



## Earth and Space Science

### RESEARCH ARTICLE

10.1002/2015EA000139

#### Special Section:

Geoscience Papers of the Future

#### Key Points:

- Reproducibility of data-intensive analyses remains a significant challenge
- Data grids are useful for reproducibility of workflows requiring large, distributed data sets
- Data and computations should be co-located on servers to create executable Web-resources

#### Correspondence to:

J. L. Goodall,  
[goodall@virginia.edu](mailto:goodall@virginia.edu)

#### Citation:

Essawy, B. T., J. L. Goodall, H. Xu, A. Rajasekar, J. D. Myers, T. A. Kugler, M. M. Billah, M. C. Whitton, and R. W. Moore (2016), Server-side workflow execution using data grid technology for reproducible analyses of data-intensive hydrologic systems, *Earth and Space Science*, 3, 163–175, doi:10.1002/2015EA000139.

Received 7 SEP 2015

Accepted 8 MAR 2015

Accepted article online 15 MAR 2016

Published online 9 APR 2016

©2016. The Authors.

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

## Server-side workflow execution using data grid technology for reproducible analyses of data-intensive hydrologic systems

Bakinam T. Essawy<sup>1</sup>, Jonathan L. Goodall<sup>1</sup>, Hao Xu<sup>2</sup>, Arcot Rajasekar<sup>3</sup>, James D. Myers<sup>4</sup>, Tracy A. Kugler<sup>5</sup>, Mirza M. Billah<sup>6</sup>, Mary C. Whitton<sup>7</sup>, and Reagan W. Moore<sup>3</sup>

<sup>1</sup>Department of Civil and Environmental Engineering, University of Virginia, Charlottesville, Virginia, USA, <sup>2</sup>Data Intensive Cyber Environments Center, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA, <sup>3</sup>School of Information and Library Science, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA, <sup>4</sup>Inter-University Consortium for Political and Social Research, University of Michigan, Ann Arbor, Michigan, USA, <sup>5</sup>Minnesota Population Center, University of Minnesota, Twin Cities, Minneapolis, Minnesota, USA, <sup>6</sup>Department of Biological Systems Engineering, Virginia Polytechnic Institute and State University, Blacksburg, Virginia, USA, <sup>7</sup>Renaissance Computing Institute, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

**Abstract** Many geoscience disciplines utilize complex computational models for advancing understanding and sustainable management of Earth systems. Executing such models and their associated data preprocessing and postprocessing routines can be challenging for a number of reasons including (1) accessing and preprocessing the large volume and variety of data required by the model, (2) postprocessing large data collections generated by the model, and (3) orchestrating data processing tools, each with unique software dependencies, into workflows that can be easily reproduced and reused. To address these challenges, the work reported in this paper leverages the Workflow Structured Object functionality of the Integrated Rule-Oriented Data System and demonstrates how it can be used to access distributed data, encapsulate hydrologic data processing as workflows, and federate with other community-driven cyberinfrastructure systems. The approach is demonstrated for a study investigating the impact of drought on populations in the Carolinas region of the United States. The analysis leverages computational modeling along with data from the Terra Populus project and data management and publication services provided by the Sustainable Environment-Actionable Data project. The work is part of a larger effort under the DataNet Federation Consortium project that aims to demonstrate data and computational interoperability across cyberinfrastructure developed independently by scientific communities.

### 1. Introduction

There is an exponential growth in data available to geoscientists. The quantity of satellite data is growing rapidly [Acharya *et al.*, 1998], and data from sensor networks are being widely used, in observatories such as the Critical Zone Observatory [Anderson *et al.*, 2008], the National Ecological Observatory Network [Cowles *et al.*, 2010], and the Ocean Observing Initiative [Keller *et al.*, 2008]. Various groups are making available large collections of model-derived data including climate projections and reanalysis products for use by scientists. Public data repositories are used in many scientific disciplines as a means for sharing data collected by the so called “long tail” of the scientific community [Dunlap *et al.*, 2008]. The number of public repositories will likely increase as funding agencies enforce requirements that scientists submit data products resulting from their funded research to these public repositories.

This exponential growth in data will impact modeling and data analysis approaches used in many geoscience disciplines. As data sets grow in complexity and resolution, there is a need for improved tools to derive information from raw data sources in support of a particular research objective. These challenges arise not only because processing large, semantically unstructured data sets can be complex and time consuming but also because capturing the computational workflows scientists complete for a particular study can be challenging. New strategies are needed so that these scientist-authored computational workflows can make use of the latest available data and be reproduced and reused by other scientists.

One strategy for dealing with the growing volume of available data has focused on creating standards for accessing remote data collections using Web service Application Programming Interfaces (APIs). The

Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI) Hydrologic Information System has created standards for both an API called Water One Flow and a data exchange language called Water Markup Language to facilitate transmission of hydrologic time series data on large repositories using Web services [Maidment, 2008]. The Open Data Access Protocol (OpenDAP) is another widely used protocol for accessing and subsetting scientific data using Web services [Cornillon *et al.*, 2003]. OpenDAP focuses in particular on gridded data and includes the concept of server-side data subsetting and format conversion that are essential for operating on large, remote files.

While the Web service approach for data access has significant benefits, it also has limitations in that the network protocol for performing the data transfers using Web services operates over Hypertext Transfer Protocol. For large files, this approach is not optimal and potentially not feasible. Data grid technology provides an alternative approach for managing distributed data and computational resources. Data grids typically include features such as authentication, replication, authorization, auditing, and metadata support that are needed to manage large, distributed data collections [Foster, 2011; Rajasekar *et al.*, 2010]. These tools are better suitable for handling large files compared to Web services because they allow for parallel data transfers and provide automated fault tolerance and restarts when connectivity is lost during a transfer. Data grid technology has been used in the atmospheric and climate sciences, notably in the Earth System Grid and Earth System Grid Federation projects [Williams *et al.*, 2008, 2011], but it has not been widely adopted in other geosciences disciplines to date. In particular, research is needed to determine best practices and approaches for leveraging the technology to address specific needs in the hydrologic modeling community, which is the focus of this research.

The objective of this research is to explore approaches for leveraging data grid technology in hydrologic modeling to support reproducible workflows using large data sets. This is some of the first research applying data grid technology for hydrologic modeling. Its primary contribution is a general methodology for analyzing large, distributed data collections, by moving processing to data and using data grids to automate data transfers and staging, in combination with automated formal publication of generated data assets. This will be important as hydrologists seek to scale up watershed models to larger river basins where data sizes and computational processing make reproducibility more challenging.

The work is focused on a use case where a scientist wishes to create a workflow automating the data processing steps required to create a publication-ready figure from a large collection of model output files, greater than 2 Gb for a single run, produced using a Variable Infiltration Capacity (VIC) [Liang and Lettenmaier, 1994] hydrologic model. The use case, which is more fully explained in section 3, demonstrates server-side data processing on large data collections, using data grid technology for data transfers, and federation with public data repositories for reproducibility of the analysis workflow. It represents one of the first applications of the newly developed Workflow Structured Object (WSO) functionality in the Integrated Rule-Oriented Data System (iRODs), which has general applicability to other scientific domains with significant data management challenges. While systems like MyExperiment [De Roure *et al.*, 2009] also focus on server-side execution of scientist-authored workflows and provide advanced features for workflow sharing and publication, they focus on using Web services for data transfer rather than grid technology.

This research also addresses the challenge of federation across different cyberinfrastructure systems. It is likely that data-intensive studies will need to access many cyberinfrastructure systems for data gathering, processing, modeling, and publication. This paper demonstrates this concept for a use case that involves three cyberinfrastructure systems: the DataNet Federation Consortium (DFC) for data storage and compute resources, the Sustainable Environment-Actionable Data (SEAD) for data publication, and Terra Populus (TerraPop) for data access. Federation across these systems requires agreed upon standards and protocols that allow for interoperability. Different types of federation are demonstrated in our solution in order to address the transfer and management of both large and small data collections.

This paper is part of a special issue on the Geoscience Paper of the Future (GPF). GPF is envisioned as a paper where all digital assets used in the study are published as open, online resource published with unique identifiers and key metadata including titles, abstracts, licenses, authors, and contacts. In this paper, the key digital assets are published through SEAD with digital object identifiers (DOIs) and key metadata attributes. The research itself is also aimed at the vision and goals of GPF focusing in particular on the use case where computation is needed on distributed data resources. It seeks to define methods for moving data from

distributed servers within a data grid automatically using federation approaches and defining workflows that aid in capturing the provenance of how data were moved and processed to create publication-ready visualizations generated using multiple reference data collections. As data volumes continue to grow, such techniques will be critical to achieve the GPF goals.

The remainder of the paper is organized as follows. In section 2 we provide background on data grid technology to orient the reader. In section 3 we present the use case in further detail, followed by the design and implementation of a prototype system for solving the use case in section 4. Finally, we provide a discussion of key aspects of our approach in section 5 before offering concluding remarks in section 6.

## 2. Data Grid Technology

Data grids are systems that enable access and sharing of large data sets that are physically distributed across the Internet but appear to the user as a single file management system. The Integrated Rule-Oriented Data System (iRODS) is a data management system that includes the capability to federate data grids [Rajasekar *et al.*, 2010]. Federation allows for the creation of virtual data collections by logically arranging data from distributed resources under a virtual collection hierarchy. Globus is another data grid technology and is used within scientific communities and includes GridFTP for fast data transfer of large files [Foster, 2011]. While iRODS and Globus are commonly used within some specific scientific domains [Allcock *et al.*, 2002; Kyriazis *et al.*, 2008], their use is not widespread within the hydrology community.

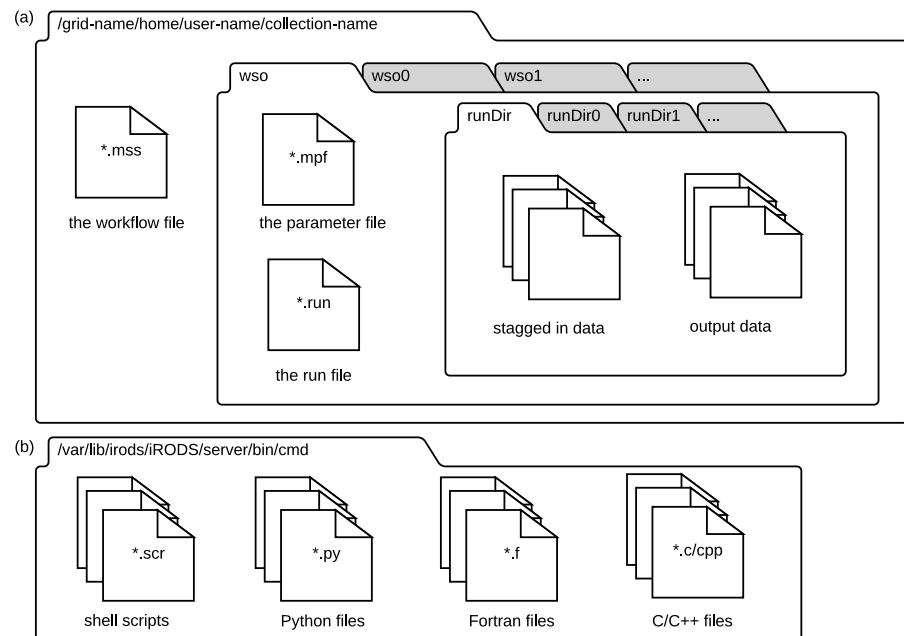
Data grids are particularly useful for scientific communities such as hydrology that rely on multiple data and computational resource providers. The iRODS-powered Data Federation Consortium (DFC) grid, which is used for this research, was developed as part of a National Science Foundation (NSF) funded project and provides support for federation of both resources and services. The work reported here is part of the DFC project and uses a DFC data grid for storage and long-term access to data sets stored across heterogeneous resources. The core iRODS software is developed and maintained by the iRODS Consortium at the Renaissance Computing Institute (RENCI), which is a partnership between the University of North Carolina at Chapel Hill (UNC-CH) and the Data Intensive Cyber Environments Center at UNC-CH. iRODS currently runs in Linux/Unix environments.

iRODS has a client-server architecture. The iRODS client software can be installed and run on any computer. Each iRODS grid installation has two types of servers: exactly one iRODS Metadata Catalog (iCAT) server and one or more iRODS resource servers, most frequently storage resource servers, e.g., data disks. Our system was developed on iRODS release 4.0, which includes software for the iRODS client, the resource server, and the iCAT server. iRODS uses the term zone as an abstraction for the physical components of an iRODS grid installation, i.e., the iCAT server and one or more resource servers that are part of the grid.

This work uses the recent development of iRODS Workflow Structured Objects (WSO), which enable workflows to be executed directly with iRODS commands. While iRODS is a mature, widely used software tool, this is some of the first work using the WSO functionality of iRODS. Therefore, this research was completed as a close collaboration between hydrologists defining the scientific workflows and the iRODS and WSO developers made possible through the DFC project. One goal of this work was to provide an example use case of applying WSO that could be beneficial for other iRODS users with interests in utilizing WSO in the future.

Figure 1a provides an overview of the file structure for a WSO. A WSO requires two primary files: a workflow file (\*.mss) and a parameter file (\*.mpf). The workflow file defines the sequence of operations to be performed by the workflow, and the parameter file lists the input arguments used when executing the WSO. The parameter file also specifies any files in iRODS that should be staged in (transferred to the physical directory on the iRODS resource server where the WSO is executed) or staged out (put into an iRODS collection) prior to and following the execution of the workflow [Rajasekar, 2014]. Examples of workflow and parameter files are provided in iRODS documentation, specifically from [https://wiki.irods.org/index.php/Workflow\\_Objects\\_\(WSO\)#Files\\_in\\_WSO](https://wiki.irods.org/index.php/Workflow_Objects_(WSO)#Files_in_WSO).

When the user creates and uploads a parameter file, iRODS automatically generates a run file (\*.run), which is then used by the client to execute the workflow. One workflow file can be used to create many instances of a WSO with each instance having a unique parameter file (see the `wso`, `wso0`, and `wso1` collections illustrated in Figure 1). The data files used by the workflow are stored in `runDir` collections. Within each WSO, there could be multiple `runDir` collections, one for each execution of the workflow. Workflows can include scripts and other scientist-authored code installed on the server in the `iRODS/server/bin/cmd` directory (Figure 1b).



**Figure 1.** (a) The structure of an iRODS Workflow Structured Object (WSO). (b) The WSO may utilize scripts installed in the iRODS/server bin/cmd directory for server-side data processing.

A WSO is executed by performing the following steps.

1. The user issues the *iput* command, which is part of the iRODS icommands client library, to transfer a workflow file (`*.mss`) from a client machine into an iRODS collection.
2. The user issues the *imkdir* command to make a new collection within the collection containing the workflow file (see the `wso` collection shown in Figure 1).
3. The user issues the *imcoll* command to mount this newly created collection.
4. The user issues the *iput* command to transfer a parameter file (`*.mpf`) into the mounted collection. This operation results in the system creating a run file (`*.run`) in the mounted collection.
5. The user issues the *iget* command on the run file to execute the workflow. The system then creates a new collection in the mounted directory (see the `runDir` collection shown in Figure 1), and the staged in and workflow-generated output files are stored in this new collection.

The same workflow can be executed for different parameter files by repeating steps 4 and 5 for a new parameter file, with each new parameter file resulting in an additional WSO collection (see `wso0`, `wso1`, ... shown in Figure 1) [Workflow Objects (WSO), 2013].

There are a number of workflow environments available to geoscientists, e.g., Kepler [Altintas et al., 2004], Taverna [Oinn et al., 2004], Triana [Harrison et al., 2008], and Pegasus [Deelman et al., 2005]. Like iRODS WSO, these workflow systems make trade-offs between power and flexibility. Many enable large-scale, parallel workflow execution on distributed resources, providing users real-time status information on the workflow execution [Vahi et al., 2013]. While workflow systems share many similarities, there are also key differences, which can often be subtle, that determine their suitability for addressing particular use cases. We used iRODS WSO in this analysis because our use case required a data processing pipeline consisting of a set of scientist-authored scripts that operate on data collections already within iRODS. Future work comparing and contrasting iRODS WSO with other workflow environments for completing this or other use cases relevant to hydrologic modeling would be a useful extension to this research [Introduction to Workflow as Objects, 2012].

### 3. Use Case Description

The prototype software described in this paper is designed to address a use case where a scientist has created a simulation using the Variable Infiltration Capacity (VIC) model for the Carolinas region of the United States. The model has been calibrated and validated for this region as part of a prior study

[Billah *et al.*, 2015] and can be used to address other hydrologic research questions as well. The scientist that created the model has published the model's input and output files on the Web for use by other scientists. A second scientist learns about the model and wishes to use the model's output files to test her own research question about drought impacts on counties within a study region. The scientist is interested in how soil moisture deficit predicted by the model varied for different populated communities within the study region. While this application is analyzing historical events, it would be relatively straightforward to set up the calibrated model to analyze current conditions and to identify populated regions vulnerable to drought conditions within the region. Such information would be valuable to resource managers in better understanding the severity of the drought and its impact on population centers within the region.

The second scientist downloads the model output files published online by the first scientist and creates the visualization by writing her own Python scripts. The scientist downloads the population data for the study counties to a local working directory. The VIC soil moisture outputs are organized in a set of "flux files," one for each node in the modeling domain. The Python scripts sort through these data extracting relevant information and summarizing the soil moisture time series. Geospatial processing tools are used to relate the coordinates of the model nodes to counties in the study region. The result of this data processing is a comma-separated values (CSV) file with the soil moisture deficit and population for each of the five counties. Finally, the scientist programs the Python script to use this CSV file to produce a publication-ready figure for visualizing the drought impacts.

In addition to publishing the scripts and data files from this analysis on a public data repository, which is now a relatively straightforward exercise given the proliferation of online data repositories, the scientist also wishes to publish the workflow used to perform the analysis as a Web executable resource. The scientist wishes to take this approach for the following reasons.

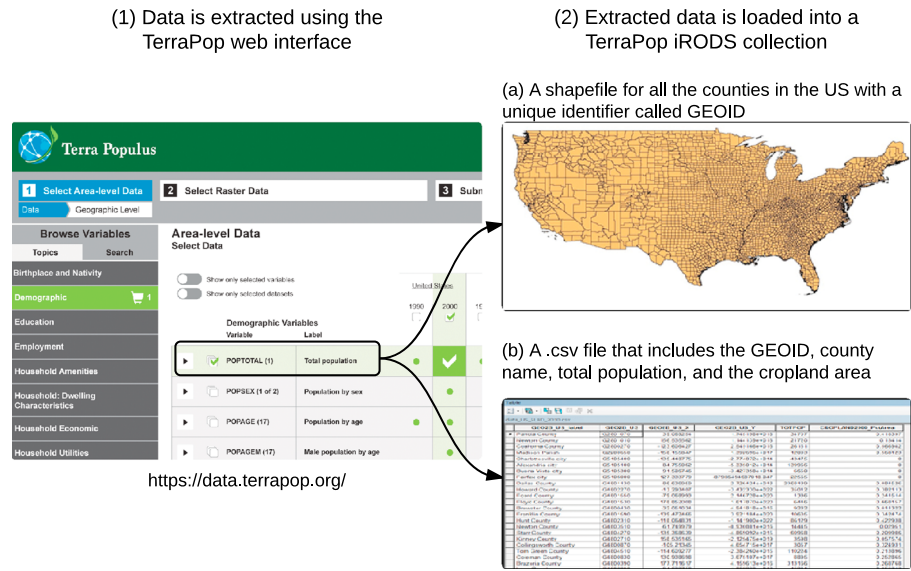
1. Having the overall workflow be executable server-side means the scripts and model output data can be colocated, removing the need to download the large model output file to the scientist's machine prior to the workflow execution.
2. By keeping data sets server-side, it is easier to ensure that the data have not been modified after making a local copy (its provenance can be proven). With the ability to publish the model and reference data once and to keep them on the server, only the visualization results need to be retrieved and published for subsequent runs.
3. Having server-side execution of the workflow controls for potential variability across different hardware and software configurations on a client machine. Even with this relatively simple use case of creating a figure, there is potential for different operating systems and versions of analysis software to result in differences in the end product. These software dependencies could result in additional time for scientists to trouble shoot errors. More critically, these dependencies could result in an end product without errors or warnings but with inconsistencies due to nonbreaking differences between dependent software versions.

Simply put, having data and processing colocated on a server as a Web executable resource results in a more controlled environment, which is critical for reproducibility.

The scientist uses iRODS WSO to create the Web executable resource. As part of the WSO, the scientist defines the steps to automatically stage in the required VIC output and population data that are stored in iRODS collections. The population data come from TerraPop, which provides global-scale data sets that focus on human population characteristics, land use, land cover, and climate change [Minnesota Population Center, 2013]. The Terra Populus data access system was used to create customized data extracts, combining variables from multiple sources into a single package. Users can browse the TerraPop collection and select the required variables; the variable required in this paper was the total population for each county in the United States. After submitting our data request, the system generated a data package that included a shapefile for all the counties in the United States, with unique codes that identify the polygon defining each county (GEOIDs), and a CSV file that includes the GEOID, name, and total population of each county (Figure 2). This data package was then automatically uploaded onto the TerraPop grid as an iRODS collection. By federating the DFC hydrology and TerraPop zones and configuring authorizations, we are able to have the population data remain on the TerraPop server and be automatically staged in for use by the WSO.

Finally, the data (including code) resulting from the analysis are published using products provided by the Sustainable Environment-Actionable Data (SEAD) project [Myers *et al.*, 2015]. The SEAD project supports publication, preservation, and sharing of data generated by scientists including data generated by running models. Using SEAD, teams of researchers can upload, share, annotate, and review input data sets and





**Figure 2.** Details on how the county-level population data are requested and extracted using the TerraPop web interface into an iRODS data collection. From this collection, iRODS stages in the required files prior to the workflow execution.

model outputs within an access-controlled Project Space and then formally publish collections of data with associated metadata and provenance for long-term preservation (generating a digital object identifier (DOI) and standards-based archival package and registering the data with the DataONE catalog for discovery). Our use of SEAD included manual entry of data and metadata via a web interface and bulk uploads of files and programmatic submission of the output figure with metadata to SEAD, which leveraged SEAD’s RESTful Web API.

#### 4. Prototype Software Design and Implementation

We present the prototype software aimed at addressing the use case by first describing the steps taken to configure the server-side software and data, next describing the steps required to configure the WSO, then describing the steps required to execute the WSO from the client machine and concluding with a summary of the results from executing the workflow.

##### 4.1. Server-Side Configuration

To perform the server-side configuration, we first installed iRODS resource server version 4.0 software on an Elastic Cloud Computing (EC2) instance in the Amazon Web Services (AWS) cloud. We chose AWS because it provides on-demand computing resources and services that can be easily scaled to meet demands. The EC2 service provided through AWS allows users to rent virtual machines (instances) with different capabilities and pay by the CPU hour. For prototyping purposes, we used a Linux-based medium-sized machine (m3) with 3.75 Gb of memory, 4 vCPU, 15 Gb of Solid State Drive (SSD)-based local instance storage, and 64 bit platform for the iRODS resource server [Amazon EC2 Instances, 2015]. Next, this new iRODS resource server was configured to be part of the DFC hydrology zone that has its iRODS Metadata Catalog (iCAT) server on a machine running at RENCI. We had to configure the AWS EC2 instance to be associated with an elastic Internet Protocol (IP) address to avoid having to update the EC2 instance’s IP addresses in the iCAT server following each restart of the EC2 instance.

We then developed a WSO on the iRODS resource server to implement the data visualization workflow described in the use case. This required that the user have an account on the server itself with read/write access to the cmd directory (Figure 1b). It was also necessary to set read/execute rights on the files associated with the WSO so that they could be executed by the iRODS user account. We uploaded to the iRODS resource server the VIC model output files from SEAD (where the original scientist had published them for use by the community), the Python scripts created by the scientist to generate the visualization, and the shell script, also created by the scientist, used to sequence the execution of the Python scripts on the iRODS resource server.

The VIC source code is not included in SEAD because the source code is available from the developer's GitHub page instead (see <https://github.com/UW-Hydro/VIC>).

#### 4.2. Client-Side Configuration

The client machine can be any computer with the iRODS client software installed. In this prototyping work, we used a second EC2 instance as the client machine simply to avoid moving data into and out of the AWS cloud. We installed the icommands iRODS client software library on the client machine. The icommands software includes a set of commands that perform operations such as make a new directory (`imkdir`) or put a file into an iRODS collection (`iput`) [Weise *et al.*, 2008]. The icommands client library includes an environment configuration file that is used to point to a particular iRODS zone and set default user credentials for accessing the iRODS zone. In our case, we configured the icommands environment to operate on the DFC hydrology zone and entered user credentials representing the scientist accessing the system.

The general file structure required for creating a WSO was described in section 2 and in Figure 1a. For our particular application, we first created a workflow file (`PopVsSm.mss`) that specifies the steps required to execute the workflow. The workflow file simply specified that the workflow should execute the scientist-authored shell script installed on the iRODS server `cmd` directory. We put the `PopVsSm.mss` file into an iRODS collection and then made a new collection named "vic\_soilmoisture." We mounted this new collection, effectively making it a WSO.

#### 4.3. Executing the Workflow

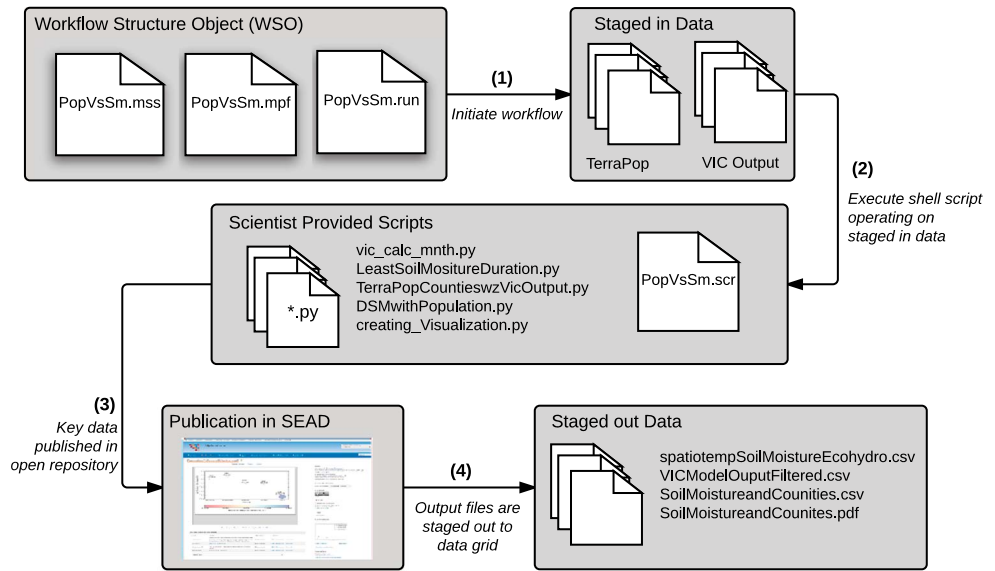
Once the WSO is mounted, it is then possible to execute the workflow. This process is described, in general, in section 2. Here we provide specifics of the WSO execution for the use case. The general flow of data and sequence of commands for executing the WSO execution for the use case are described in Figure 3.

1. The user initiates execution of the workflow by issuing an `iget` command on the `PopVsSm.run` file that is in the mounted WSO collection. The `PopVsSm.mpf` parameter file defines the data required by the workflow and stages these files from different iRODS collections into the directory on the iRODS resource server where the WSO is executed. In our case, we staged in the VIC model output data stored in the DFC hydrology grid and county-level population data from the TerraPop grid. While these two data sets are stored within different grids, it is possible to gain access to the data directly using iRODS authentication because the grids are federated.
2. Once all required data are staged into the iRODS resource server directory where the workflow is executed, the workflow file specifies that the scientist-authored shell script stored on the iRODS server should be executed. This shell script then calls a series of scientist-authored Python scripts that process the staged in data to create the output figure.
3. A final step in the shell script is publishing the figure resulting from the workflow automatically to a SEAD project space for sharing with colleagues and subsequent publication. The SEAD API is used for this purpose and allows for the submission of the file along with associated metadata to a SEAD project space.
4. Upon completion of the workflow, key output data are staged out into iRODS collections according to specifications in the parameter file. This allows the files to be accessible to authorized users in the grid.

Figure 4 shows the steps for executing a WSO from a user's perspective when working with the icommands client library. The user must know which iRODS collection contains the script files required for executing the WSO to be able to execute it. Once the user has logged into the client machine, the user changes the working directory to the iRODS logical path where the WSO has been mounted. In this case, the WSO was mounted as the "vic\_soilmoisture" collection. The user next issues an `iput` command to put the parameter file (`PopVsSm.mpf`) into the mounted WSO. This step is not illustrated in Figure 3 for brevity but results in the generation of a run file (`popvssm.run`) in the collection. Finally, the client executes the workflow by issuing an `iget` command on the `popvssm.run` file.

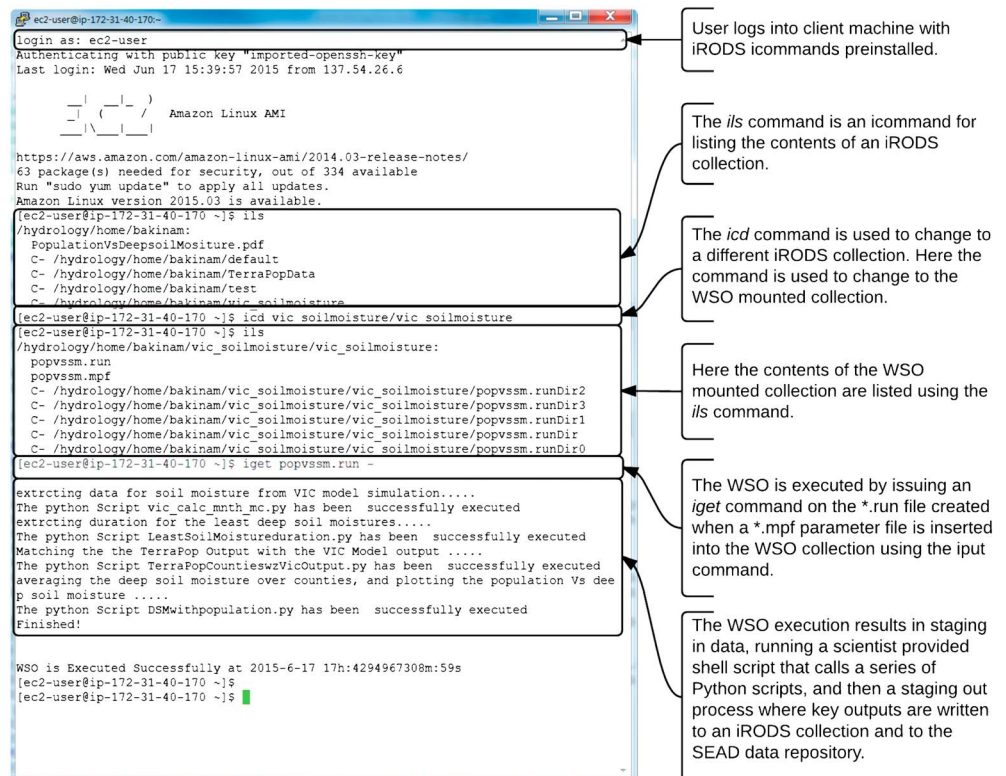
#### 4.4. Results From the Workflow Execution

When the workflow is executed, the output messages are written to the console, even though computation is performed on the server-side and no data (other than the output messages) are transferred to the client machine. Once the workflow execution has completed, the user can access the output collection called `runDir` resulting from the workflow execution. The `runDir` file contains by default the `stdout` from the execution of the workflow along with any staged in and derived data from the workflow [Workflow Objects (WSO), 2013].



**Figure 3.** The steps that occur on the server side when a user executes the WSO. Data are staged in from iRODS collections, scientist-authored scripts are run to create the figure, data are published through a SEAD project space using the SEAD API, and key output data are staged out back into iRODS collections.

The workflow also results in publication of the workflow results to a SEAD project space. Figure 5 shows the data collections as they appear through the SEAD project space website. Most data were uploaded using the SEAD web interface. Figure 6 shows the figure resulting from the WSO execution that was automatically written to the SEAD project space using the SEAD API as a final step in the WSO execution.



**Figure 4.** The steps required from a client machine in order to execute the WSO using the icommands client library.



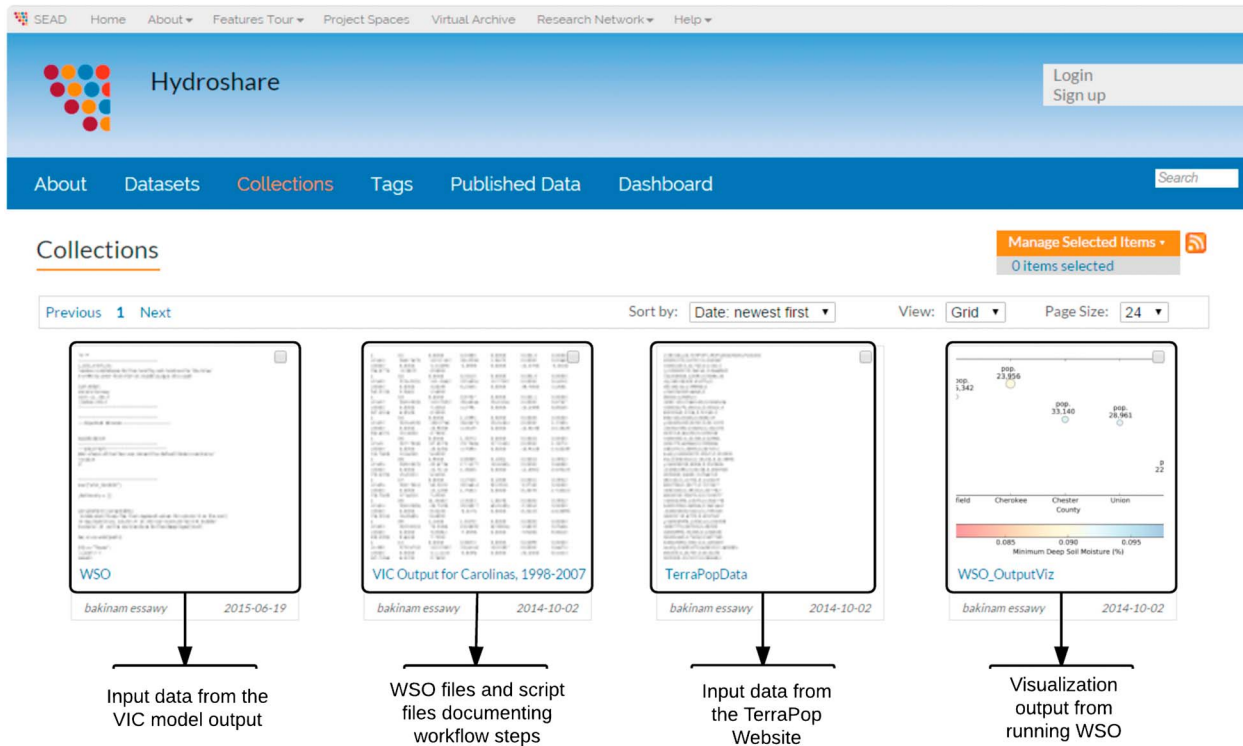


Figure 5. Contents of the Sustainable Environment Actionable Data (SEAD) project space used for storing and accessing data used in the workflow.

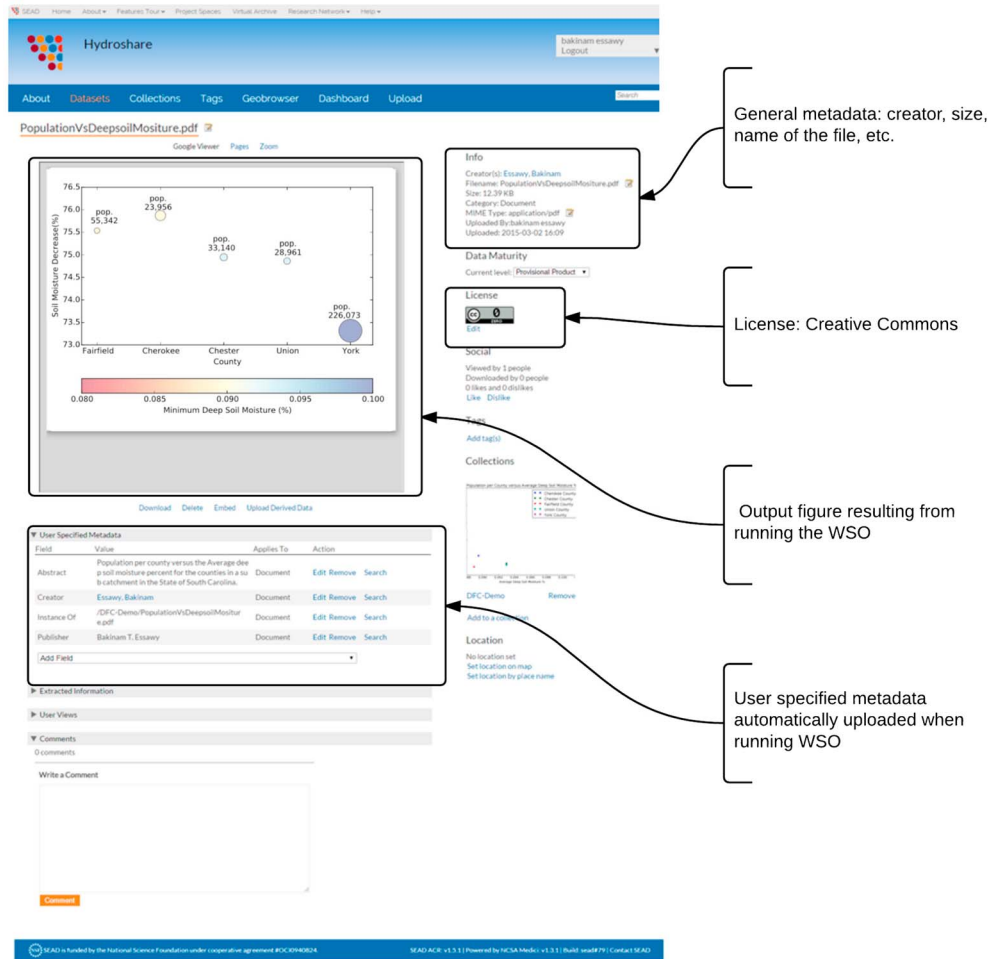
## 5. Discussion

### 5.1. Reproducibility

To support transparency and reproducibility of this work as envisioned by the Geoscience Paper of the Future (GPF) project, the data collections in the use case (e.g., the VIC output files, the TerraPop data, the WSO files, and the output figure) were published in SEAD. As part of this publication process, each collection was given metadata including a brief abstract, creators, and the publisher and then published to generate a digital object identifier (DOI) (Table 1). The output figure resulting from the WSO execution was first written to a SEAD project space along with basic metadata as a final step in the WSO execution using the SEAD API. From there, the scientist logged into the SEAD web interface and set additional metadata fields to publish the resource with an assigned DOI. Any combination of automated and manual entry is supported, and researchers can choose which data to publish. In our case, we automatically captured outputs from multiple test runs before manually selecting, annotating, and publishing (including creating a DOI for) only the final run.

Use of an open, metadata-aware repository makes it simple to capture additional derived data and provenance information as research continues. By publishing the reference data, scripts, and output data separately in SEAD, we also demonstrate the ability for larger reference data to be published once and then referenced via provenance links from the derived output files that could be generated by many researchers over time. For example, the VIC output files used in this workflow may be used in other research studies. If each publication using these VIC output files references its DOI, it will be possible to track the impact of the model output files through citation counts similar to what is done now for tracking citation counts of research papers.

Other end points could be used for publishing key digital assets from the WSO workflows. For example, the Consortium of Universities for the Advancement of Hydrologic Sciences, Inc. (CUAHSI) HydroShare system is in development and could serve as an alternative or secondary end point for publishing results with more discipline-specific metadata [Horsburgh et al., 2015; Morsy et al., 2014; Tarboton et al., 2014], as could systems such as FigShare or Zenodo. We anticipate a growing number of such repositories and for federation between them (e.g., SEAD is already a member node in DataOne [Michener et al., 2012], advertising our WSO publications through DataONE’s catalog). This research shows how iRODS WSO could play an important



**Figure 6.** View of figure, produced by executing the WSO, within the SEAD project space. The workflow uses the SEAD API to upload this resource along with metadata to the SEAD project space.

role in moving data resources within such data repositories to and from computational resources to support data computation use cases.

Using a public cloud offers further opportunities for reproducibility. It is possible to quickly set up virtual machines (VMs) with a variety of operating systems to reproduce computational analyses. It is also possible to capture images of VM instances that can be stored for future reproducibility. Exploring the use of virtual containers (e.g., the Docker project) rather than VMs would be a useful extension to this work. Virtual containers can reduce setup time and storage costs compared to VMs for software, like what was used in this work, which run in a Linux operating system.

### 5.2. Federation

Federation across cyberinfrastructure systems is a key aspect of this work. Federation describes how distinct and formally disconnected systems interoperate. There is a growing set of cyberinfrastructure systems available to scientists, and many studies will benefit from the use of more than one of these systems. Effective ways for federating across these systems will result in powerful tools that save scientists' time and encourage reproducibility through automatic data transfers handled directly by systems. This concept was illustrated in our study by showing how distinct cyberinfrastructure systems can be federated and used collectively within a single workflow execution.

Figure 7 provides a depiction of the workflow that emphasizes different data collections and approaches for federating between DFC, TerraPop, and SEAD. The use case in this study represents two levels of federation that we believe are relevant for most scientific studies. The federation between the AWS machine where the

**Table 1.** Key Digital Assets Used in the Study That Are Published Through SEAD With Basic Metadata

Title	DOI	Author	Contact	Abstract	License
TerraPopData Extract	10.5967/M08P5XH5	Essawy, Bakinam	Goodall, Jonathan	Population data extracted from TerraPop ( <a href="https://data.terrapop.org">https://data.terrapop.org</a> ) for the study region	Creative Commons (CC)
VIC Output for Carolina, 1998–2007	10.5967/M0DF6P6F	Essawy, Bakinam	Goodall, Jonathan	Output from a VIC model for the Carolinas, USA, calibrated for the period 1998–2007 to study drought impacts	Creative Commons (CC)
WSO	10.5967/M0J67DXR	Essawy, Bakinam	Goodall, Jonathan	The scripts and related files used to create the iRODS Workflow Structured Object (WSO)	Creative Commons (CC)
WSO_OutputViz	10.5967/M0513W51	Essawy, Bakinam	Goodall, Jonathan	Impact of 2007 drought on five counties in the study region	Creative Commons (CC)

workflow was executed and the TerraPop reference data is what we term a strong federation, while the federation between the AWS machine and SEAD is what we term a weak federation. A strong federation is based on a strong trust model where one data grid administrator can add credentials of users of other data grid and grant access to resources based on authentication through other data grids. One primary benefit of this level of federation is that data grid technology can be used to transfer files between the two systems. For large files, this level of federation will be important because of the functionality provided by data grids like iRODS that are designed specifically to ensure rapid and successful transfer of large files over a network. Weak federation, based on federation through Web service APIs, allows for greater flexibility and less required trust between systems, because all operations are through services. Transferring large data through Web services, however, is not ideal for the reasons we outlined in section 2.

### 5.2. Adoption

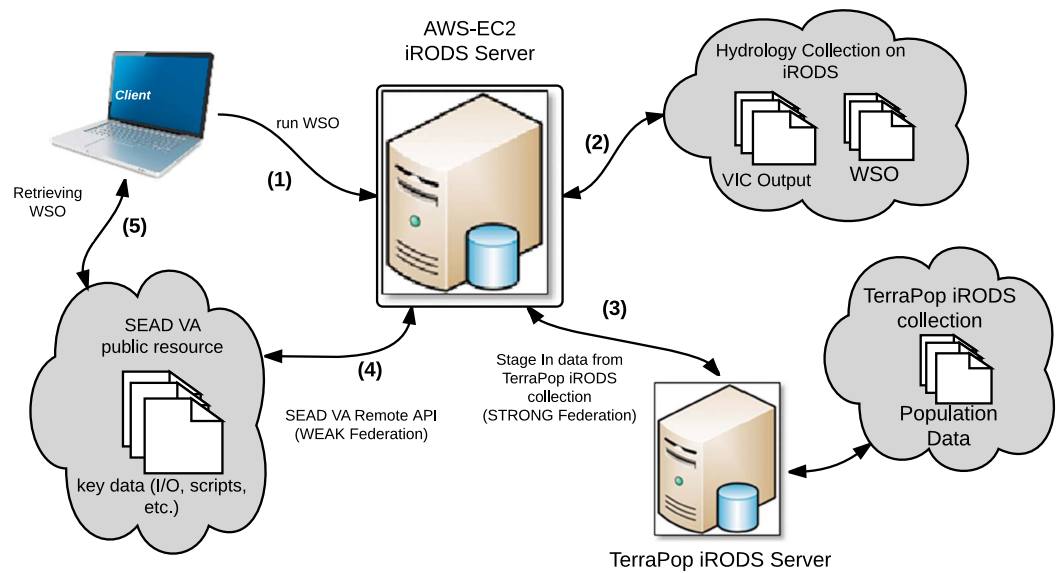
While there are many advantages to the approach described in this paper, there are also important barriers to adoption, especially in terms of the current prototype system. Currently, users of the system need to be familiar with an iRODS client (e.g., the icommands client library used in this study). They must also be aware of steps for executing a WSO. Developers need an understanding of how to structure new WSOs and will need access to the server running the iRODS resource server software for installation and configuration of the WSO.

There are opportunities for abstracting the complexity of directly interfacing with iRODS WSO for end users in order to encourage broader adoption of the technology. One way to do this would be to have someone familiar with iRODS WSO take input from the scientist including the scripts needed to execute the workflow and the location (iRODS logical path name) of the input data for the scripts. The administrator would then mount a WSO with an example parameter file and make it available through the system to end users. The user could then execute the workflow either using the icommands client library, as described in the paper, or through other tailored client applications able to operate on iRODS collections including executing WSOs stored within iRODS collections. We believe this would be a fairly straightforward process for moving scientist-authored codes into a form that is Web executable.

### 5.3. Data Size and Heterogeneity Challenges

This work only begins to illustrate the potential benefit of using data grid technology for executing workflows that require heterogeneous data from distributed data sources. We showed how WSOs allow for automatically staging in of required data distributed across a data grid. We also showed how data produced from the workflow can be staged out, meaning written to collections in the data grid where it can be accessible to other users. While it was not demonstrated in this use case, one can execute a distributed workflow across the network on multiple iRODS resource server using WSO.

This approach allows the location of the input and output files for a computational tool to be independent of the location where the processing is done. However, unlike approaches that rely only on Web service APIs for data staging prior to workflow execution, iRODS provides a more robust data staging approach that leverages grid technology. While the use case demonstrated the concept using fairly small file sizes, the solution we used can be applied to larger terabyte scale data as well. Given that modeling in many geoscience



**Figure 7.** Main components and data flow in the workflow emphasizing data collections and federation approaches

disciplines requires access to large, distributed data, data grid technology provides a powerful way for data staging associated with workflow execution.

### 6. Conclusions

The focus of this paper is on creating scientist-authored workflows as Web executable resources in data grids. The iRODS WSO provides researchers with the ability to publish their research methods for computational studies as workflows that specify the tools, data, and sequence of steps taken to complete the study. All of these digital objects (data, software, model outputs, etc.) can be made accessible to other users of the data grid as well as to nongrid users through publication in SEAD.

There are many challenges in reaching the ultimate goal of reproducibility, especially when dealing with data-intensive modeling analyses that require a large, diverse set of input data and generate a large, diverse set of output data. Through this work, we argue that reproducibility will require more data processing server-side, i.e., where reference data and models are managed together, than what is common now. This is due to the large and increasing size of data sets used by geoscientists and the growing complexity of software and software dependencies that require constrained environments to ensure reproducibility.

We also argue for multiple federation approaches as means for providing interoperability across the variety of cyberinfrastructure systems needed for data access, analysis, modeling, and publication services. Federation approaches most often used in geoscience disciplines emphasize Web service APIs; however, to support large data sets, the community should have broader adoption of data grid federation approaches as well. The use of both approaches was demonstrated for a use case that leveraged four federated but heterogeneous cyberinfrastructure systems: DFC, TerraPop, SEAD and via an existing connection with SEAD and DataONE.

Any approach for making scientific computations into Web executable resources must have a low barrier to entry for users. We have proposed an approach that allows scientists to write scripts as is typically done now for data analysis using languages familiar to scientists. These scripts can be then made available as Web executable resources to scientists using iRODS WSO technology. Future work should explore embedding of iRODS WSOs into systems that include tailored interfaces for scientific communities. Then, rather than performing the steps described in the paper for executing WSO that include the use of the icommand client library, the end user could have a simpler and more tailored interface for viewing and executing workflows that abstracts technical details from the end user.

There are encouraging trends toward increased publication of data (including code) used in scientific studies. It is important that the momentum behind these trends result in scripts and workflows as Web executable resources to capture their full potential in advancing reproducibility goals. The advantages of Web executable resources include the increased ability to share, reproduce, and collaborate on scientists-authored workflows. While the

potential of scientific scripts and workflows as Web executable resources is clear, important issues remain related to managing large data and computation collections. We have demonstrated here an approach using data grids for addressing this challenge and have argued for moving processing to reference data stored within data grids as a method for creating reproducible scientific workflows on large data sets.

#### Acknowledgments

This work was supported by the National Science Foundation (NSF) under awards ACI-0940841, ACI-0940824, and ACI-0940818 and by Amazon Web Services (AWS) through an Education Research Grant award. This research would not have been possible without assistance from the larger iRODS, DFC, SEAD, and TerraPop teams. The data used are listed in Table 1 and can be found in the SEAD repository at the DOIs provided in Table 1.

#### References

- Acharya, A., M. Uysal, and J. Saltz (1998), Active disks: Programming model, algorithms and evaluation, *ACM SIGPLAN Not.*, *33*, 81–91, doi:10.1145/291006.291026.
- Allcock, B., J. Bester, J. Bresnahan, A. L. Chervenak, I. Foster, C. Kesselman, S. Meder, V. Nefedova, D. Quesnel, and S. Tuecke (2002), Data management and transfer in high-performance computational grid environments, *Parallel Comput.*, *28*, 749–771, doi:10.1016/S0167-8191(02)00094-7.
- Altintas, I., C. Berkley, E. Jaeger, M. Jones, B. Ludäscher, and S. Mock (2004), Kepler: An extensible system for design and execution of scientific workflows, 16th International Conference On. IEEE, 2004. pp. 423–424. doi:10.1109/SSDM.2004.1311241.
- Amazon EC2 Instances [WWW Document] (2015), [Available at <http://aws.amazon.com/ec2/instance-types/>, accessed 6.7.15.]
- Anderson, S. P., R. C. Bales, and C. J. Duffy (2008), Critical zone observatories: Building a network to advance interdisciplinary study of Earth surface processes, *Mineral. Mag.*, *72*, 7–10, doi:10.1180/minmag.2008.072.1.7.
- Billah, M. M., J. L. Goodall, U. Narayan, J. T. Reager, V. Lakshmi, and J. S. Famiglietti (2015), A methodology for evaluating evapotranspiration estimates at the watershed-scale using GRACE, *J. Hydrol.*, *523*, 574–586, doi:10.1016/j.jhydrol.2015.01.066.
- Cornillon, P., J. Gallagher, and T. Sgouros (2003), OPeNDAP: Accessing data in a distributed, heterogeneous environment, *Data Sci. J.*, *2*, 164–174, doi:10.2481/dsj.2.164.
- Cowles, T., J. Delaney, J. Orcutt, and R. Weller (2010), The Ocean Observatories Initiative: Sustained ocean observing across a range of spatial scales, *Mar. Technol. Soc. J.*, *44*(6), 54–64, doi:10.4031/MTSJ.44.6.21.
- De Roure, D., C. Goble, and R. Stevens (2009), The design and realisation of the virtual research environment for social sharing of workflows, *Futur. Gener. Comput. Syst.*, *25*, 561–567, doi:10.1016/j.future.2008.06.010.
- Deelman, E., G. Singh, M. Su, J. Blythe, Y. Gil, and C. Kesselman (2005), Pegasus: A framework for mapping complex scientific workflows onto distributed systems, *Sci. Program.*, *13*, 219–237.
- Dunlap, R., L. Mark, S. Rugaber, V. Balaji, J. Chastang, L. Cinquini, C. DeLuca, D. Middleton, and S. Murphy (2008), Earth system curator: Metadata infrastructure for climate modeling, *Earth Sci. Inf.*, *1*, 131–149, doi:10.1007/s12145-008-0016-1.
- Foster, I. (2011), Globus online: Accelerating and democratizing science through cloud-based services, *IEEE Comput. Soc.*, *15*, 70–73, doi:10.1109/MIC.2011.64.
- Harrison, A., et al (2008), WS-RF workflow in Triana, *Int. J. High Perform. Comput. Appl.*, *22*, 268–283, doi:10.1177/1094342007086226.
- Horsburgh, J. S., M. M. Morsy, A. M. Castronova, J. L. Goodall, T. Gan, H. Yi, M. J. Stealey, and D. G. Tarboton (2015), Hydroshare: Sharing diverse environmental data types and models as social objects with application to the hydrology domain, *J. Am. Water Resour. Assoc.*, doi:10.1111/1752-1688.12363.
- Introduction to Workflow as Objects [WWW Document] (2012), [Available at [https://wiki.irods.org/index.php/Introduction\\_to\\_Workflow\\_as\\_Objects](https://wiki.irods.org/index.php/Introduction_to_Workflow_as_Objects), accessed 6.7.2015.]
- Keller, M., D. S. Schimel, W. W. Hargrove, and F. M. Hoffman (2008), A continental strategy for the National Ecological Observatory Network, *Front. Ecol. Environ.*, *6*, 282–284, doi:10.1890/1540-9295(2008)6[282:ACSFTN]2.0.CO;2.
- Kyriazis, D., K. Tserpes, G. Kousiouris, A. Menychtas, and T. Varvarigou (2008), Data aggregation and analysis: A grid-based approach for medicine and biology. Int. Symp. on. IEEE 841–848.
- Liang, X., and D. P. Lettenmaier (1994), A simple hydrologically based model of land surface water and energy fluxes for general circulation models, *J. Geophys. Res.*, *99*, 14,415–14,428, doi:10.1029/94JD00483.
- Maidment, D. R. (2008), Bringing water data together, *J. Water Resour. Plann. Manage.*, *134*, 95–96.
- Michener, W. K., S. Allard, A. Budden, R. B. Cook, K. Douglass, M. Frame, S. Kelling, R. Koskela, C. Tenopir, and D. A. Vieglais (2012), Participatory design of DataONE—Enabling cyberinfrastructure for the biological and environmental sciences, *Ecol. Inf.*, *11*, 5–15, doi:10.1016/j.ecoinf.2011.08.007.
- Minnesota Population Center (2013), *Terra Populus: Beta Version [Machine-Readable Database]*, Univ. of Minnesota, Minneapolis.
- Morsy, M. M., J. L. Goodall, C. Bandaragoda, A. M. Castronova, and J. Greenberg (2014), Metadata for describing water models International Environmental Modelling and Software Society (iEMSs) 7th International Congress on Environmental Modelling and Software, doi:10.13140/2.1.1314.6561.
- Myers, J., et al. (2015), Towards sustainable curation and preservation: The SEAD Project's data services approach. Proc. IEEE 11th Int. e-Science Conf. Munich, Ger., doi:10.1109/eScience.2015.56.
- Oinn, T., et al. (2004), Taverna: A tool for the composition and enactment of bioinformatics workflows, *Bioinformatics*, *20*, 3045–3054, doi:10.1093/bioinformatics/bth361.
- Rajasekar, A. (2014), Workflows [WWW document]. 6th Annu. iRODS User Gr. Meet. June 2014 Inst. Quant. Soc. Sci. MA. [Available at <http://irods.org/wp-content/uploads/2014/06/Workflows-iRUGM-2014.pdf>, accessed 8.12.15.]
- Rajasekar, A., et al. (2010), iRODS Primer: Integrated rule-oriented data system. Synthesis lectures on information concepts, retrieval, and services. doi:10.2200/S00233ED1V01Y200912ICR012.
- Tarboton, D. G., et al. (2014), HydroShare: Advancing collaboration through hydrologic data and model sharing, in *International Environmental Modelling and Software Society (iEMSs) 7th International Congress on Environmental Modelling and Software*, San Diego, Calif., edited by D. P. Ames, N. W. T. Quinn, and A. E. Rizzoli, doi:978-88-9035-744-2.
- Vahi, K., et al. (2013), A general approach to real-time workflow monitoring. In *High Performance Computing, Networking, Storage and Analysis (SCC)*, pp. 108–118. doi:10.1109/SC.Companion.2012.26.
- Weise, A., M. Wan, W. Schroeder, and A. Hasan (2008), Managing groups of files in a Rule Oriented Data Management System (iRODS), *Comput. Sci.*, *5103*, 321–330, doi:10.1007/978-3-540-69389-5\_37.
- Williams, D. N., et al. (2008), Data management and analysis for the Earth System Grid, *J. Phys. Conf. Ser.*, *125*, 012072, doi:10.1088/1742-6596/125/1/012072.
- Williams, D. N., B. N. Lawrence, M. Lautenschlager, D. Middleton, and V. Balaji (2011), The Earth System Grid Federation: Delivering globally accessible petascale data for CMIP5, Proceedings of the 32nd Asia-Pacific Advanced Network Meeting, pp. 121–130. doi:10.7125/APAN.32.15.
- Workflow Objects (WSO) [WWW Document] (2013), [Available at [https://wiki.irods.org/index.php/Workflow\\_Objects\\_\(WSO\)](https://wiki.irods.org/index.php/Workflow_Objects_(WSO)), accessed 6.7.15.)]