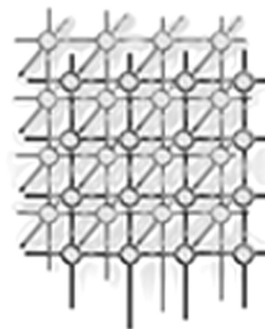


Virtual data Grid middleware services for data-intensive science



Yong Zhao^{1,*,\dagger}, Michael Wilde², Ian Foster^{1,2}, Jens Voeckler¹,
James Dobson³, Eric Gilbert⁴, Thomas Jordan⁴
and Elizabeth Quigg⁴

¹*Department of Computer Science, University of Chicago, Chicago, IL 60637, U.S.A.*

²*Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL 60439, U.S.A.*

³*Department of Psychological and Brain Sciences, Dartmouth College, Hanover, NH 03955, U.S.A.*

⁴*Fermi National Accelerator Laboratory, Batavia, IL 60510, U.S.A.*

SUMMARY

The GriPhyN virtual data system provides a suite of components and services for data-intensive sciences that enables scientists to systematically and efficiently describe, discover, and share large-scale data and computational resources. We describe the design and implementation of such middleware services in terms of a virtual data system interface called Chiron, and present virtual data integration examples from the QuarkNet education project and from functional-MRI-based neuroscience research. The Chiron interface also serves as an online ‘educator’ for virtual data applications. Copyright © 2005 John Wiley & Sons, Ltd.

KEY WORDS: virtual data; workflow; data Grid; data integration; portal; provenance

INTRODUCTION

Next-generation data-intensive applications pose unprecedented demands on information technology. Communities of users distributed worldwide need to extract information from large collections of data, and would like to share and reuse both data products and the resources needed to produce and store those data products [1,2].

*Correspondence to: Yong Zhao, Department of Computer Science, University of Chicago, 1100 East 58th Street, Chicago, IL 60637, U.S.A.

[†]E-mail: yongzh@cs.uchicago.edu

Contract/grant sponsor: National Science Foundation GriPhyN project; contract/grant number: ITR-086044

Contract/grant sponsor: National Science Foundation (QuarkNet project)

Contract/grant sponsor: Office of High Energy Physics, Office of Science, U.S. Department of Energy



So-called virtual data mechanisms virtualize data with respect to location, representation, and materialization. A common data model is used to describe data types and representations, and the recipes for derivations of data are specified in a declarative manner. Requests for data products can be transparently mapped into computation and/or data access operations across multiple Grid computing and storage [3].

The GriPhyN virtual data system (VDS) [4] implements such virtual data mechanisms. Definition and query operations are expressed in a virtual data language (VDL), and information about data and computational procedures is stored in a virtual data catalog (VDC).

VDS has been applied to a variety of data-intensive applications, including physics event simulation [5], galaxy cluster finding [6], genome analysis [7], and biomedical image analysis [8]. However, the initial virtual data interfaces (a Java API with command-line wrappers), while powerful for scripting-level integration, proved cumbersome for both deeper application integration and direct end-user use. Many user communities want to both integrate virtual data capabilities more deeply into science applications and to encapsulate the tasks of setting up and configuring the VDS and its associated Grid compute and storage resources. End users need an interactive environment in which they can easily discover and share virtual data products, compose workflows, and monitor workflow executions across multiple Grid sites.

These considerations have led us to develop, apply, and evaluate a set of interfaces that provide more powerful access to VDS data structures. On top of these, we developed the Chiron portal for convenient interaction with the VDS. Chiron allows users to:

- manage user accounts and track user activities;
- describe and publish applications and datasets;
- discover, validate, share, and reuse applications and datasets published by other users;
- configure and discover Grid resources;
- compose workflows and dataset derivations;
- execute workflow either locally or in a Grid environment and monitor job progress; and
- record and query provenance information.

With these enhancements, the VDS provides both user-level and service (middleware)-level functionality that portal users can integrate to build customized virtual data interfaces for specific user communities and science applications.

The rest of this paper is organized as follows. We first briefly review the GriPhyN VDS and VDL and then describe the Chiron system architecture. Next, we talk about the dimensions of the virtual data discovery space and then present virtual data integration examples from the QuarkNet project and from functional magnetic resonance imaging (fMRI) experiments. We also describe our work on dataset type and dataset iteration, and lastly we compare with related work and draw our conclusions.

This paper builds on and extends an earlier paper [9], updating its description of various system components and detailing the mechanisms used to organize and access datasets.

VDS

VDS supports the description, discovery, and reuse of information on how data is derived by computations. A common data model is used to describe datasets and procedures. This data model



defines representations for information about datasets, transformations, derivations, invocations, and metadata.

A dataset is a unit of data managed by the VDS and is associated with a type. A transformation (TR) is a typed definition for a computation procedure, which takes certain types of datasets as inputs and outputs. A compound TR can be defined to comprise other TRs in an acyclic fashion; a derivation (DV) instantiates a TR by giving the actual arguments and other necessary information, a DV record can serve both as a historical record of what was done and as a recipe for operations that can be performed in the future; an invocation (IV) is the actual execution of a DV. Datasets can be replicated at multiple locations for performance and availability purposes, and IV keeps track of the specific replicas used or produced in the execution. Datasets and transformations can be annotated with metadata (attributes) to facilitate discovery, access, and composition.

Virtual data lifecycle

The virtual data lifecycle begins when a virtual data definition, described in VDL, is entered in the VDC. This definition can then be discovered by virtual data query operations, and the associated scientific analysis procedures can be executed and validated. New algorithms and procedures can be derived from the knowledge gathered; another round of the lifecycle starts when the derived virtual data is published into the VDC or recorded by the VDS.

Virtual data API

In our earlier development of VDS, we only provided users with a limited set of command-line interfaces. These interfaces limited the applicability of the system as only specific functionalities were exposed. Subsequently, we have defined and developed a complete set of Java APIs so that users can have fine-grained control over the system, integrate VDS with their own applications and analysis environments, and build virtual data services and portals based on the APIs.

This virtual data API provides support for:

- system property setup such as database backend/driver configuration, schema locations, and logging;
- VDL conversion and parsing, i.e. conversion between textual and XML syntaxes;
- VDC manipulation, i.e. support for inserting, updating, deleting, and searching virtual data in the VDC;
- querying and editing the catalogs that VDS maintains for mapping logical to physical entities;
- Grid resource management;
- workflow composition, execution, and monitoring; and
- provenance tracking.

These interfaces allow customized implementations for certain components. For example, we implement a number of database backends that variously support file-based, relational database, and native XML database storage, all conforming to the database interface. Similarly, the site selection interface allows developers to enforce Grid site selection policies. The API also makes it easy to develop and deploy portals, and to integrate VDS into specific scientific domains.

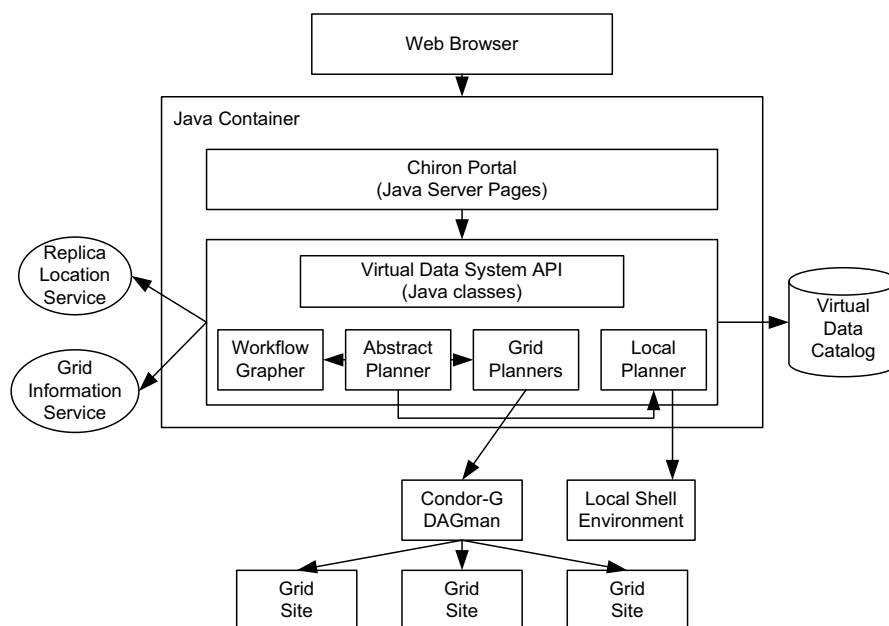


Figure 1. Chiron system architecture.

CHIRON SYSTEM ARCHITECTURE

The Chiron portal is an example of a system implemented in terms of our VDS APIs. The Chiron implementation combines those APIs and libraries with standard commodity Web technologies, databases, and Grid technologies. We show the system architecture in Figure 1.

A Chiron client only needs a Web browser that supports SSL for security, DHTML and Java scripts for interactivity, and cookies for customization. The Chiron server comprises an Apache Tomcat Web server, Java Server Pages (JSP), the VDS API and libraries (jar files), backend databases for the VDS, GraphViz for graph visualization, workflow engines for local and Grid execution, and Condor-G [10] for Grid job submission and execution.

Chiron also uses the open source Apache Axis SOAP toolkit to provide SOAP-based Web services interfaces to the underlying VDS. Chiron is configured to start Axis automatically and the Axis servlet listens on `http://localhost:8080/chiron/services`.

System configuration

The Chiron portal code and libraries are deployed into a Tomcat servlet engine. Chiron extends the VDS 'properties' file to configure various runtime parameters, such as database driver and schema, monitoring and discovery service (MDS) [11] and replica location service (RLS) [12] host information,



file-transfer mode, logging options, and scheduler-specific configuration. Backend database support is provided by a uniform data-management interface on top of JDBC that leverages Java reflection and dynamic class loading to make multiple database management systems available transparently to the higher-level code.

User management

All portal users must create a user account and should have a valid Grid certificate. We use user ID and password authentication via secure HTTP connections, so that the client side only needs a Web browser to access the portal. Different users are assigned different roles at the portal, such as guest, administrator, and super user, and different levels of controls are associated with each different role. The portal also keeps track of sessions and user preferences. User activities are tracked as provenance data in the VDC, where analysis patterns could be uncovered by leveraging data-mining approaches.

Job submission

Once users have registered and published data derivation recipes in the VDC, they can request that these virtual datasets be materialized. We term the process of mapping these requests to the execution of computation procedures ‘planning’ as it is suggestive of query planning in database systems. Following a request, the VDS finds all datasets and procedures required to derive the requested dataset, and generates an ‘abstract’ derivation workflow—abstract in the sense that it is not tied to specific executables, dataset replicas, or Grid sites. This workflow is represented in an XML format that can be interpreted and executed by different planners and workflow engines.

Chiron supports two different planners. The *local planner* takes a workflow definition and produces a multi-level shell script that invokes the required graph of programs. This planner is designed to run derivations on a single host and thus is useful for initial debugging of applications and for determining whether specific errors are caused by the Grid, VDS, or the application.

The *Pegasus Grid planner* [13] uses information from the Grid information service (MDS) and RLS to convert the abstract plan into a concrete execution plan. That concrete plan is then submitted to DAGMan and executed on remote Grid sites via Condor-G and Globus. A Grid site usually has compute elements (CEs) managed by a local scheduler (Condor, PBS, LSF, etc.), and storage elements (SEs) accessed via GridFTP.

The portal also displays execution progress and job status once the workflow is converted into a concrete plan and executed. Upon successful completion, requested datasets can be displayed as images or text, if their format permits.

Workflow visualization

Chiron uses GraphViz, a graph-drawing tool developed by AT&T Laboratories, to visualize workflows. The XML representation of a workflow is first converted into GraphViz’s ‘dot’ format and then transformed into a displayable image. Some sample workflows are presented in Figures 3 and 4. Parallelograms and ellipses represent datasets and computation procedures, respectively. When displayed in Chiron, a user can click on a dataset node to show all derivations that process the dataset;



clicking on a procedure node displays detailed information about the corresponding transformation. These graphs provide insights into both individual steps and the workflow process in the large.

Web services interface

While the VDS Java API provides a powerful interface, it is suitable only for programmatic integration. However, many science communities require deeper integration of virtual data capabilities into their applications and problem-solving environments, and the encapsulation of system setup and Grid resource configurations. The Chiron servlet and JSP scripts expose higher-level functionalities and facilities for portal and Web interface development.

Chiron also provides a Web service interface for more flexible, loosely coupled integration. This interface is provided via an embedded Apache Axis toolkit. A set of virtual data Web services is deployed in the Axis engine, which provides the same functionalities as the JSP pages for virtual data discovery, composition and integration. The service client stub classes have been generated from the WSDL service description and distributed with the portal, but other Web services invocation mechanisms can be used to communicate with the services.

This combination of VDS API, service-oriented JSP scripts, and Web services interface allow the Chiron portal to provide science communities and applications with a suite of Grid middleware services, of different levels and granularities, to explore the power of virtual data and Grid computing.

VIRTUAL DATA DISCOVERY

The two keys to virtual data reuse and sharing are efficient description and discovery. To this end, Chiron allows users to search the VDC for virtual data entities along three different dimensions: namespace and/or attributes, metadata annotations, and provenance history. Definition namespaces organize virtual data procedures into a structured hierarchy; metadata associates arbitrary annotations with virtual data products; provenance information records the derivation lineage of virtual data products.

Definition

A definition is a declarative specification of some data analysis procedure. It can be either a transformation ('TR' in VDL), which is similar to a function definition in that it defines a function name and formal parameters, or a derivation ('DV' in VDL), which acts like a function call, with the actual parameters (scalars and data files) identified [3,4].

Definitions can be organized hierarchically using namespaces, which might be associated with, for instance, institutions, projects, groups, or users. Chiron displays VDL definitions in a tree view that can be customized to show namespaces, transformations, and/or derivations (Figure 2). A user can explore different namespaces and display the detailed definitions of each transformation and derivation by exploring the tree. Definitions can also be shown in VDL [3,4] for advanced users.

Chiron also provides a search bar for quick exploration. A user can search for transformations, derivations, datasets, and metadata using key-word and wildcard searches. These definitions provide an interface-level understanding of data analysis procedures.

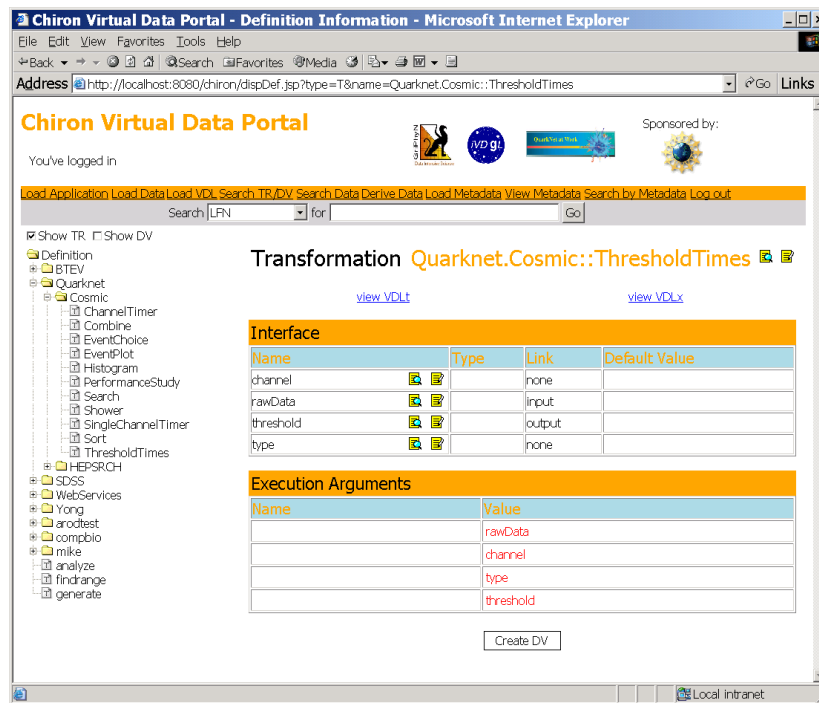


Figure 2. Chiron screenshot.

Metadata

Metadata annotation is essential to facilitate the discovery and understanding of the datasets and procedures stored in the VDC. We allow annotations to be associated with five categories of virtual data objects: datasets, transformations, derivations, parameters, and the nested calls in compound transformations. (A compound transformation consists of a graph of calls to other transformations.) Annotations are represented by tuples of the form (attribute_name, type, attribute_value), where type can be integer, float, string, date, or boolean.

Metadata about a dataset might include cataloging information, such as description, creator, creation date, and publisher; content information such as data format, structure, and size; and provenance information, i.e. how data is acquired, transformed, and transferred. For code and procedures, in addition to information about their inputs, outputs, and execution environments, metadata might describe functionality, prerequisites, and parameter constraints.

We define a metadata query language to search for virtual data objects that have specific metadata. A user does not need to specify the types of the attributes being searched: a query parser takes care of type conversion and parses it into an abstract syntax tree, which is then optimized and translated into



SQL or, if available, XML XPath or XQuery statements. The query statements are submitted directly to the backend database and the results are interpreted and then returned in the appropriate format.

The Chiron portal provides Web interfaces for associating metadata with different virtual data objects, displaying metadata about a specific virtual data objects, and searching for virtual data objects with specific metadata.

Provenance

Provenance information captures the derivation history of data [14]. In our case, it consists of an ‘invocation’ record for a data derivation, which records when, how, and where the data is produced, what kind of resources are used, and so on. Combined with the abstract workflow produced for the data derivation, provenance information gives the audit trail of the data production process. A user can use this information to repeat a derivation, to explain or validate a result, and so forth.

An invocation record is recorded for each task that is executed within a workflow. These records can be used to provide reports on data quality, runtime, resource consumption, or even computation anomalies; to re-execute past workflow or to modify/improve future workflow; and to answer *ad hoc* queries.

VIRTUAL DATA INTEGRATION

Science communities and scientific applications must often integrate virtual data into an existing analysis environment, a portal, or a Web interface. Chiron facilitates the development of application-specific portals and Web interfaces, and provides reusable core services for accessing virtual data and Grid resources. We illustrate virtual data integration scenarios drawn from the QuarkNet science education project and a neuroscience research project.

QuarkNet cosmic ray application

The Chiron portal is being applied in the QuarkNet project [15], an NSF- and DOE-funded project that aims to create a research community of physicists, high-school teachers, and their students. QuarkNet engages teachers and students in scientific investigations about the structure of matter and the fundamental forces of nature. Students learn fundamental physics as they analyze live, online data and participate in inquiry-oriented investigations, and teachers join research teams with physicists at local universities or laboratories. Currently, the project involves six large experimental physics collaborations, 52 research groups, 158 physicists, 514 high-school teachers, and thousands of students.

The QuarkNet Cosmic Ray Project is an experimental activity designed to give students hands-on experience in scientific collaboration and discovery. Students learn to construct, deploy and test cosmic-ray detectors, and analyze real-time data gathered by the detectors. About 40 detectors have been installed in high-school classrooms across the U.S. QuarkNet helps students understand and explore the origin of highly energetic cosmic rays, the source of their energy, and the clues they hold about the early universe.

QuarkNet has created a specialized Web site and Web interfaces for science education [16]. This portal, driven by virtual data middleware, allows students to launch, configure, and control remote

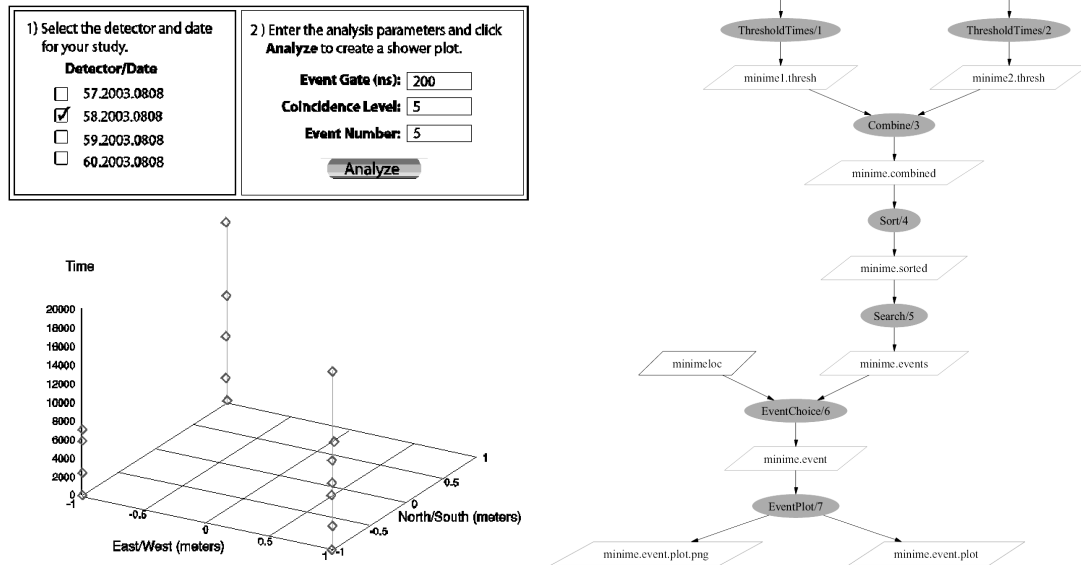


Figure 3. Cosmic-ray-shower study.

applications as though they are using a desktop computing environment. The key requirements for the Web interfaces are to give students the means to: (a) discover and apply datasets, algorithms, and data analysis methods; (b) collaborate by developing new transformations and by sharing results and observations; and (c) learn data analysis methods that will motivate and prepare them for a scientific career. These requirements match those for Chiron, but here it is the Web applications that need to interact with Chiron. Thus, we have encapsulated relevant Chiron portal functions into APIs so that these Web applications can integrate virtual data mechanisms seamlessly into their framework.

For a specific data analysis task, students interact directly with the QuarkNet Web interface, which in turn accesses Chiron services to upload data and Perl code, process data using simple form-based tools, annotate and share data and results, and create posters to describe and share analysis processes. Students can use metadata queries to discover each other's workflows and data products.

In Figure 3, we show the customized Web page for a cosmic-ray-shower study. This study involves a data-analysis process that combines detector and GPS data gathered from distributed geographical locations and determines whether a cosmic-ray shower has occurred in the sampled period and location. The plot is generated by clicking the 'analyze' button, which sends the request to the derivation script in Chiron. The script first generates the abstract workflow (shown to the right in Figure 3), then calls the shell planner to create a concrete, executable workflow, and finally executes the corresponding graph of analysis code. A link to a visual rendering of the workflow graph is returned. These steps all occur transparently: the user only interacts with the QuarkNet Web page.

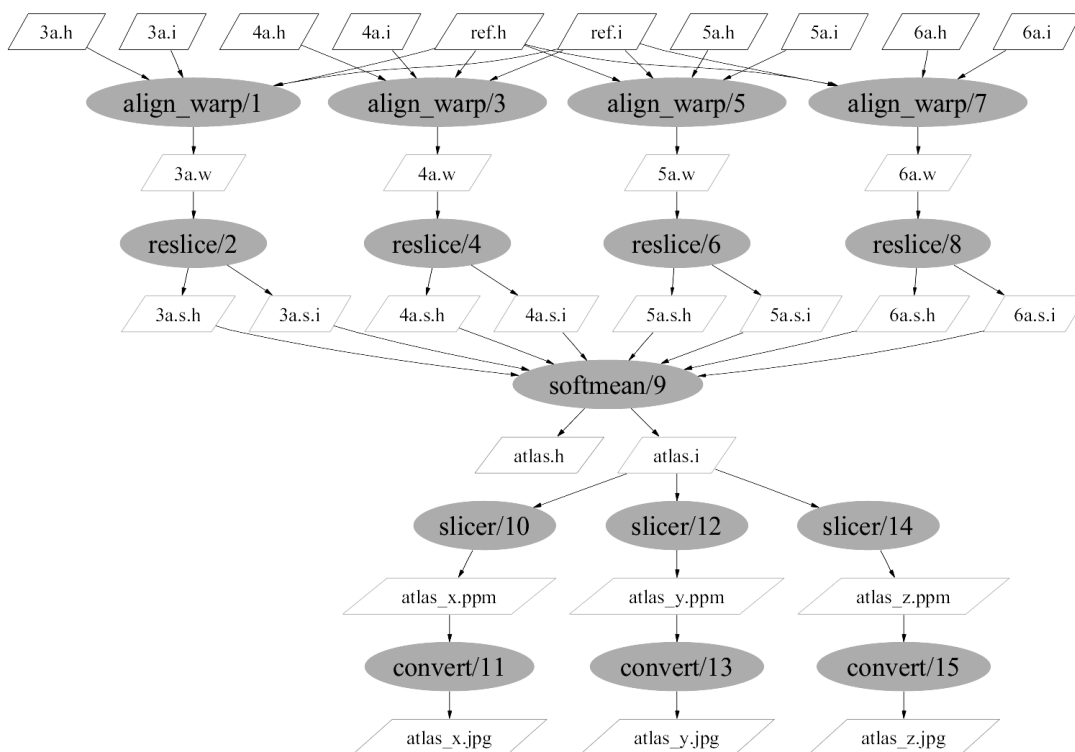


Figure 4. fMRI workflow.

Chiron also provides mechanisms to generate application-specific Web pages automatically, by using HTML templates, transformation definitions, and annotations. For example, the title of a transformation can be rendered into the title of the Web page, default/eligible values for parameters can be rendered as options, input datasets can be filtered by dataset types, and output filenames can be generated by context.

Cognitive neuroscience application

Automated data provenance and workflow management is a critical requirement for the fMRI research community [17]. Our virtual data middleware has been used to create population-based atlases from the fMRI Data Center's archive [18] of high-resolution anatomical magnetic resonance data, using a multi-stage I/O intensive pipeline that 'warps' the anatomical features of the brains of multiple subjects into a standard space.

The VDL workflow for the creation of such a brain 'atlas' (Figure 4) involves several derivations that employ the automated image registration (AIR) suite [19,20] to create an 'averaged' brain from



a collection of high-resolution anatomical data. The workflow warps each anatomical volume to a reference volume, using a high-order nonlinear multi-parameter model. The brain structures are warped to define a similar shape, then each volume is resliced using the warp values, and the resulting high-resolution volumes are averaged using the AIR suite's 'softmean' application.

In the context of this atlas-generation use case, we experimented with applying the middleware services of the VDS to various science-driven discovery problems. A few of these are illustrated in the following queries and corresponding responses.

- Which TRs can process a subject image?

```
Q: xsearchvdc -q tr_meta dataType subject_image input
A: fMRIDC.AIR::align_warp
```

- Which TRs can create an ATLAS?

```
Q: xsearchvdc -q tr_meta dataType atlas_image output
A: fMRIDC.AIR::softmean
```

- List anonymized subject-images for young subjects. This query searches for files based on their metadata.

```
Q: xsearchvdc -q lfn_meta dataType subject_image privacy
      anonymized subjectType young
A: 3472-4_anonymized.img
```

- For a specific patient image, 3472-3, show all DVs and files that were derived from this image, directly or indirectly.

```
Q: xsearchvdc -q lfn_tree 3472-3_anonymized.img
A: 3472-3_anonymized.img
   3472-3_anonymized.sliced.hdr
   ...
   3472-3_anonymized.sliced.img
```

Similar queries can be used for data validation processes and to construct scripts that compose new virtual data workflows based on existing ones, and then catalog, run, and record them.

MANAGING DATASET COLLECTIONS

A commonly observed pattern in scientific computation is that a procedure to process a large dataset is usually specified by higher-level scripts that iterate over the dataset, applying a transformation to each sub-component. For example, in the fMRI spatial normalization application (AIRSN), a single study consists of groups of subjects, each subject has a series of fMRI scans (runs), and each run contains a series of images (volumes), where a script would iterate over each volume in a run, each run in a subject, and each subject in a group.



Another observation is that data are often represented in different physical formats, which may change over time. Computational procedures designed for a given format may no longer work correctly once formats change. Such format dependency prevents code reuse over datasets that are conceptually identical but have different physical representations.

To address these two issues, the XML dataset typing and mapping (XDTM) model [21] was developed. XDTM allows virtual datasets to be defined, and procedures to be written that select subsets of the datasets and iterate over those subsets, in a manner that is independent of the datasets' concrete physical representations. XDTM employs a two-level description of datasets, defining separately via a type system (XML Schema) the abstract structure of datasets, and the mapping of that abstract data structure to physical representations. The use of XML Schema as a type system has the benefit of supporting the same powerful standardized query language that we use in our selection methods.

We have implemented a prototype of the XDTM model by extending the current VDS, and successfully applied the system to the fMRI AIRSN workflow.

RELATED WORK

The GridAnt [22] system provides a flexible workflow model based on Jakarta Ant and Grid execution. However, it does not address provenance tracking, planning support, or the declarative representation of application logic, nor does it contain the sophisticated mechanisms that can transparently perform a transformation on a set of input data objects in a location-independent manner, anywhere in a distributed Grid or local desktop environment.

The SAM system deployed by the Collaboration for Multiscale Chemical Science [23] is based on the uniform action model of the WebDAV protocol [24], and has a well-developed 'laboratory notebook' model of provenance tracking. However, its middleware approach for Grid execution has the same location-dependence limitations as GridAnt.

The Grid Access Portal for Physics Applications [25] supports job submission, job-status checking, and resource management. Each application is packaged as a *notebook* comprising Web pages and editable parameterized scripts. Notebooks can be published and stored in Web-based archives for others to retrieve and modify. The notebook concept is similar to our transformation and derivation descriptions, but our virtual data descriptions are more general and are decoupled from physical entities to support community-wide sharing and reuse.

The MyGrid project [26] has developed powerful distributed execution semantics and enactment engines, visual workflow specification capabilities, and provenance tracking. However, these capabilities are only provided for the relatively 'regular' execution model of Web services, but not for the standard application programs that constitute the majority of today's scientific computing toolkits.

CONCLUSIONS

We have described a VDS interface, Chiron, that allows users to describe, discover, share, and reuse scientific processes and data. Chiron provides user-level Web interfaces, service-oriented scripts, and Web-services-based middleware functionalities that facilitate integration of virtual data mechanisms



into specialized data-analysis applications and problem solving environments. Chiron has been applied successfully within the QuarkNet education project and tested in the context of a fMRI research project.

Chiron also serves as a virtual data training tool. We have loaded procedures and workflows from several virtual data applications, including high-energy physics, cosmic-ray detection, bioinformatics, and fMRI, as examples into the portal, thus allowing first-time and novice virtual data users to explore core virtual data concepts and technologies without being distracted by software configuration details.

We have also extended the virtual data language to use XDTM mechanisms to abstract dataset types from their physical representations, and to enable selection and iteration on large composite datasets.

In future work we plan to leverage knowledge representation techniques to enhance process definition and metadata annotation, and to enable semantics-driven virtual data discovery, composition, and integration.

ACKNOWLEDGEMENTS

This work is supported in part by the National Science Foundation GriPhyN project under contract ITR-086044. The QuarkNet project is supported in part by the National Science Foundation and the Office of High Energy Physics, Office of Science, U.S. Department of Energy. The GriPhyN VDS was implemented by Gaurang Mehta, Karan Vahi, Jens Voeckler, and Yong Zhao. Yong Zhao implemented the Chiron portal; Tom Jordan, Liz Quigg, Paul Nepwoda, Evgeni Peryshki, Nick Dettman, Yun Wu, and Eric Gilbert implemented the QuarkNet Cosmic Ray Collaboration portal. James Dobson provided the fMRI workflows and assisted in their development and testing. The Pegasus planner is the work of Ewa Deelman and her group at the USC Information Sciences Institute.

REFERENCES

1. Avery P, Foster I. The GriPhyN project: Towards petascale virtual data Grids. <http://www.griphyn.org> [June 2005].
2. Chervenak A, Foster I, Kesselman C, Salisbury C, Tuecke S. The data Grid: Towards an architecture for the distributed management and analysis of large scientific datasets. *Journal of Network and Computer Applications* 2001; **23**:187–200.
3. Foster I, Voeckler J, Wilde M, Zhao Y. The Virtual Data Grid: A new model and architecture for data-intensive collaboration. *Proceedings of the 1st Biennial Conference on Innovative Data Systems Research*, January 2003.
4. Foster I, Voeckler J, Wilde M, Zhao Y. Chimera: A virtual data system for representing, querying, and automating data derivation. *Proceedings of the 14th Conference on Scientific and Statistical Database Management*. IEEE Computer Society Press: Los Alamitos, CA, 2002.
5. Arbee A, Avery P, Bourilkov D, Cavanugh R, Graham G, Katageri S, Rodriguez J, Voeckler J, Wilde M. Virtual data in CMS production. *Proceedings of Computing in High Energy and Nuclear Physics*, March 2003. Available at: <http://arxiv.org/abs/cs.DC/0306009> [June 2005].
6. Annis J, Zhao Y, Voeckler J, Wilde M, Kent S, Foster I. Applying Chimera virtual data concepts to cluster finding in the Sloan Sky Survey. *Proceedings of Supercomputing 2002 (SC2002)*, November 2002. ACM Press: New York, 2002.
7. Rodriguez A, Sulakhe D, Marland E, Nefedova N, Wilde M, Maltsev N. Grid enabled server for high-throughput analysis of genomes. *Proceedings of the Workshop on Case Studies on Grid Applications*, March 2004.
8. Leung K *et al.* Analysis of serial MR images of joints. *Proceedings of the IEEE International Symposium on Biomedical Imaging (ISBI)*. IEEE Press: Piscataway, NJ, 2004; 221–224.
9. Zhao Y, Wilde M, Foster I, Voeckler J, Dobson J, Jordan T, Quigg E. Grid middleware services for virtual data discovery, composition, and integration. *Proceedings of the 2nd Workshop on Middleware for Grid Computing*. ACM Press: New York, 2004.
10. Frey J, Tannenbaum T, Foster I, Livny M, Tuecke S. Condor-G: A computation management agent for multi-institutional Grids. *Cluster Computing* 2002; **5**(3):237–247.
11. Czajkowski K, Fitzgerald S, Foster I, Kesselman C. Grid information services for distributed resource sharing. *Proceedings of the 10th IEEE International Symposium on High-Performance Distributed Computing (HPDC-10)*, August 2001. IEEE Press: Piscataway, NJ, 2001.



12. Chervenak A *et al.* Giggie: A framework for constructing scalable replica location services. *Proceedings of Supercomputing 2002 (SC2002)*, November 2002.
13. Deelman E *et al.* Pegasus: Mapping scientific workflows onto the Grid. *Proceedings of the 2nd EU Across Grids Conference*, 2004.
14. Buneman P, Khanna S, Tan W-C. Why and where: A characterization of data provenance. *Proceedings of the International Conference on Database Theory*, 2001.
15. The QuarkNet Project. <http://quarknet.fnal.gov/> [June 2005].
16. Bardeen M, Gilbert E, Jordan T, Nepywoda P, Quigg E, Wilde M, Zhao Y. The QuarkNet/Grid collaborative learning e-Lab. *Proceedings of the 2nd International Workshop on Collaborative and Learning Applications of Grid Technology and Grid Education*. IEEE Computer Society Press: Los Alamitos, CA, 2005.
17. Van Horn J, Dobson J, Woodward J, Wilde M, Zhao Y, Voeckler J, Foster I. Grid-based computing and the future of neuroscience computation. *Methods in Mind*. MIT Press: Cambridge, MA, 2006.
18. Van Horn J. Online availability of fMRI results images. *Journal of Cognitive Neuroscience* 2003; **15**(6):769–770.
19. Woods R. Automated image registration: I. General methods and intrasubject, intramodality validation. *Journal of Computer Assisted Tomography* 1998; **22**:139–152.
20. Woods R. AIR 5 suite. <http://bishopw.loni.ucla.edu/AIR5> [June 2005].
21. Moreau L, Zhao Y, Foster I, Voeckler J, Wilde M. XDTM: XML dataset typing and mapping for specifying datasets. *Proceedings of the European Grid Conference on Advances in Grid Computing (EGC 2005)*, Amsterdam, The Netherlands, February 2005 (*Lecture Notes in Computer Science*, vol. 3470). Springer: Berlin: 2005.
22. Amin K *et al.* GridAnt: A client-controllable grid workflow system. *Proceedings of the 37th Hawaii International Conference on System Sciences*. IEEE Computer Society Press: Los Alamitos, CA, 2004.
23. Myers J *et al.* A collaborative informatics infrastructure for multi-scale science. *Proceedings of the Challenges of Large Applications in Distributed Environments (CLADE) Workshop*, June 2004. IEEE Computer Society Press: Los Alamitos, CA, 2004.
24. Stein G. Web Digital Authoring and Versioning (WebDAV) Resources Community Web site. <http://www.webdav.org/> [June 2005].
25. Krishan S *et al.* The XCAT science portal. *Proceedings of Supercomputing 2001 (SC2001)*, November 2001. ACM Press: New York, 2001.
26. Goble G, Pettifer S, Stevens R. Knowledge integration: In silico experiments in bioinformatics. *The Grid: Blueprint for a New Computing Infrastructure*, Foster I, Kesselman C (eds.). Morgan Kaufmann: San Mateo, CA, 2004; 121–134.