Graves, "Techniques and Experience in Mining Remotely Sensed Satellite Data," *Artificial Intelligence Review: Issues on the Application of Data Mining*, vol. 14, pp. 503-531, 2000.

Lynnes, C. and R. Mack, "KDD Services at the Goddard Earth Sciences Distributed Active Archive Center," in *Data Mining for Scientific and Engineering Applications*, R. L. e. a. Grossman, Ed.: Kluwer Academic Publishers, 2001, pp. 165-181.

Lynnes, C. S., 2006. The Simple, Scalable, Script-based Science Processor (in press), in *Earth Science Satellite Remote Sensing*, Springer-Verlag.

Ramachandran, R., H. Conover, S. J. Graves, and K. Keiser, "Algorithm Development and Mining (ADaM) System for Earth Science Applications," presented at Second Conference on Artificial Intelligence, 80th AMS Annual Meeting, Long Beach, Long Beach, CA, 2000.

Ramachandran, R., S. A. Christopher, S. Movva, X. Li, H. T. Conover, K. R. Keiser, S. J. Graves, and R. T. McNider, "Earth Science Markup Language: A Solution to Address Data Format Heterogeneity Problems in Atmospheric Sciences," *Bulletin of the American Meteorological Society*, accepted 2004a.

Ramachandran, R., S. Graves, H. Conover, and K. Moe, "Earth Science Markup Language (ESML): a solution for scientific data-application interoperability problem," *Computers & Geosciences*, vol. 30, pp. 117-124, 2004b.

Rushing, J., R. Ramachandran, U. Nair, S. Graves, R. Welch, and A. Lin, "ADaM: A Data Mining Toolkit for Scientists and Engineers," *Computers & Geosciences*, vol. 31, pp. 607-618, 2005.

SOAP, SOAP::Lite for Perl: <http://www.soaplite.com/>, accessed October 2006.

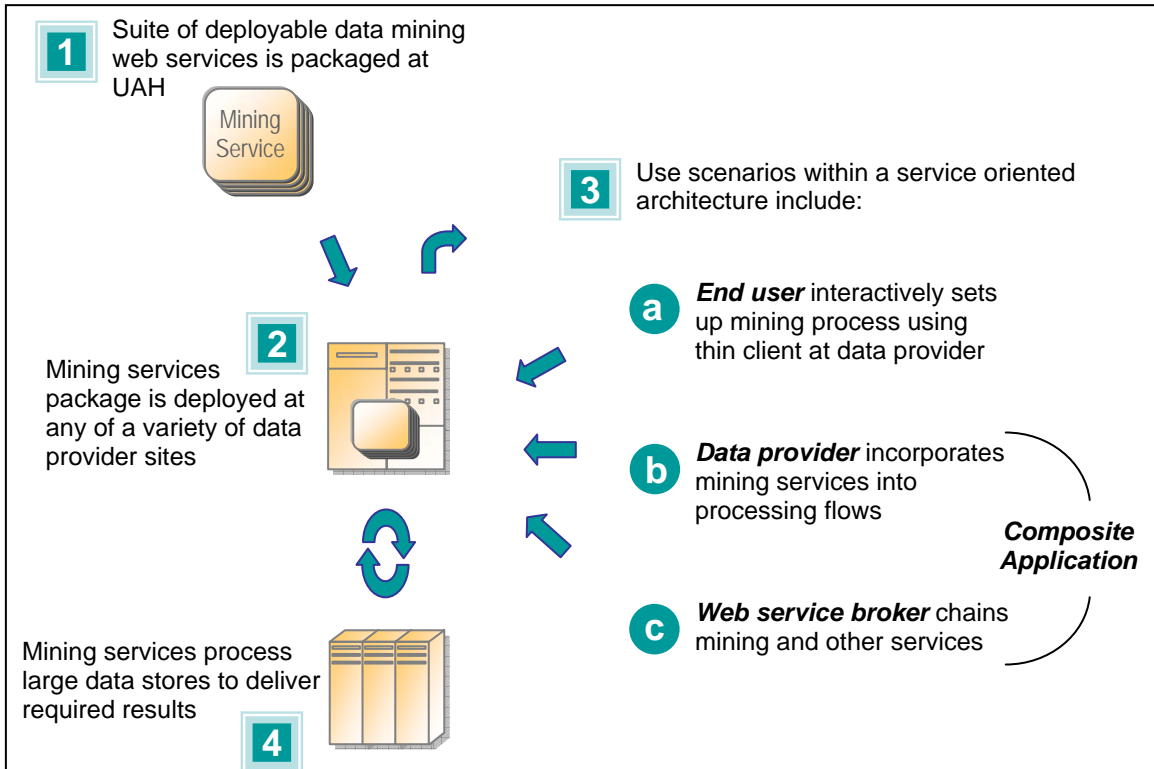


Figure 1: Overall Architecture of Mining Web Services

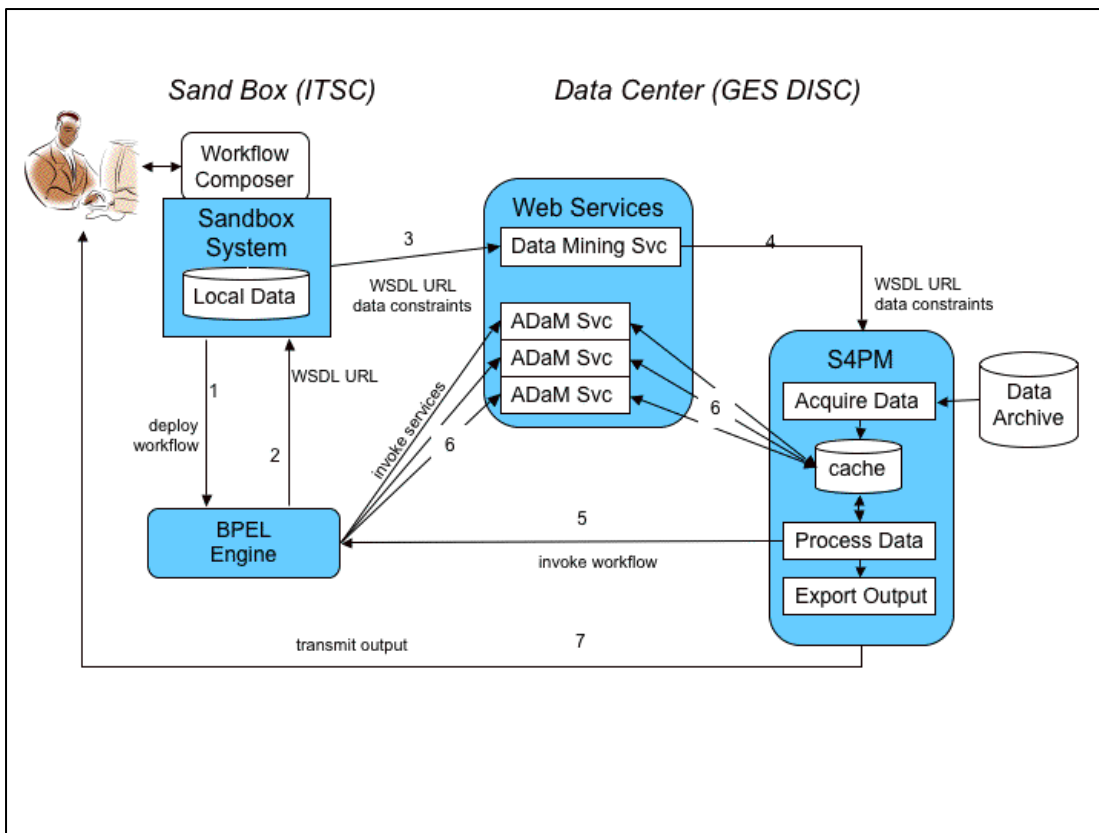


Figure 2: Demonstration Architecture