

SPECIAL ISSUE PAPER

A science data gateway for environmental management

Deborah A. Agarwal¹, Boris Faybishenko¹, Vicky L. Freedman², Harinarayan Krishnan¹, Gary Kushner¹, Carina Lansing², Ellen Porter², Alexandru Romosan¹, Arie Shoshani¹, Haruko Wainwright¹, Arthur Weidmer¹ and Kesheng Wu^{1,*},[†]

¹*Lawrence Berkeley National Laboratory, Berkeley, CA, USA*

²*Pacific Northwest National Laboratory, Richland, WA, USA*

SUMMARY

Science data gateways are effective in providing complex science data collections to the world-wide user communities. In this paper we describe a gateway for the Advanced Simulation Capability for Environmental Management (ASCEM) framework. Built on top of established web service technologies, the ASCEM data gateway is specifically designed for environmental modeling applications. Its key distinguishing features include (1) handling of complex spatiotemporal data, (2) offering a variety of selective data access mechanisms, (3) providing state-of-the-art plotting and visualization of spatiotemporal data records, and (4) integrating seamlessly with a distributed workflow system using a RESTful interface. ASCEM project scientists have been using this data gateway since 2011. Copyright © 2015 John Wiley & Sons, Ltd.

Received 21 July 2015; Revised 31 August 2015; Accepted 1 September 2015

KEY WORDS: data gateway; environmental management; web service

1. INTRODUCTION

Applications in diverse areas such as biology, chemistry, and earth science are collecting large amounts of complex data. Often, there is a worldwide community of users who want to make use of these data sets. The science gateway is an effective technique for concentrating the data sets at a small number of professionally managed sites accessible by the user community as web services [1–5]. Aside from numerous research prototypes, we are aware of a number of different science gateways currently in production at major computer centers such as National Energy Research Scientific Computing Center[‡] [6] and Argonne Leadership Computing Facility (ALCF) [7, 8]. Most of these science gateways treat the data they manage in a relatively simple way, for example, assuming the data to be a sequence of letters or by restricting the data to be from static meshes. In this paper, we describe an implementation of a science data gateway that handles complex geometries and data types associated with spatiotemporal data records. We create a comprehensive data analysis capability by designing a querying scheme, an interactive visual exploration interface, and a way of integrating the portal with a distributed simulation and analysis environment.

Our web-based data service is designed for a specific type of spatiotemporal data from an important application domain. In the next section, we describe this application and the demands on the data management system. Following that, we devote four sections to discuss each of the aforementioned four features: Section 3 covers spatiotemporal data objects to be captured, Section 4 describes how we achieve selective access to spatiotemporal objects, Section 5 discusses plotting and visualizing accessed data objects, and Section 6 explains integration of this data gateway with the ASCEM distributed workflows. A brief summary is given in Section 7.

*Correspondence to: Kesheng Wu, LBNL, Berkeley, CA 94720, USA.

[†]E-mail: KWu@lbl.gov

[‡]<https://www.nersc.gov/users/data-analytics/science-gateways/>

2. APPLICATION DRIVER

The US Department of Energy (DOE) is in charge of the nation's effort to clean up nuclear contamination at various fuel processing sites; remediation of these legacy wastes is one of the most complex and technically challenging cleanup efforts in the world, with costs over the next few decades projected to be \$265–305 billion [9]. To understand the propagation of the contaminants and the effectiveness of the remediation procedures, DOE has made significant investments in developing computational tools to predict the long-term behavior of subsurface contaminant plumes. Our work is part of a project called Advanced Simulation Capability for Environmental Management (ASCEM) [10–12]. One of the key features of ASCEM is the user environment, Akuna, which is a custom-built interface for managing subsurface modeling workflows [13]. Akuna provides users with a range of tools to utilize data with simulations, analyze data, translate conceptual models to numerical models, execute simulations, and visualize results. Additional toolsets provide users with methods for sensitivity analysis, model calibration, and uncertainty quantification. The data gateway described in this paper stores the persistent data used as input by Akuna workflows. It is also a central place for users to retrieve data for further analysis on their own local resources, a platform for interactive remote explorations, as well as a resource for carrying out in-depth analyses. Next, we describe some of the key tasks of ASCEM to illustrate the demand on the data management system.

In environmental management, scientists use numerical models to assess anticipated risks, support remediation and monitoring program decisions, and assist in the design of specific remedial actions for complex systems. These decisions often need to be made with incomplete information and the impact of knowledge gaps needs to be quantified. Scientists rely on numerical models capable of providing the missing information, thus enabling them to quantify the uncertainty and explore different scenarios [14–16]. Figure 1 shows levels of complexity of the different types of problems the ASCEM platform intends to manage. The interdependencies between the various tasks of an ASCEM workflow can be highly dynamic as the simulations progress. The input and output from these modeling operations may be only megabytes in size, but the demand on these data records is dynamic in nature. Having the data constantly available is a good strategy to simplify the data management tasks for these modeling efforts.

The ASCEM project is developing computer software capabilities for both an integrated platform and a high-performance computing (HPC) multiprocess simulator. The integrated platform provides the user interface and tools for end-to-end model development, starting with the definition of the conceptual model, the management of data for model input, the model calibration and uncertainty analysis, and the processing of model output, including visualization. HPC can be a critical part of the ASCEM workflow. Typical HPC resources are in high demand, and therefore, their availability is hard to predict. When they become available, we may not want to spend the time to gather all the necessary data to the computation site. Instead, we centralize the data artifacts into a web-based service. This science data gateway is always available and can be optimized for a wide range of application scenarios.

The Platform and HPC capabilities are being tested and evaluated for environmental management applications in a set of demonstrations as part of the Site Applications Thrust Area activities, with the Savannah River site (F-Area) serving as one of the real-world sites for an end-to-end demonstration of these capabilities [10–12]. This demonstration requires a wide variety of data from both simulation and experimental observations. The data include complex geographical information, which is usually not supported by web-based data services. Effective support for such geographical information and associated operations is the key part of the ASCEM data management gateway.

ASCEM has taken the approach of building a distributed web-based modeling architecture. This web-based distributed modeling approach has also emerged as an effective mechanism to connect the high-quality data available on the web with HPC resources to assist modelers in quantifying uncertainty [15], understanding chemical reactions [4], and so on [17]. ASCEM data gateway needs to support all these diverse sets of operations.

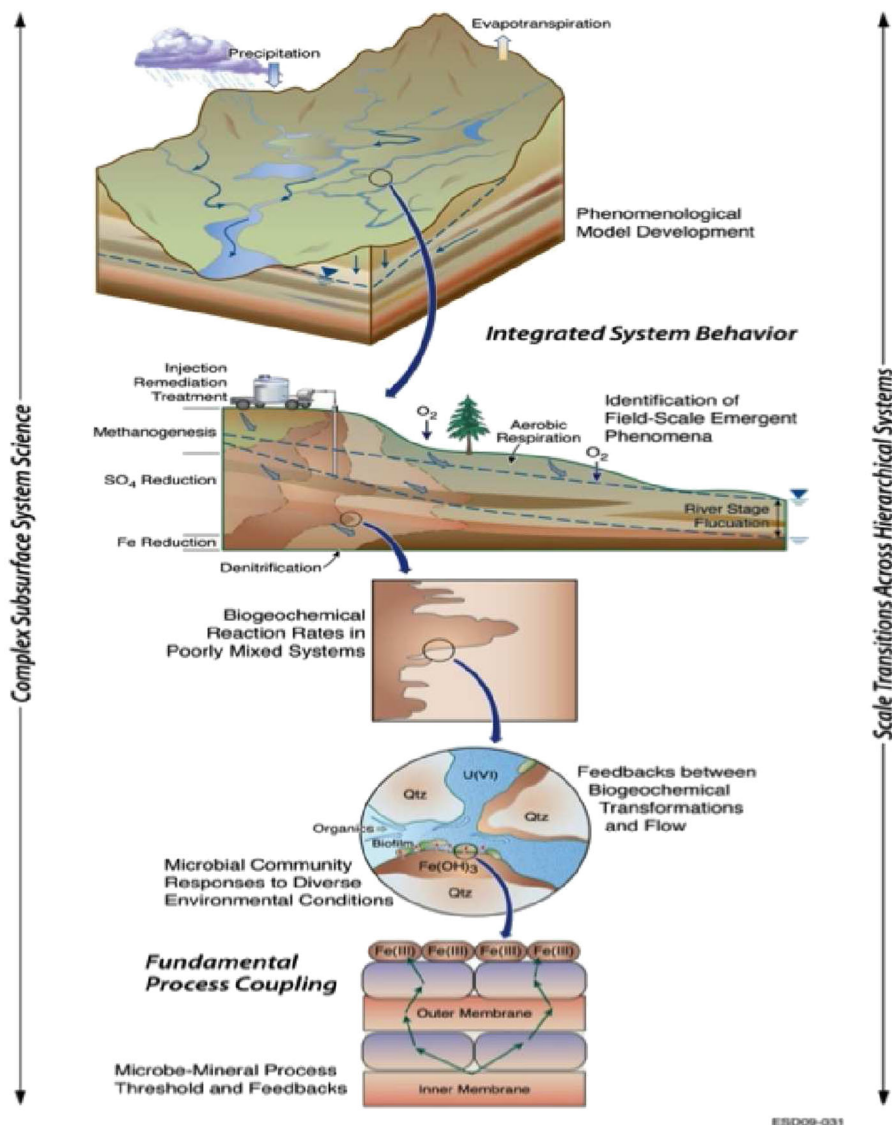


Figure 1. An illustration of problem hierarchy for Advanced Simulation Capability for Environmental Management platform.

3. SPATIOTEMPORAL DATA

The data management toolset is at the core of the ASCEM software framework, consisting of a data store and information management infrastructure, and is accessible by all the toolsets on the ASCEM platform. The data stored in the data management infrastructure include all measured site data. An illustration of the ASCEM data management component is shown in Figure 2. At most environmental management sites, different types of data are typically organized into separate databases such as a well and borehole database, a tank database, and a contaminants database. In addition, many data sets are available as files on individual users' systems or on shared disk systems. However, in the ASCEM project, the application scientists have identified the ability to access collections of these previously disparate data sets to be the most important aspect of the data management system. To achieve this goal, we organized the data sets into a single database containing all data collected at the site relevant to analysis and simulation. Each data object stored in the ASCEM system has additional information describing its type, place of measurement, site

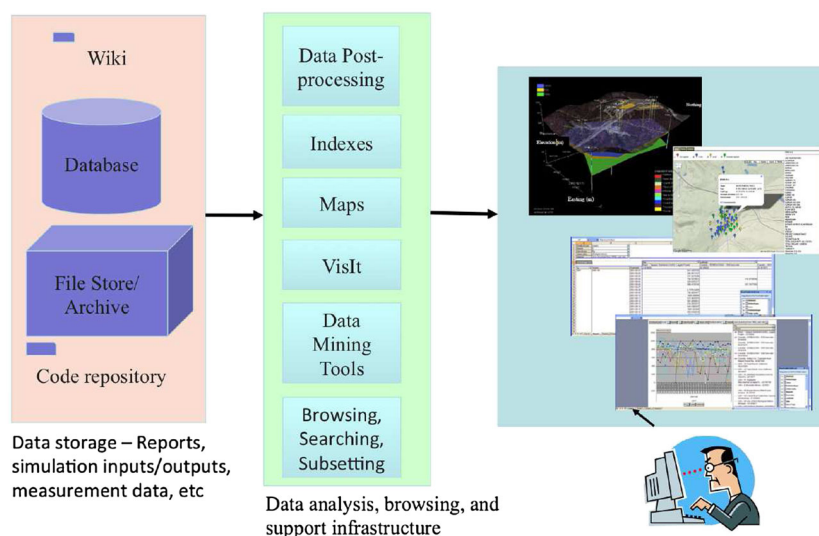


Figure 2. An illustration of the key components of Advanced Simulation Capability for Environmental Management data management.

designation, original source of the data, and format and context-dependent information important for searching.

The Platform core framework enables sharing of data sets, association of analyses with simulation results, flexible data access methods, automation of management of the large volumes of data associated with sensitivity and validation studies, and maintenance and viewing of provenance between the various data items and simulation runs. It is also capable of extracting data from established data management systems at user sites, along with necessary metadata and provenance information.

Data sets involved in environmental management are spatiotemporal in nature. A data set might include points such as locations of wells where measurements were taken, regions such as outlines of buildings, outlines of storage tanks, and other shapes including lakes and rivers. Figure 3 is a part of the Savannah River Site with certain landmarks outlined on the map. Effective support for complex geometric shapes and geographical features is fundamental to the usability of the ASCEM data management system.

The content and the organization of the data records can also critically affect the operations on the data. Common operations on such a set of environmental data might include placing locations of wells on a map, located correctly relative to existing buildings and other landmarks, identifying a known building, finding all the wells containing a certain type of information in or around a building, and so on. An important task our data management system needs to perform is to harmonize the representations of data in order to work with the different coordinate systems for measuring the locations of wells, buildings, and sensors at various sites. There are also variations in representing the shape and outline of buildings, rivers, lakes, and tanks. Our data management system converts all of the varying formats into a common format selected by the ASCEM project.

The measured data is organized into curves with each curve representing a time series or a depth series of a particular measurement type at a particular location. This organization was developed based on user interviews to understand access patterns needed in the system. It supports highly efficient retrieval of the data for plotting, searching, and downloads.

The current implementation of the ASCEM data management system uses a combination of PostgreSQL and PostGIS to support such operations [18–20]. More details of these operations are explained in the next sections.

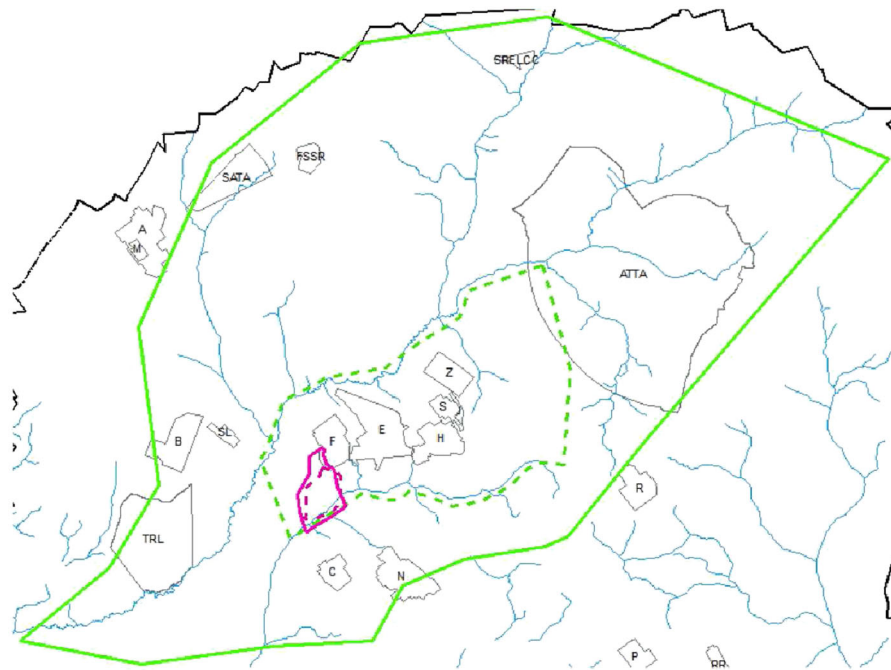


Figure 3. An example of a legacy waste site.

4. SELECTIVE ACCESS TO LARGE DATA SETS

In ASCEM, we take a holistic approach to managing the data records. All model input, output, and control files are managed collectively, versioned and supported by data access tools. This approach reduces the complexity of the data management tasks and makes the modeling task more effective and efficient. The ASCEM data management toolset uses the concept of a ‘collection of data objects’ to support such groupings. Such a collection can be defined statically by groups of application scientists or generated dynamically through selection operators. Our data management toolset also contains a significant amount of metadata on the files and data collections.

The ASCEM data management toolset follows the typical web service approach by offering a command-line interface through a RESTful API, and an interactive web interface [7, 21, 22]. Here, we briefly describe the interactive web interface called the Map Tool 4. The RESTful API offers similar functionality.

Map Tool makes use of the Google Maps service to provide the basic geographical features such as topography, roads, streams, ponds, and so on. The Google Map service also allows us to perform complex visualization tasks to be described in Section 5. In this section, we provide an overview of the selection operations.

There are four different ways to select objects with Map Tool

- (1) Selecting a single well or object, by clicking on the well or object.
- (2) Selecting all wells or objects within a single rectangular bounding box, by drawing a bounding box on the map interface.
- (3) Selecting objects within a polygon defined by an arbitrary number of corner points within the web interface.
- (4) Selecting with more complex criteria through arbitrary Structured Query Language (SQL) style conditions entered through the filter tab on the map interface.

Figure 5 shows a rectangular bounding box selection using Map Tool. The different wells within the bounding box are marked by filled circles as opposed to the unselected wells marked by open circles. A list of analytes being measured by the selected wells is displayed adjacent to the map. Information about selected analytes can be displayed either individually or collectively (if more than

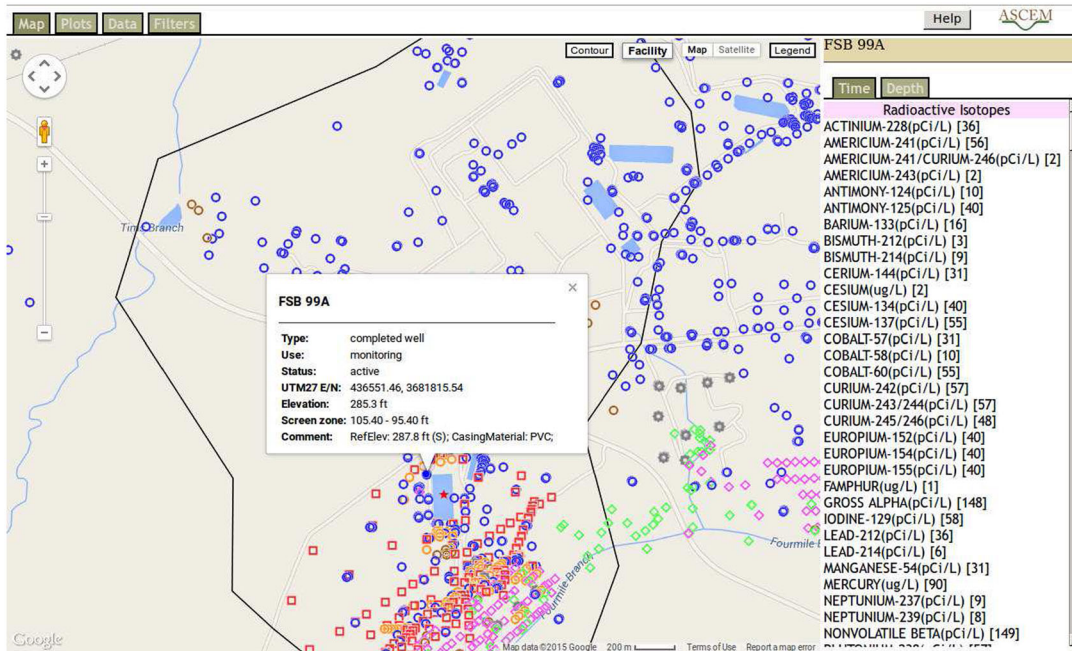


Figure 4. The ASCEM Data Management Map Tool.

one analyte is selected), and can be directly imported into a simulation model or saved to a file for off-line analysis.

5. PLOTTING AND VISUALIZATION

Map Tool is able to access a number of different visual analytics tools such as VisIt[§] and ROOT[¶]. It supports a range of plotting and rendering tools through the interactive web interface. Here, we briefly describe two examples: a simple line graph in Figure 6 and a set of advanced contours in Figure 7.

Figure 6 shows a simple line graph. There are a number of different forms of such plotting functions that produce scatter plots, bar graphs, and so on. This plotting capability can work with a single curve, or a collection of curves (e.g., measurements of all Uranium isotopes across a set of wells). This screen capture also shows the option of saving the data used in the plot to a text file. In general, this functionality is available for users to save data records for further analysis operations.

Figure 7 shows a set of advanced contours [23]. This advanced feature allows users to visualize concentration over time and space for any listed analyte for the selected wells; for a given aquifer and analyte, users can also overlay a contour plot of the analyte concentration as measured by all the wells screened in that particular aquifer for any calendar year in the measured range. This particular set of figures are produced with VisIt [23].

The Visualization Toolset of the ASCEM Platform has many more features that produce 2D and 3D images and movies on both primary and derived quantities. We refer interested readers to the software manual for a more detailed description.^{||}

[§] <https://wci.llnl.gov/codes/visit/>

[¶] <http://root.cern.ch>

^{||} <http://babe.lbl.gov/ascem/maps/SRDataBrowser.php>

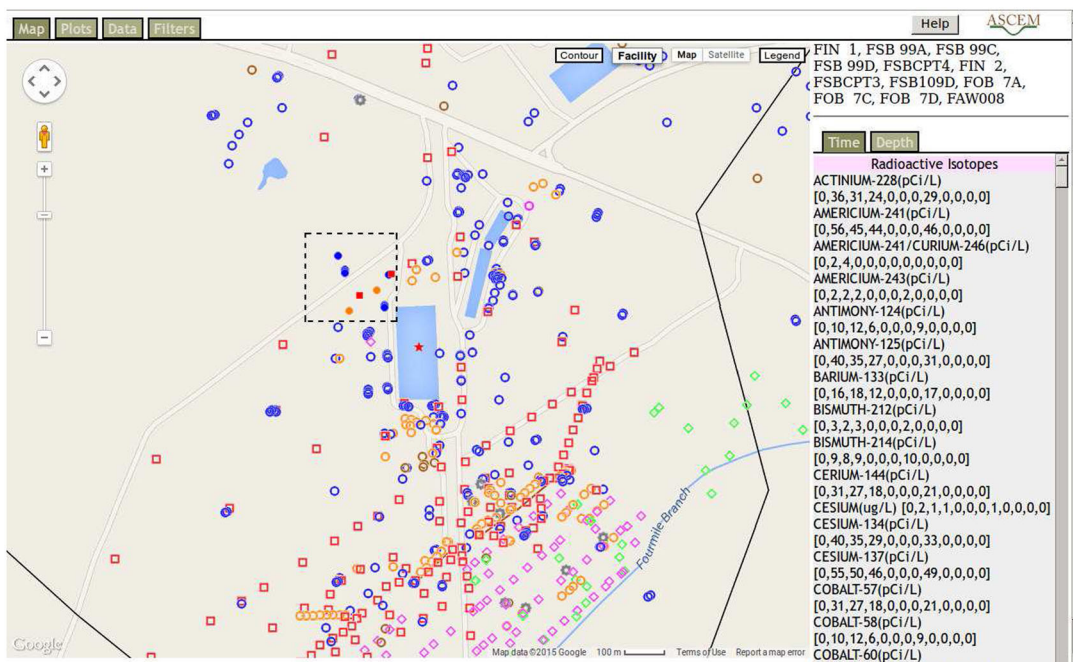


Figure 5. A sample selection of wells in a region using Map Tool.

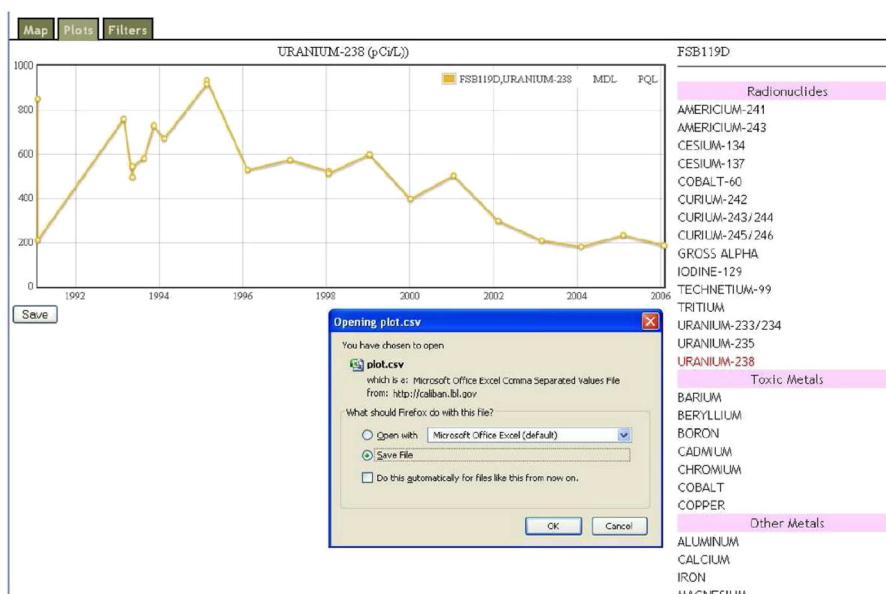


Figure 6. An example of a simple graph from Map Tool with the option to save the data for off-line analyses.

6. INTEGRATION WITH DISTRIBUTED WORKFLOWS

Several generic scientific workflow systems exist, for example, Kepler [24], Pegasus [25], Tigres [26], and Triana [27]. There are also a number of more specialized workflow systems designed for specific application domains [28]. For example, Turuncoglu and colleagues have developed a portable and replicable simulation workflows to create self-describing models with common model component interfaces for coupling an Earth System Modeling Framework with the Regional Ocean Modeling System and Weather Research and Forecasting Model [29, 30]. ASCEM has developed its own workflow system for subsurface modeling. The user interface to this workflow system is called

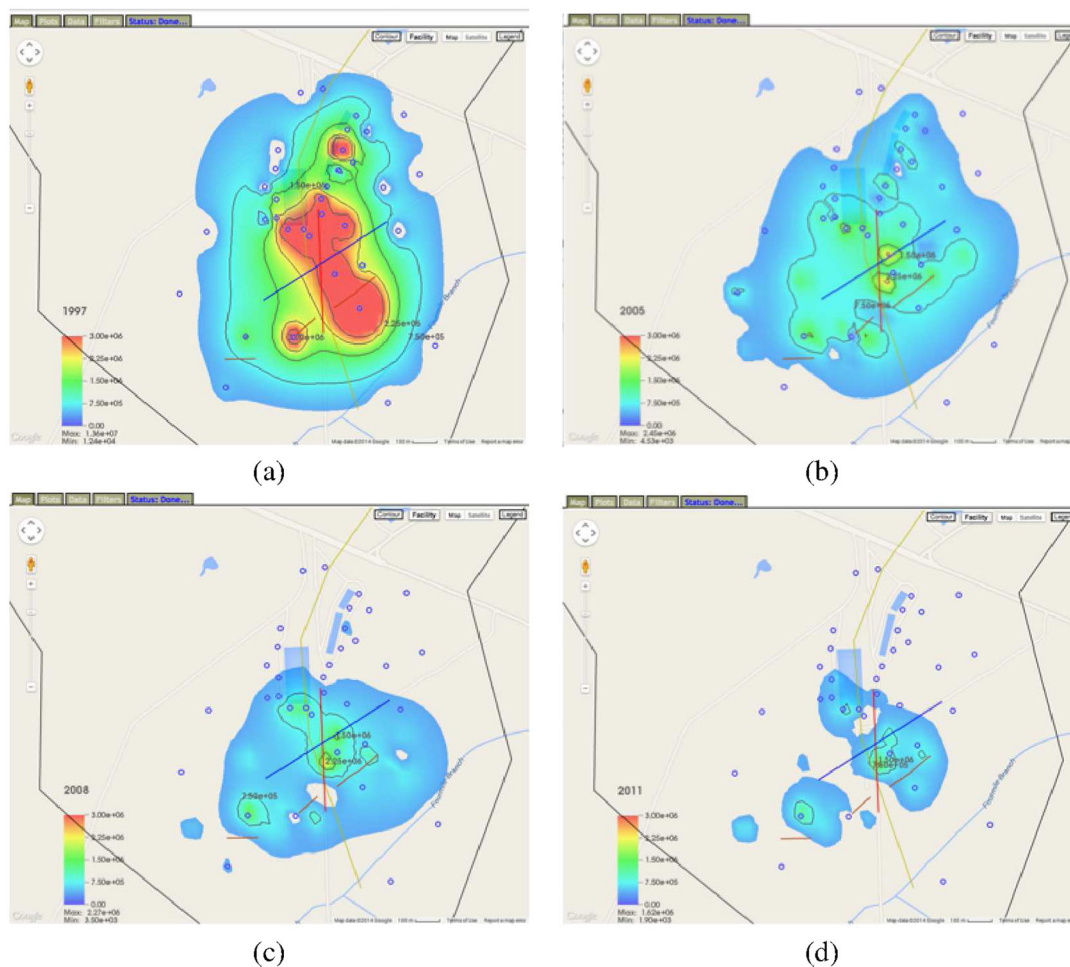


Figure 7. Web-based 2D plume visualization on the ASCEM data management website. A colored contour map of the tritium concentrations (in pCi/L) is created on top of the Google Map of the Savannah River Site F-Area: (a) 1997, (b) 2005, (c) 2008, and (d) 2011.

Akuna [13]. Here, we provide a brief overview of Akuna, and then describe the key feature of our data management tool that allows for a straightforward integration with the Akuna workflow engine.

Akuna is designed to connect effectively a diverse set of capabilities required for environmental management tasks including management of complex spatiotemporal data objects, visualization of many different quantities on geographical features, HPC simulation of transport and chemical reactions, and assessment of uncertainty in the simulation results. It is an open-source, platform-independent user environment. It includes features for basic model setup, sensitivity analysis, parameter estimation, uncertainty quantification, launching and monitoring simulations, and visualization of both model setup and simulation results. Features of the model setup tool include visualizing wells and lithologic contacts, generating surfaces or loading surfaces produced by other geologic modeling software and specifying material properties, initial and boundary conditions, and model output. Figure 8 shows the architecture of the current implementation [13].

Akuna workflow engine has implemented the capability of remote data access web services. This allows our data management component to seamlessly integrate with the Akuna workflow engine through the RESTful API. Because our web service follows the established standards, our approach allows the distributed Akuna workflows to access the necessary data whenever and wherever they are needed. The data management component is implemented with modern web technologies to provide intuitive interfaces for transparent data accesses. The selection capability allows the users to access only the data records needed for any analysis operation.

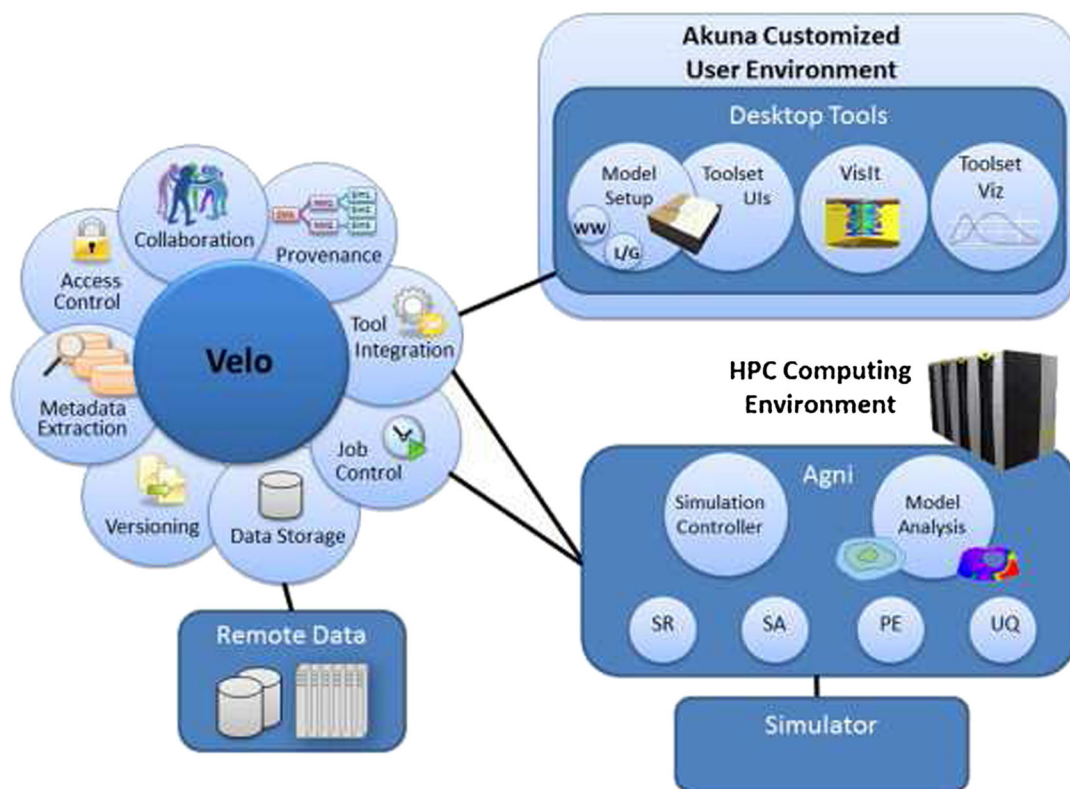


Figure 8. An architectural view of the Akuna workflow engine.

7. SUMMARY AND FUTURE WORK

The ASCEM project is developing computational tools to predict the long-term behavior of subsurface nuclear contamination. Part of this effort include collecting and disseminating important data about the experimental observations and simulations. ASCEM uses a specialized data gateway to satisfy the demand from the distributed community of scientists. The ASCEM data gateway is designed to serve the complex spatiotemporal data to a diverse set of applications. As such, the ASCEM data gateway includes a number of features that are uncommon in other implementations of science gateways. Here is a brief summary of the four key features explained in Sections 3–6.

- (1) The ASCEM data gateway supports complex geographical data types, including various coordinate systems for measuring objects on the globe, complex shapes on the globe, overlapping shapes, and irregular objects. It also supports time series of various kinds for capturing periodic observations and simulation output, such as the concentration values shown in Figure 6. These complex data types are supported through PostgreSQL.
- (2) It supports a set of advanced selection operations including points, rectangles, polygons, and SQL-style selection conditions. Through PostGIS, it also supports intersection of complex shapes. The data selected can be imported into the simulation model, or could be directed to another web service, saved to files, or used in plotting and other visual analytics operations. Common types of selection operations are explained in Section 4, and an illustration of a simple selection operation is shown in Figure 5.
- (3) It implements visualization functions by invoking the state-of-the-art features from VisIt and ROOT. In a number of special cases, we have worked with visualization experts to develop new algorithms to better handle the complex geometry involved [23].

- (4) It can be used by ASCEM workflows seamlessly to support large-scale data analysis and risk assessment operations. This is a natural outcome of using the web service standards in data consumers and producers within ASCEM project.

These key features of the ASCEM data gateway have now been implemented. The web service has been available online since 2011.** It has been used to produce a number of ASCEM publications [10–12, 23]. Our next set of tasks will be focused on exercising this gateway with a number of real-world risk assessment workflows. Based on the experiences from these use cases, we may expand this gateway to include additional features needed.

ACKNOWLEDGEMENTS

This work is supported in part by the Director, Office of Laboratory Policy and Infrastructure Management of the US Department of Energy under contract no. DE-AC02-05CH11231. This work used resources of The National Energy Research Scientific Computing Center (NERSC).

REFERENCES

- Allen B, Bresnahan J, Childers L, Foster I, Kandaswamy G, Kettimuthu R, Kordas J, Link M, Martin S, Pickett K, Tuecke S. Software as a service for data scientists. *Communications of the ACM* 2012; **55**(2): 81–88.
- Bernstein FC, Koetzle TF, Williams GJB, Meyer EF, Brice MD, Rodgers JR, Kennard O, Shimanouchi T, Tasumi M. The protein data bank: a computer-based archival file for macromolecular structures. *Archives of Biochemistry and Biophysics*. 1978; **185**(2): 584–591.
- Wruck W, Peuker M, Regenbrecht CRA. Data management strategies for multinational large-scale systems biology projects. *Briefings in Bioinformatics*. 2014; **15**(1): 65–78.
- Wainwright HM, Molins S, Davis JA, Arora B, Faybishenko B, Krishnan H, Hubbard S, Flach G, Denham M, Eddy-Dilek C, Moulton JD, Lipnikov K, Gable C, Miller T, Freshley M. Optimizing monitoring and remediation strategies at the Savannah River Site F–Area, using the Advanced Simulation Capability for Environmental Management (ASCEM). *Complex Soil Systems Conference*, Berkeley, USA, September 2014.
- Wang H, Huang JZ, Qu Y, Xie J. Web services: problems and future directions. *Web Semantics: Science, Services and Agents on the World Wide Web*. 2004; **1**(3): 309–320.
- Cholia S, Skinner D, Boverhof J. 2010. NEWT: a RESTful service for building high performance computing web applications. In *Gateway Computing Environments Workshop (GCE)*, 2010, IEEE; 1–11.
- Alameda J, Christie M, Fox G, Futrelle J, Gannon D, Hategan M, Kandaswamy G, von Laszewski G, Nacar MA, Pierce M, Roberts E, Severance C, Thomas M. The open grid computing environments collaboration: portlets and services for science gateways. *Concurrency and Computation: Practice and Experience* 2007; **19**(6): 921–942.
- Raicu I, Foster I, Szalay A, Turcu G. 2006. Astroportal: a science gateway for large-scale astronomy data analysis. In *Teragrid conference*; 12–15.
- Pierce EM, Freshley MD, Hubbard SS, Looney BB, Zachara JM, Liang L, Lesmes D, Chamberlain GH, Skubal KL, Adams V, et al.. Scientific opportunities to reduce risk in groundwater and soil remediation. *Technical Report PNNL-18516*, Pacific Northwest National Laboratory, Richland Washington, USA, 2009. (Available from: https://www.pnl.gov/main/publications/external/technical_reports/PNNL-18516.pdf.)
- Agarwal D, Wiedmer A, Faybishenko B, Hunt J, Kushner G, Romosan A, Shoshani A, Whiteside T. A methodology for management of heterogeneous site characterization and modeling data. In *The XIX International Conference on Computational Methods in Water Resources (CMWR 2012)*, Urbana-Champaign, IL, 2012.
- Seitz RR, Flach G, Hubbard S, Faybishenko B, Finsterle S, Steefel C, Dixon P, Moulton D, Freshley M, Freedman V, Gorton I, Marble J. 2013. Advanced simulation capability for environmental management-current status and phase II demonstration results-13161. In *Wm2013 conference*, 2013; 24–28.
- Williamson M, Meza J, Moulton D, Gorton I, Freshley M, Dixon P, Seitz R, Steefel C, Finsterle S, Hubbard S, Zhu M, Gerdes K, Patterson R, Collazo YT. Advanced simulation capability for environmental management (ASCEM): an overview of initial results. *Technology & Innovation* 2011; **13**(2): 175–199.
- Freedman VL, Chen X, Finsterle S, Freshley MD, Gorton I, Gosink LJ, Keating EH, Lansing CS, Moeglein WAM, Murray CJ, Pau GSH, Porter E, Purohit S, Rockhold M, Schuchardt KL, Sivaramakrishnan C, Vessilinov VV, Waichler SR. A high-performance workflow system for subsurface simulation. *Environmental Modelling & Software* 2014; **55**(0): 176–189.
- Arnette AN. Integrating rooftop solar into a multi-source energy planning optimization model. *Applied Energy* 2013; **111**: 456–467.

** <http://ascemdoe.org/>

15. Bastin L, Cornford D, Jones R, Heuvelink GBM, Pebesma E, Stasch C, Nativi S, Mazzetti P, Williams M. Managing uncertainty in integrated environmental modelling: the UncertWeb framework. *Environmental Modelling & Software* 2013; **39**: 116–134. Thematic Issue on the Future of Integrated Modeling Science and Technology.
16. Doherty J, Brebber L, Whyte P. PEST: Model-independent parameter estimation. *Watermark Computing*, Corinda, Australia, 1994; **122**.
17. Berners-Lee T, Cailliau R, Groff JF, Pollermann B. World-wide web: the information universe. *Internet Research* 2010; **20**(4): 461–471.
18. Obe R, Hsu L. *Postgis in Action*. Manning Publications Co.: Shelter Island, NY 2011.
19. Stones R, Matthew N. *Beginning Databases with PostgreSQL: from novice to professional*. Apress: New York, NY, 2005.
20. Worsley J, Drake JD. *Practical PostgreSQL*. "O'Reilly Media, Inc.": Boston, MA, 2002.
21. Pautasso C, Zimmermann O, Leymann F. RESTful web services vs. big web services: making the right architectural decision. In *Proceedings of the 17th International Conference on World Wide Web*, ACM: Beijing, China, 2008; 805–814.
22. Richardson L, Ruby S. *RESTful Web Services*. "O'Reilly Media, Inc.": Boston, MA, 2008.
23. Krishnan H, Meyer J, Romosan A, Childs H, Bethel EW. Enabling advanced environmental management via remote and distributed visual data exploration and analysis. *Computing and Visualization in Science* 2012; **15**(3): 123–133.
24. Altintas I, Berkley C, Jaeger E, Jones M, Ludascher B, Mock S. Kepler: an extensible system for design and execution of scientific workflows. In *Proceedings of the 16th International Conference on Scientific and Statistical Database Management*, IEEE: Santorini Island Greece, 2004; 423–424.
25. Deelman E, Singh G, Su MH, Blythe J, Gil Y, Kesselman C, Mehta G, Vahi K, Berriman GB, Good J, Laity A, Jacob JC, Katz DS. Pegasus: a framework for mapping complex scientific workflows onto distributed systems. *Scientific Programming* 2005; **13**(3): 219–237.
26. Ramakrishnan L, Poon S, Hendrix V, Gunter D, Pastorello GZ, Agarwal D. Experiences with user-centered design for the Tigres workflow API. In *2014 IEEE 10th International Conference on e-Science (e-Science)*, vol. 1, Guarujá SP, Brazil, October 2014; 290–297.
27. Churches D, Gombas G, Harrison A, Maassen J, Robinson C, Shields M, Taylor I, Wang I. Programming scientific and distributed workflow with Triana services. *Concurrency and Computation: Practice and Experience* 2006; **18**(10): 1021–1037.
28. Taylor IJ, Deelman E, Gannon DB, Shields M. *Workflows for e-Science: Scientific Workflows for Grids*. Springer Publishing Company, Incorporated: New York, NY 2014.
29. Turuncoglu UU, Murphy S, DeLuca C, Dalfes N. A scientific workflow environment for earth system related studies. *Computers & Geosciences* 2011; **37**(7): 943–952.
30. Turuncoglu UU, Dalfes N, Murphy S, DeLuca C. Toward self-describing and workflow integrated earth system models: a coupled atmosphere-ocean modeling system application. *Environmental Modelling & Software* 2013; **39**: 247–262. Thematic Issue on the Future of Integrated Modeling Science and Technology.