

# 1    **Using a Data Grid to Support Regional-Scale Hydrologic** 2    **Modeling**

3    Mirza M. Billah<sup>1</sup>, Jonathan L. Goodall<sup>2,3,\*</sup>, Ujjwal Narayan<sup>4</sup>, Bakinam T. Essawy<sup>2</sup>,  
4    Venkat Lakshmi<sup>5</sup>, Arcot Rajasekar<sup>6</sup>, Reagan W. Moore<sup>6</sup>

5

6    <sup>1</sup> *Department of Biological Systems Engineering, Virginia Tech, Blacksburg, Virginia*

7    <sup>2</sup> *Department of Civil and Environmental Engineering, University of Virginia,*  
8    *Charlottesville, Virginia*

9    <sup>3</sup> *Department of Civil and Environmental Engineering, University of South Carolina,*  
10    *Columbia, SC*

11    <sup>4</sup> *CMNS-Earth System Science Interdisciplinary Center, University of Maryland, College*  
12    *Park, MD*

13    <sup>5</sup> *Department of Earth and Ocean Sciences, University of South Carolina, Columbia, SC*

14    <sup>6</sup> *School of Library and Information Science, University of North Carolina, Chapel Hill,*  
15    *NC*

16    \* *To whom correspondence should be addressed (E-mail: [goodall@virginia.edu](mailto:goodall@virginia.edu);*  
17    *Address: University of Virginia, Department of Civil and Environmental Engineering,*  
18    *PO Box 400742, Charlottesville, Virginia 22904; Tel: (434) 243-5019)*

## Abstract

Modeling a regional-scale hydrologic system introduces major data challenges related to the access and transformation of heterogeneous datasets into the information needed to execute a hydrologic model. These activities are difficult to automate, making the reproducibility and extensibility of model simulations conducted by others difficult or even impossible. This study addresses this challenge by demonstrating how the integrated Rule Oriented Data Management Systems (iRODS) can be used to support workflow execution needed when using data-intensive models to simulate regional-scale hydrologic systems. Focusing on the Variable Infiltration Capacity (VIC) model as a case study, data preparation steps are written using micro-services and rules within iRODS. VIC and iRODS are applied to study hydrologic conditions in the Carolinas, USA during the period 1998-2007 to better impacts of drought within the region. The application demonstrates how automated workflows in iRODS can support hydrologic modelers to create more reproducible and extensible “end-to-end” model simulations.

*Keywords:* data management; workflows; hydrologic modeling

*Software availability:* Software is available by request to the corresponding author.

# 1 INTRODUCTION

## 2 Motivation

3       Application of regional-scale hydrologic models presents a number of challenges  
4 associated with handling and processing large datasets. The models are data intensive and  
5 require a significant amount of time and effort in order to transform available datasets  
6 into the form required by the hydrologic model (Leonard and Duffy, 2013). The data  
7 challenges also include collecting datasets from various data providers that have  
8 inconsistent data access protocols, file formats, and semantics (Horsburgh et al., 2014).  
9 The result is that the datasets required to setup, calibrate, and validate hydrologic models  
10 are not made available by data providers in a form that is ready for application directly  
11 into models, rather each model requires specific transformations of the available  
12 information before it can be used within a model (Leonard and Duffy, 2014). Due to the  
13 level of heterogeneity between data sources, these data preparation steps are difficult to  
14 automate and, with the exception of a few models with robust data preparation tools,  
15 often require significant manual intervention. Even for models with data preparation tools  
16 available, they are often tied to pre-specified data sources that may not be the best  
17 available information for a specific region or modeling objective. The end result is that  
18 modelers (i) consume significant time on tasks that could be automated and (ii) lack the  
19 ability to easily reproduce and extend past work completed by others to explore new  
20 scientific questions or water management goals.

21       Within the information and computer science communities, there has been work  
22 to create advanced data management and scientific workflow software that has the  
23 potential to address these challenges facing the hydrology community (Gil et al., 2007;

1 Ludäscher et al., 2006; Oinn et al., 2006). There has been uptake of these tools within  
2 some scientific communities, in particular bioinformatics, but with only a few exceptions  
3 (Fitch et al., 2011; Guru et al., 2009; Perraud et al., 2010; Piasecki and Lu, 2010), there  
4 has been minimal adoption of workflow tools for hydrologic modeling. This past work  
5 has clearly demonstrated the utility of workflow environments for implementing the data  
6 preparation tasks and allowing the software to coordinate and automate processing steps,  
7 as well as track the provenance of the datasets generated through the processing steps.  
8 The potential of this software to provide authenticated access to external data sources,  
9 along with procedures for automatically transforming data products to those required by a  
10 model, makes these software systems particularly well suited to automating hydrologic  
11 modeling workflows. However these workflow tools have challenges when used for  
12 hydrologic analysis and modeling that include but are not limited to the lack of common  
13 data models within hydrology (Perraud et al., 2010).

14 Building on this paper work, we postulate that a key reason for the lack of update  
15 of workflow environments within hydrology is the data-centric nature of hydrology, and  
16 the vast heterogeneity of datasets used within hydrologic models. For this reason, we  
17 believe that data management must play a central role within hydrologic workflows.  
18 Hydrologic data sources are distributed and maintained by different governmental and  
19 nongovernmental agencies. They are provided using different file types, structures, and  
20 semantics. We believe that centralizing and standardizing this data is not a sustainable  
21 long-term approach. Rather the long-term solution, we believe, will be a federated data  
22 grid approach with decentralized data management supplemented with server-side  
23 processing that allows for on-demand access to derived data produces required by

specific hydrologic models. Thus, when we speak of workflows in this paper, we are describing processing pipelines that are able to operate on servers that house large data achieves and that can deliver to client machines, where hydrologic models are executed, derived data products at are much closer to the information required as input by the model.

This paper illustrates our concept by focusing on the DataNet Federation Consortium (DFC) grid, which is built using the data management system Integrated Rule-Oriented Data System (iRODS), and the regional-scale hydrologic model Variable Infiltration Capacity (VIC). These systems are briefly introduced in the background section that follows this introduction section. Then, in the design and implementation section, the approach for automating VIC data preparation and post-processing workflows using iRODS is presented. The workflows were demonstrated for an example of modeling drought in the Carolinas region of the United States. Finally, this paper concludes with a discussion of applying iRODS to support hydrologic modeling.

## **Background**

The DFC grid (<http://www.datafed.org>) was built as part of an NSF-funded research project to provide storage and compute resources that allow for long-term access to the stored datasets. DFC is enabled by the iRODS, an open source, policy-based cyberinfrastructure developed by the Data Intensive Cyber Environments (DICE) group for distributed data management (Rajasekar et al., 2010b) (<http://www.irods.org>). It is used by a wide variety of end users including the bioinformatics community (Goff et al., 2011). Data management tasks or policies are implemented within iRODS as rules. Rules specify a sequence of lower-level micro services that are used to operate on

1 datasets within the data grid. The user can specify the sequences of data collection,  
2 transformation, curation, preservation, and processing steps in a workflow within one or  
3 more rules that use micro-services developed by users or administrators. The system is  
4 capable of remote execution of workflows as well, meaning data preparation steps can be  
5 executed on remote servers rather than on a client machine, which is especially beneficial  
6 for large-datasets.

7       The Variable Infiltration Capacity (VIC) model is a large-scale hydrologic model  
8 that applies water and energy balances to simulate terrestrial hydrology at a regional  
9 spatial scale (Liang et al., 1996a). The scientific background of the model is summarized  
10 in the case study section, while here the model is presented from a data management  
11 perspective. The model requires several input datasets to be generated by applying  
12 multiple data processing steps before a simulation run can be executed (Figure 1). The  
13 datasets for preprocessing include precipitation, maximum and minimum temperature,  
14 wind speed, topography, soil, and vegetation information. Each dataset is processed using  
15 scripts that require the execution of data processing routines to generate new datasets  
16 used in the model simulation. Each data processing script performs a certain task that is a  
17 prerequisite to the following script in the workflow. Running these data collection and  
18 processing scripts currently requires significant manual intervention and time to  
19 complete. In addition, a large amount of data is generated during the procedure,  
20 introducing issues with storing intermediate datasets for reproducibility of model results  
21 as well as further analysis of key model inputs.

## Study Objective and Scope

Given these data challenges in running VIC, which are not dissimilar from any other data-intensive hydrologic simulation model, and given the potential advantages of advanced data management systems such as iRODS, the objective of this study is to apply iRODS for automating the execution of VIC. The study is built on prior work where the VIC model was used for modeling water balances for South Carolina river basins (Billah and Goodall, 2011; Billah et al., 2015). VIC has also been applied by others for a number of climate conditions and basins to estimate several hydrologic variables with high accuracy (Abdulla et al., 1996; Lakshmi et al., 2004; Lohmann et al., 1998; Sheffield and Wood, 2007; Sheffield et al., 2004). However, there has been less work on the data challenges associated with running VIC. Therefore, this work addressed these data challenges by integrating VIC with iRODS in the DFC grid using federated data access and data processing for executing VIC for regional-scale hydrologic analysis.

The resulting system focused on developing data management workflows that automate pre and post-processing of data for VIC model and visualize the model results for state water managers for drought decision-making. The work demonstrated a case study for the drought in South Carolina during 1998 to 2002 (Badr, A. W., Wachob, A., Gellici, 2004) that had a significant effect on water resources across the Carolinas. The availability of detailed spatial and temporal information of hydrologic systems across the Carolinas, and in particular the ability to forecast future conditions, is a useful tool for water resources management. Soil moisture in particular is a difficult parameter to observe at state-level spatial scales (Sheffield et al., 2004), but it is an important indicator of drought at various scales (Lakshmi et al., 2004; Sheffield et al., 2004). The modeling





1 system created through this study provides estimates of soil moisture across the Carolinas  
2 for the period of 1998 to 2007. The analysis demonstrated insights of the methodological  
3 advancement in terms of accessibility, reproducibility, human intervention, time  
4 consumption, and storage management for the VIC model application for hydrologically-  
5 based drought analysis at regional spatial scales

## 6 **SYSTEM DESIGN AND IMPLEMENTATION**

7 The presentation of the system design and implementation is organized within  
8 two broad categories. The first category discusses the server-side application for  
9 workflow implementation and managing datasets, while the second category focuses on  
10 the client-side application for workflow execution using the server-side setup. Both the  
11 server and client-side tool are designed to manage datasets using iRODS for VIC model.  
12 Tasks associated with these categories are described in the following subsections.

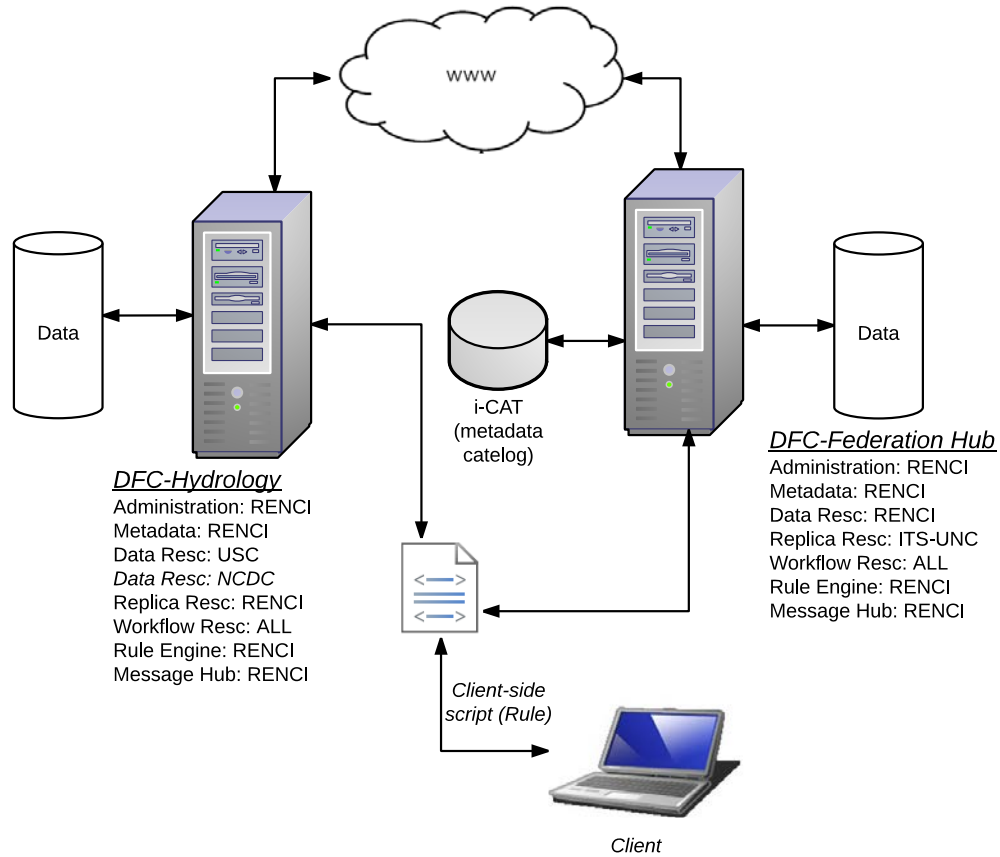
### 13 **Server-side Design and Implementation**

14 The server-side application for iRODS is deployed on a remote server for  
15 hydrologic analysis that is part of the DFC federated grid. The application uses input  
16 datasets either from remote web services or a federated grid. The federated grid is a  
17 cyberinfrastructure that remotely connects several discipline specific resource servers.  
18 The server accommodates data processing sequences to produce output that can be  
19 transferred to a client.

#### 20 *The DFC-Hydrology Grid*

21 The DFC-Hydrology Grid is deployed as part of the larger DFC (Figure 2), which is a  
22 project funded by the National Science Foundation (NSF) to support data collection,

1 analysis, preservation, sharing, and publication of scientific data and models  
2 (<http://www.datafed.org>). The DFC- Hydrology grid consists of an iRODS server that  
3 communicates with the iRODS metadata catalog (iCAT) database located in the DFC-  
4 Federation Hub of the DFC grid. The server is administered by RENCI (Renaissance  
5 Computing Institute) in Chapel Hill, North Carolina. The server contains a Rule Engine  
6 (RE) that interprets rules (executed from the client-side) using micro-services on the  
7 DFC-Hydrology data grid or the DFC-Federation Hub. The RE also connects to the  
8 catalog server in the DFC-Federation Hub and updates the iCAT database. The hydrology  
9 related data collected from various external sources are stored and preserved in the DFC-  
10 Hydrology server. The workflow can be implemented on any of the federated servers for  
11 data collection and transformation, however, transformed datasets are sent back to the  
12 client-side. These transformed datasets are replicated using RENCI operated resources  
13 within the DFC federated grid after completing data processing tasks in the client- side.  
14 For instance, the collection of precipitation data from a remote location is done by the  
15 DFC-Federation Hub and transformation is performed in the DFC-Hydrology server. The  
16 iRODS server-side application automates the process of accessing remote sources using  
17 different protocols (e.g., HTTP or FTP) for retrieving datasets. The iRODS server-side  
18 application is responsible for managing resource storage and federating these with other  
19 DFC grids, thereby providing seamless access to the data stored in DFC. These grids  
20 provide continuous support for data access and, in the case of connection failure, the  
21 federated grids provide data access from one or more other storage resources.



1

2 Figure 2: Schematic diagram of the NSF supported Data Federation Consortium (DFC)  
3 data management system showing the connections between the DFC-Hydrology with the  
4 DFC-Federation Hub.

## 5 *Micro-service Implementation*

6 Micro-services are the building blocks for implementing policy-based data  
7 management within DFC grid (Rajasekar et al., 2010a). A well-defined function in a  
8 micro-service performs a specific task as part of a distributed workflow system. A  
9 number of micro-services are available to automate data collection, processing, and  
10 storage in the DFC federated resource servers. These micro-services are primarily  
11 developed by system or application programmers and are compiled into the iRODS  
12 server code. The micro-services that are applied for this study are listed in Table 1. These

1 micro-services are chained together to provide a higher-level of functionality to  
2 implement multiple tasks. Although the flexibility to chain a number of micro-services  
3 provides multiple ways to complete a series of tasks, iRODS applies priorities and  
4 validation conditions to select the best micro-service to complete a given task.  
5 Application of multiple micro-services in iRODS makes it possible to chain routines to  
6 perform multiple tasks within a single workflow.

## 7 **Client-side Design and Implementation**

8 Client-side application is used to execute data management workflow within the  
9 federated grid. This consists of several command-line utilities known as i-commands to  
10 manage datasets and data processing rules to execute data management workflow  
11 (Rajasekar et al., 2010a). These i-commands are used to download and upload data  
12 from/to the DFC grids and execute micro-services in the server-side. Furthermore, the  
13 client-side application provides opportunity to create and execute rules to gain access to  
14

15 Table 1: Micro-services applied for VIC model application using iRODS.

| No. | Micro-service               | Purpose                             |
|-----|-----------------------------|-------------------------------------|
| 1   | msiExecCmd                  | Execute commands                    |
| 2   | msiCollCreate               | Make data collection                |
| 3   | msiDataObjCreate            | Create data object                  |
| 4   | msiDataObjWrite             | Write data object                   |
| 5   | msiDataObjClose             | Close data object                   |
| 6   | msiAddSelectFieldToGenQuery | Make query for data using field     |
| 7   | msiAddConditionToGenQuery   | Make query for data using condition |
| 8   | msiExecGenQuery             | Execute Query                       |
| 9   | msiGetValByKey              | Extract value from query result     |
| 10  | msiSplitPath                | Get directory path                  |
| 11  | msiDataObjUnlink            | Delete temporary file               |
| 12  | msiRmColl                   | Remove data collection              |
| 13  | msiGetSystemTime            | Get time stamp                      |

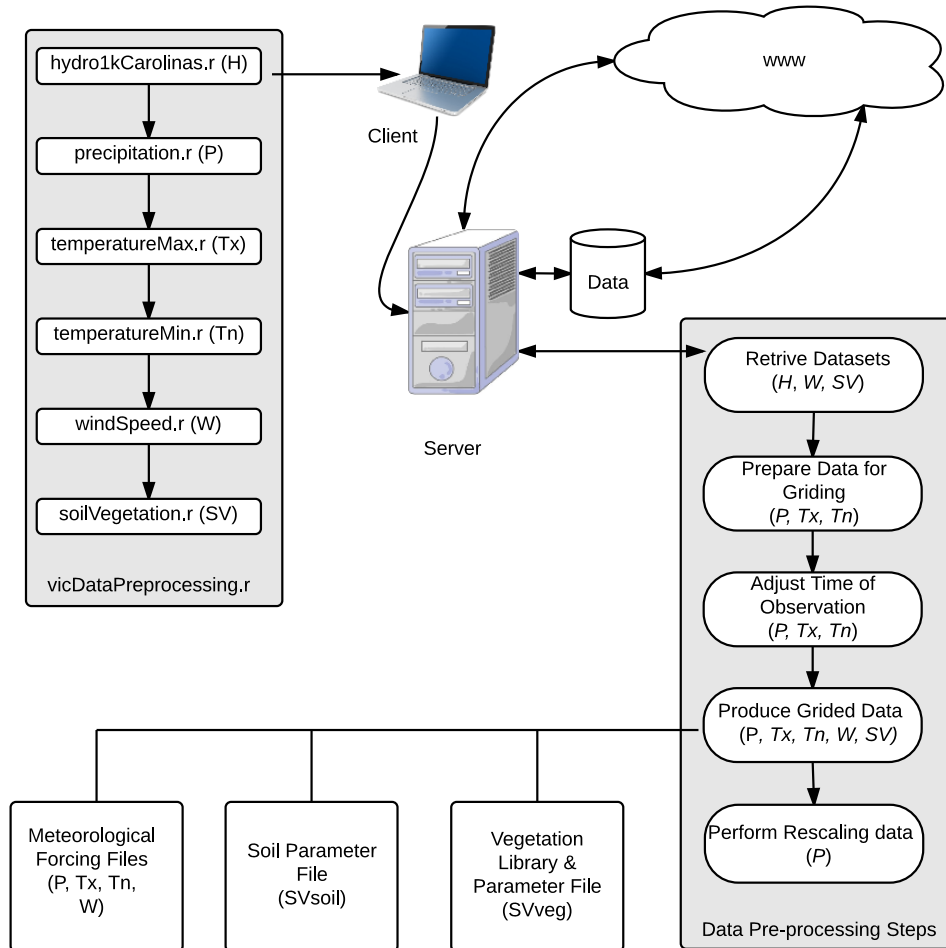
heterogeneous data sources, to process datasets into the formats required by models or scientific communities. The rule or action is a critical and fundamental component for iRODS. It provides a flexible mechanism to integrate external systems for specialized processing and metadata management (Hedges et al., 2009, 2007). The rules implement data processing steps as policies on the DFC federated grid and respond to various requests and conditions by integrating related micro-services from the server-side. These data processing rules, specifically for the VIC model, are generally categorized into two separate workflows: i) data pre-processing and ii) data post-processing. The pre-processing workflow is responsible for data collection and transformation into standard VIC format, while the post-processing workflow is used for model results visualization. Details of these categories are discussed in the following.

#### *VIC Pre-processing Workflows*

Data pre-processing involves collecting and transforming datasets from heterogeneous sources. For our purpose, data are collected from the United States Geological Survey (USGS), National Climatic Data Center (NCDC), National Center For Atmospheric Research (NCAR), National Centers for Environmental Prediction (NCEP), and Land Data Assimilation System (LDAS). These datasets are then processed in the DFC-Hydrology grid and stored in the DFC-Federation Hub. The metadata catalog in the DFC-Federation Hub is updated automatically while data are uploaded with recent information. This metadata catalog functions as an information center and enables the discovery of data that are stored into the grid.

Data processing workflows also include several data-specific processing rules that transform the collected datasets from the external sources into model readable inputs

1 (Figure 3). This data-specific rule is a combination of multiple step-based routines each  
 2 of which performs a particular task. The step-based routines are simple, such as retrieving



3  
 4 Figure 3: Model pre-processing workflows showing the major steps for transforming  
 5 datasets to set up the VIC model for a specific study area. Rules are initiated from a client  
 6 but executed on a server using micro-services.

7 data, preparing data for gridding, adjusting observation times, and transforming gridded  
 8 and rescaled datasets. For our study, these step-based routines are integrated into data-  
 9 specific rules used to complete a series of tasks from collection and transformation of the  
 10 datasets into model inputs. Therefore, for each dataset, a separate rule is created to

1 perform the data transformation tasks including the data collection from respective  
2 sources.

3

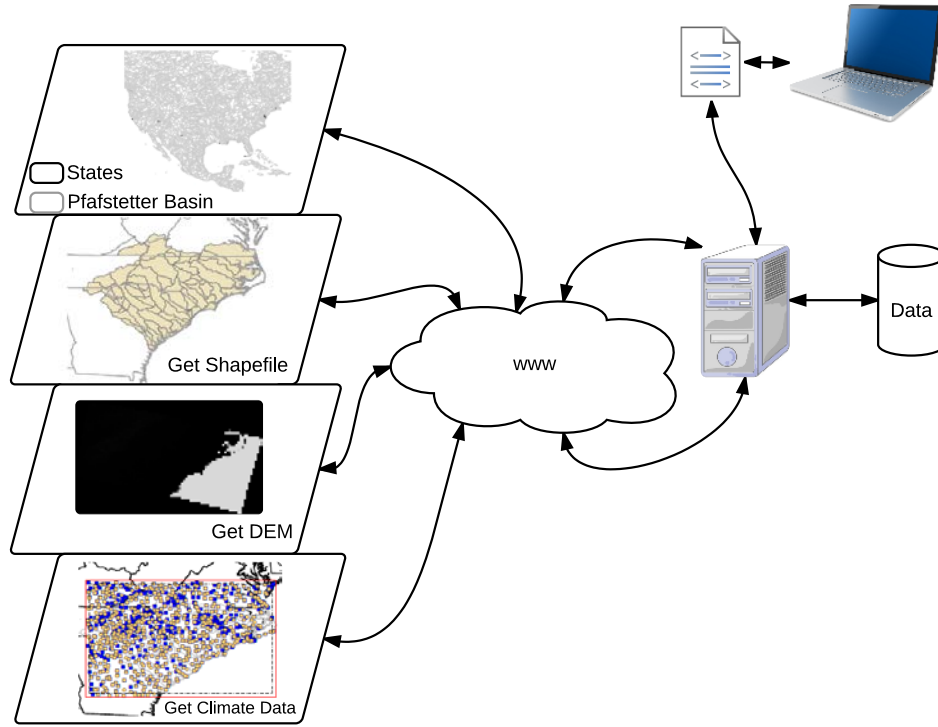
4 Table 2: Rules created within iRODS for VIC data pre-processing.

| No. | Micro-service      | Purpose   |
|-----|--------------------|---|
| 1   | hydro1kCarolinas.r | Collect HYDRO1k and Climate data for the study area |
| 2   | precipitation.r    | Collect and process daily precipitation data        |
| 3   | temeperatureMax.r  | Collect and process daily maximum temperature data  |
| 4   | temeperatureMin.r  | Collect and process daily minimum temperature data  |
| 5   | windSpeed.r        | Collect and process annual wind speed data          |
| 6   | soilVegetation.r   | Collect and process soil and vegetation data        |

5

6       The sequential implementation of the data-specific rules can be integrated into a  
7 single model-specific rule called *vicDataPreprocessing.r* to process complete data inputs  
8 for the VIC model. However, we have used six separate data-specific rules to process the  
9 datasets for VIC model for the period of 1998 to 2007 (Table 2). In the data-specific  
10 rules, step-based calculations are chained in such a way that data processing steps are not  
11 violated and perform only designated tasks (Figure 3). For example, preparation for data  
12 gridding is not executed without collecting and storing data in the DFC grid or rescaling  
13 is not performed before gridding the datasets. Also, not all of the data processing routines  
14 are performed for all the datasets. For instance, preparation for data processing is not  
15 executed for wind, soil and vegetation datasets because these data do not require grid  
16 preparation. Overall, data retrieval is performed by hydro1kCarolinas.r (H), windSpeed.r  
17 (W), and soilVegetation.r (SV), while precipitation.r (P), temperatureMax.r (Tx), and  
18 temperatureMin.r (Tn) implement respective workflows for data trans-formation using  
19 the retrieved climate data.

1 Data-specific rules are mostly inherent data processing scripts from the VIC that  
2 execute a series of routines that are associated with tasks (Table 2). For instance, we have  
3 retrieved climate data (precipitation, maximum and minimum temperature data) from the



4  
5 Figure 4: Data flow in the hydro1kCarolinas.r rule that extracts climate data from NCDC  
6 GHCND using HYDRO1K basin/DEM datasets to define a study region.

7 DFC federated grid by implementing a rule hydro1kCarolinas.r (Figure 4). This rule uses  
8 a series of micro-services in the DFC federated grid and executes a series of tasks  
9 sequentially to extract and register data over a defined area. We used the Pfafstetter basin  
10 numbering system as described in Furans and Olivera (2001) for defining study area from  
11 HYDRO1k basin and DEM datasets. The climate datasets for the defined study area were  
12 downloaded via FTP from the NCDC Global Historical Climatology Network (GHCND)  
13 database. While downloading the climate data, a buffer of 0.25° was considered around



the defined study area to collect sufficient climate data. The climate data contained precipitation, maximum and minimum temperature, and wind speed datasets. The precipitation data was processed using the precipitation.r rule, which used the GHCND data downloaded through DFC server and converted the station specific datasets into gridded datasets with a spatial resolution of  $1/8^\circ$ . Similarly, temperatureMax.r and temperatureMin.r rules downloaded and converted station specific temperature values into gridded datasets with  $1/8^\circ$  spatial resolution. Furthermore, we also collected wind speed, soil and vegetation data from their respective sources while executing respective data-specific rules in the DFC-Hydrology. The annual wind data were collected from NCAR/NCEP and processed to generate gridded datasets of  $1/8^\circ$  resolution for the study area using windSpeed.r. The soilVegetation.r rule was applied to transform the LDAS soil and vegetation information into information required for the model.

### *VIC Post-processing Workflows*

A VIC model simulation outputs hydrologic flux and state variables for a gridded discretization of the landscape. The hydrologic flux and state variables include evapotranspiration, soil moisture, baseflow, and snow depth. These variables are output into text files and require additional processing in order to visualize the model results and gain a better understanding of water movement within the system being studied.

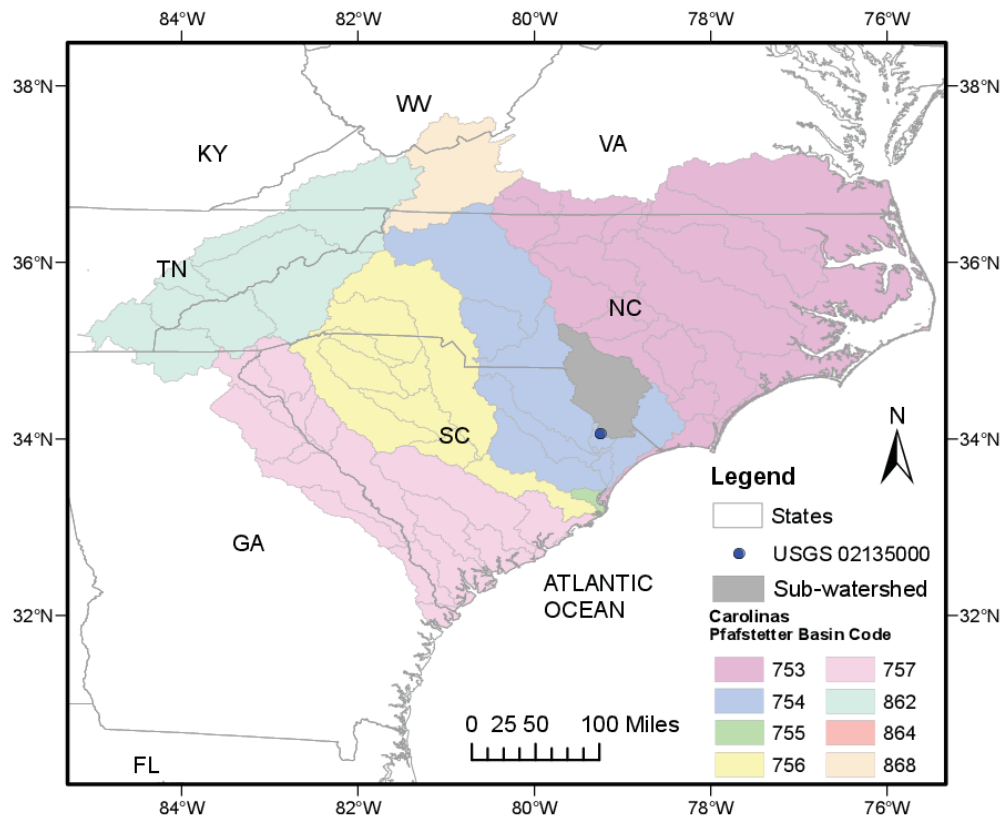
Workflows can play an important role in this regard because they are able to automate the steps required to transform model output data into more useful visualizations such as maps or graphs. Scientists can use workflows to more easily visualize the VIC model predictions and can share their approaches for creating model visualizations with other VIC modelers. Because rules are written in a simple scripting language, it is possible for

1 multiple stakeholders to make use of and potentially modify data post-processing rules to  
2 create useful visualizations and analysis routines that summarize model predictions.

### 3 **EXAMPLE APPLICATION**

#### 4 **Study Area**

5 The region used for the case study application covers all the major river basins in  
6 both North and South Carolina with a total area of 280,736 km<sup>2</sup> (108,393 mi<sup>2</sup>) (Figure 5).



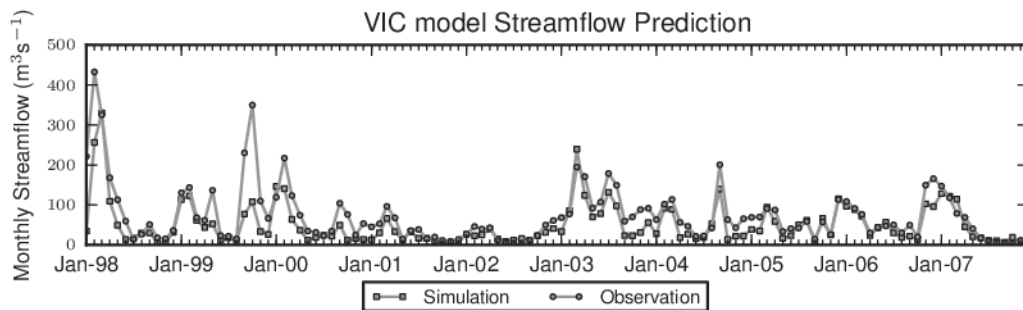
7 Figure 5: Study area with major river basins in North and South Carolina, USA. These  
8 subwatersheds were extracted using selected Pfafstetter Basin Code from the Hydro1K  
9 dataset. The sub-watershed and stream gage station used for the VIC model calibration  
10 are also shown in the Figure.  
11

1 The Pfafstetter Basin Code system was used to define this “Carolinas” study region from  
2 the HYDRO1k basin dataset. The codes were incorporated in the hydro1kCarolina.r rule  
3 to extract the area during the automation of the workflow.

#### 4 **Model Setup**

5 VIC is able to simulate the land surface portion of the hydrologic cycle by solving  
6 the full water and surface energy balance equations (Liang and Lettenmaier, 1994; Liang  
7 et al., 1996b) in the Carolinas on a daily time step. VIC was calibrated using the  
8 following seven parameters: variable infiltration curve (b), maximum base flow ( $D_{smax}$ ),  
9 fraction of base flow where base flow occurs ( $D_s$ ), fraction of maximum soil moisture  
10 content above which nonlinear base flow occurs ( $W_s$ ), mid ( $d_2$ ) and deep ( $d_3$ ) soil layer  
11 depth, and minimum stomatal resistance ( $r_0$ ) (Abdulla and Lettenmaier, 1997a, 1997b;  
12 Crow, 2003; Troy et al., 2008). The range investigated and the final values of the  
13 parameters applied in this study were described in (Billah et al., 2015) and a comparison  
14 of the monthly average streamflow for the calibrated VIC model and streamflow from the  
15 USGS streamflow station are provided in (Figure 6). A trial and error approach was  
16 followed for calibrating the model for a selected portion of the overall study region for  
17 different parameter values. This approach was used because the model execution time for  
18 the overall study area was too long to use the entire study region for calibration. We  
19 compared the VIC model prediction for streamflow at the Little Pee Dee River at  
20 Galivants Ferry stream gage station that is part of the USGS National Water Information  
21 System (NWIS) network (USGS 02135000; Figure 5) for the period of 1998 to 2007.  
22 This streamflow station has a drainage area of 7257 km<sup>2</sup> and includes portions of both  
23 North and South Carolina. This station was selected based on its available time series

1 record and because it is on an unmanaged portion of a river network. Nash-Sutcliffe  
2 Efficiency (NSE) index between simulated streamflow using VIC routing scheme and  
3 observed streamflow for the station, which is the commonly used approach for evaluating  
4 VIC models, was used as an objective function in the calibration. The NSE index of the  
5 final calibration is 0.6, a value that is considered to be a satisfactory calibration by  
6 watershed-scale hydrologic modelers (Moriassi et al., 2007).

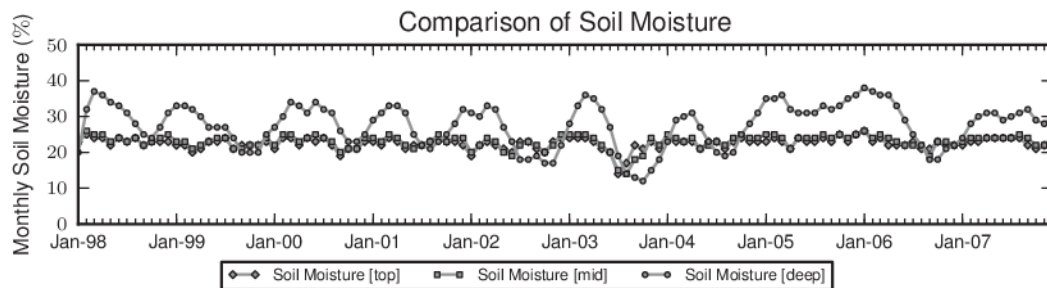


7  
8  
9 Figure 6: Streamflow comparison between the VIC model predictions and USGS  
10 observations. The comparison is performed at the USGS station Little Pee Dee River at  
11 Galivants Ferry, SC (Station Number 02135000).  
12

### 13 **Model Results**

14 VIC model runs generate a number of grid-based hydrologic flux and state  
15 variable outputs. One of these outputs is soil moisture estimated for each grid cell in the  
16 simulation domain and for each soil layer within the model. Soil moisture is an important  
17 indicator of drought, so it was used as the example model output for creating a post-  
18 processing rule. This post-processing rule (soilmoisture.r) works by extracting model  
19 results from the various model output files, then summarizing the soil moisture over the  
20 study region for each soil layer. The rule results in a time-series plot of monthly soil  
21 moisture within the three VIC soil layers (Figure 7). The plot of soil moisture estimates  
22 provides understanding of the impact of drought on soil moisture, particularly the deep

1 soil layer which is more sensitive to long term trends in water availability compared to  
2 the middle and upper soil layers. The plot represents the final product from the modeling  
3 work that would be used in a journal publication describing the work. Because generation  
4 of the plot can be easily repeated by re-running the soilmoisture.r rule, the results  
5 presented in the journal publication can be more easily reproduced and extended by other  
6 researchers.



7  
8 Figure 7: Comparison of monthly averaged soil moisture in the three soil layers predicted  
9 by the VIC model in the Carolinas for the periods of 1998 to 2007.

## 10 SUMMARY, DISCUSSION, AND CONCLUSIONS

11 The hydrologic modeling process can involve many steps from data access and  
12 transformation, to model setup, calibration, and validation, to analysis and visualization  
13 of model outputs. This entire “end-to-end” process involves some steps that are easily  
14 automated and others that often require intervention by expert modelers. The goal in this  
15 and related work is to automate those steps that are straightforward but tedious, while still  
16 allowing experts to guide the process and intervene when needed. Currently too many  
17 steps that can be automated are not. As a result, modelers are unable to focus on the  
18 important tasks that require their expertise and insights because time must be spent on  
19 more basic data gathering and transformation steps. Furthermore, the steps that could be  
20 automated are typically not thoroughly documented and, even if they are thoroughly

1 documented, are time consuming to repeat. This makes independent reproducibility of  
2 model results, a requirement for scientific progress and water resource management  
3 objectives, difficult or even impossible.

4       This work addresses these challenges by leveraging the iRODS technology and  
5 the DataNet Federation Consortium (DFC) cyberinfrastructure to create workflows that  
6 automate pre and post-processing routines for the hydrologic model VIC. VIC is a widely  
7 used hydrologic model that is typically applied to large spatial regions to address  
8 questions related to water resource availability during periods of drought. The workflows  
9 developed for VIC include sufficient information to allow others to independently  
10 reproduce the model results, from raw data products to visualizations of model outputs  
11 used in publications. The workflows therefore act as a means for forced documentation of  
12 the steps used to create model input files including the provenance of data as it is  
13 transformed from the form provided by federal and academic data repositories to the  
14 form output by the models. In the case of VIC, and likely in the case of other hydrology  
15 applications as this work is extended to include other hydrologic models, the workflows  
16 were created by leveraging existing scripts written to complete specific data pre or post-  
17 processing tasks. The approach used provides a means for placing these scripts in larger  
18 workflows and removes the need to access and understand the original scripts to  
19 reproduce model results or to reuse the tools for a new study. In the latter case,  
20 modification is only required when selecting a new area of interest. The Pfafstetter Basin  
21 Codes are replaced in the shared workflow for hydrologic systems analysis workflow  
22 hydro1kCarolinas.r. However, all other workflows are used without modification to

1 automate the remaining data processing steps, which effectively build a complete VIC  
2 model that is ready for simulation.

3 A key contribution of this work is demonstrating a methodological approach to  
4 assist in data-intensive hydrologic modeling. Using iRODS has advantages that include  
5 workflow automation, datasets curation and preservation, data replication, and workflow  
6 sharing for results reproduction. iRODS is flexible and robust and it was possible to  
7 extend the software to develop rules specific for the data processing workflows  
8 associated with running the VIC model. Because of the data grid concept used by iRODS,  
9 it was possible to design the workflows in such a way that both server-side and client-  
10 side applications could be leveraged to perform the data gathering and preparation steps.  
11 It was also possible to create shared workflows that emphasized the interoperability of the  
12 DFC federated grids for sharing rules and datasets among distributed users. For instance,  
13 we applied ecohydroworkflow (Miles, B., Band, 2013), a shareable workflow for data  
14 management for hydrologic models, to collect and register GHCND and HYDRO1k  
15 datasets from NCDC and USGS respectively in the DFC-Federation Hub. This  
16 implementation provided an example of how datasets collected for multiple models and  
17 applications can be easily transferred between grids within an iRODS federation using i-  
18 commands. The ability to execute data-specific rules on the client-side allowed for end-  
19 to-end data management, reduction of time and potential for human error during data  
20 processing, model prediction reproducibility, large data management, and opportunity to  
21 analyze model predictions before registering and sharing datasets within the DFC  
22 federated grid.

1       The workflows created within iRODS to enable end-to-end execution of model  
2 pre and post-processing steps are a key step to achieving reproducible hydrologic model  
3 runs. It was possible to chain all the data pre-processing routines by including all datasets  
4 in model-specific rule vicDataPreprocessing.r using micro-services in iRODS. Despite  
5 the opportunity to combine all data processing routines into a model-specific rule, the  
6 system was purposely designed in such a way that each dataset had its own sub-workflow  
7 to transform the raw data into model readable information. Doing so provided a level of  
8 granularity so that the data-specific rules could be later re-used within other hydrologic  
9 modeling applications in the DFC federated grid. Data post-processing workflows were  
10 also created to automate the tasks required to create visualizations and publication-quality  
11 figures based on model outputs.

12       By running data pre-processing, model execution, and data post-processing rules  
13 in sequence, it is possible to go from raw data from federal data repositories to figures  
14 summarizing model outputs that are used in journal publications and conference  
15 presentations. Having such capability would reduce human errors that could occur by not  
16 correctly performing a transformation step during a manual execution of the work. It  
17 would also free researchers to devote more time to enhancing, calibrating, and validating  
18 models, rather than on tedious steps required to set-up first iterations of the model. A  
19 longer term goal could be for researchers to publish such workflows that can be used to  
20 recreate publication figures as supplemental resources with the journal paper itself.

21       It is important to note the specific data challenges required for a hydrologic model  
22 application. First, VIC requires a large amount of data when applied to a region the size  
23 of the Carolinas, and of course even more data would be required for Continental scale



1 model executions, which are not uncommon when applying VIC. The data used in the  
2 VIC model pre-processing steps included meteorological datasets at stations and on grids,  
3 topography datasets available as grids, and soil and vegetation datasets also available as  
4 grids. Transformation of these raw input data resulted in intermediate datasets with  
5 different spatial projections, filled gaps, and other modifications required before initiating  
6 the model simulation. Over the years, researchers have created scripts for completing  
7 many of these data pre-processing steps, and so it was relatively straight forward to wrap  
8 these scripts as iRODS rules using i-commands. One advantage of using iRODS rather  
9 than just running the scripts outside of iRODS is that iRODS provides a metadata server  
10 in the DFC-Federation Hub that is automatically updated with information tracking the  
11 provenance of the datasets during the transformation process. Also, the large datasets  
12 used in the VIC modeling exercise were automatically curated and preserved in the DFC  
13 federated grid for future uses and can be referenced through publications describing the  
14 analysis.

15       Creating data post-processing workflows provided the opportunity to visualize  
16 model results more rapidly and interactively as part of the modeling process. The data  
17 post-processing workflow combined results acquisition, cleaning and analyzing, and  
18 sharing findings to expose the relationships contained within the data. Visual observation  
19 of model results in the form anticipated for the final publication helped in understanding  
20 how changes to both the data pre-processing steps and the model execution steps  
21 impacted key model results. It is common that stakeholders will have unique interests for  
22 each model application, so creating general visualization tools will not always be  
23 possible. Therefore, creating workflows that leverage lower-level micro-services to

visualize model results in customized ways is a powerful tool provided by iRODS. This process of creating visualizations is another time saving strategy that allows the modeler to focus on model specific tasks rather than technique approaches for transferring data between the model and a visualization system.

Questions remain regarding the level of reuse that will be practical when micro-services and rules are used across hydrologic simulation models. There are classes of hydrologic models that can be grouped based on the use cases considered when developing the model. VIC falls into the “macro-scale” class of hydrologic models, meaning it is typically applied to regional, continental, or even global scale hydrologic systems. Other hydrologic models focus on more local scale systems such as a single catchment. Different classes of hydrologic models will likely require different schemes for pre-processing tasks such as discretizing the landscape, and will likely make use of different raw datasets to setup and parameterize the model. A key challenge moving forward, therefore, will be to ensure a flexible workflow environment where new data access and transformation tools specific to certain models can be easily developed and shared within focused hydrologic communities.

## **ACKNOWLEDGMENTS**

The authors wish to acknowledge the National Science Foundation (NSF) under the project DataNet Full Proposal: DataNet Federation Consortium (Award Number:094084.

## REFERENCES

- Abdulla, F. a., Lettenmaier, D.P., 1997a. Application of regional parameter estimation schemes to simulate the water balance of a large continental river. *J. Hydrol.* 197, 258–285. doi:10.1016/S0022-1694(96)03263-5
- Abdulla, F. a., Lettenmaier, D.P., 1997b. Development of regional parameter estimation equations for a macroscale hydrologic model. *J. Hydrol.* 197, 230–257. doi:10.1016/S0022-1694(96)03262-3
- Abdulla, F. a., Lettenmaier, D.P., Wood, E.F., Smith, J. a., 1996. Application of a macroscale hydrologic model to estimate the water balance of the Arkansas-Red River Basin. *J. Geophys. Res.* 101, 7449. doi:10.1029/95JD02416
- Badr, A. W., Wachob, A., Gellici, J.A., 2004. South Carolina, South Carolina Water Plan. Second Edition. South Carolina Department of Natural Resources. Land, Water and Conservation Division, Columbia, SC.
- Billah, M.M., Goodall, J.L., 2011. Annual and interannual variations in terrestrial water storage during and following a period of drought in South Carolina, USA. *J. Hydrol.* 409, 472–482. doi:10.1016/j.jhydrol.2011.08.045
- Billah, M.M., Goodall, J.L., Narayan, U., Reager, J.T., Lakshmi, V., Famiglietti, J.S., 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
- Crow, W.T., 2003. Multiobjective calibration of land surface model evapotranspiration predictions using streamflow observations and spaceborne surface radiometric temperature retrievals. *J. Geophys. Res.* doi:10.1029/2002JD003292

1 Fitch, P., Perraud, J.-M., Cuddy, S., Seaton, S., Bai, Q., Hehir, D., Sims, J., Merrin, L., Ackland, R.,  
2 Herron, N., 2011. The Hydrologists Workbench: more than a scientific workflow tool, in:  
3 Proceedings, Water Information Research and Development Alliance Science Symposium.

4 Furnans, J., Olivera, F., 2001. Watershed Topology - The Pfafstetter System. Esri User Conf. Vol. 21.

5 Gil, Y., Deelman, E., Ellisman, M., Fahringer, T., Fox, G., Gannon, D., Goble, C., Livny, M., Moreau, L.,  
6 Myers, J., 2007. Examining the Challenges of Scientific Workflows. Computer (Long. Beach. Calif).  
7 40, 24–32. doi:10.1109/MC.2007.421

8 Goff, S.A., Vaughn, M., McKay, S., Lyons, E., Stapleton, A.E., Gessler, D., Matasci, N., Wang, L.,  
9 Hanlon, M., Lenards, A., Muir, A., Merchant, N., Lowry, S., Mock, S., Helmke, M., Kubach, A.,  
10 Narro, M., Hopkins, N., Micklos, D., Hilgert, U., Gonzales, M., Jordan, C., Skidmore, E., Dooley, R.,  
11 Cazes, J., McLay, R., Lu, Z., Pasternak, S., Koesterke, L., Piel, W.H., Grene, R., Noutsos, C.,  
12 Gendler, K., Feng, X., Tang, C., Lent, M., Kim, S.-J., Kvilekval, K., Manjunath, B.S., Tannen, V.,  
13 Stamatakis, A., Sanderson, M., Welch, S.M., Cranston, K.A., Soltis, P., Soltis, D., O’Meara, B., Ane,  
14 C., Brutnell, T., Kleibenstein, D.J., White, J.W., Leebens-Mack, J., Donoghue, M.J., Spalding, E.P.,  
15 Vision, T.J., Myers, C.R., Lowenthal, D., Enquist, B.J., Boyle, B., Akoglu, A., Andrews, G., Ram, S.,  
16 Ware, D., Stein, L., Stanzione, D., 2011. The iPlant Collaborative: Cyberinfrastructure for Plant  
17 Biology. Front. Plant Sci. 2, 34. doi:10.3389/fpls.2011.00034

18 Guru, S.M., Kearney, M., Fitch, P., Peters, C., 2009. Challenges in using scientific workflow tools in the  
19 hydrology domain, in: 18th World IMACS Congress and MODSIM09 International Congress on  
20 Modelling and Simulation. Cairns, Qld., pp. 3514–3520.

21 Hedges, M., Blanke, T., Hasan, A., 2009. Rule-based curation and preservation of data: A data grid  
22 approach using iRODS. Futur. Gener. Comput. Syst. 25, 446–452. doi:10.1016/j.future.2008.10.003

1 Hedges, M., Hasan, A., Blanke, T., 2007. Management and preservation of research data with iRODS.  
2 Proc. ACM first Work. CyberInfrastructure Inf. Manag. eScience CIMS 07 17–22.  
3 doi:10.1145/1317353.1317358

4 Horsburgh, J.S., Tarboton, D.G., Hooper, R.P., Zaslavsky, I., 2014. Managing a community shared  
5 vocabulary for hydrologic observations. *Environ. Model. Softw.* 52, 62–73.  
6 doi:10.1016/j.envsoft.2013.10.012

7 Lakshmi, V., Piechota, T., Narayan, U., Tang, C., 2004. Soil moisture as an indicator of weather extremes.  
8 *Geophys. Res. Lett.* 31, 2–5. doi:10.1029/2004GL019930

9 Leonard, L., Duffy, C.J., 2013. Essential Terrestrial Variable data workflows for distributed water  
10 resources modeling. *Environ. Model. Softw.* 50, 85–96. doi:10.1016/j.envsoft.2013.09.003

11 Leonard, L., Duffy, C.J., 2014. Automating data-model workflows at a level 12 HUC scale: Watershed  
12 modeling in a distributed computing environment. *Environ. Model. Softw.* 61, 174–190.  
13 doi:10.1016/j.envsoft.2014.07.015

14 Liang, X., Lettenmaier, D.P., 1994. A simple hydrologically based model of land surface water and energy  
15 fluxes for general circulation models. *J. Geophys. Res.* 99, 14,415–14,428. doi:10.1029/94JD00483

16 Liang, X., Lettenmaier, D.P., Wood, E.F., 1996a. One-dimensional statistical dynamic representation of  
17 subgrid spatial variability of precipitation in the two-layer variable infiltration capacity model. *J.*  
18 *Geophys. Res.* 101, 21403. doi:10.1029/96JD01448

19 Liang, X., Wood, E.F., Lettenmaier, D.P., 1996b. Surface soil moisture parameterization of the VIC-2L  
20 model: Evaluation and modification. *Glob. Planet. Change* 13, 195–206. doi:10.1016/0921-  
21 8181(95)00046-1

1 Lohmann, D., Raschke, E., Nijssen, B., Lettenmaier, D.P., 1998. Regional scale hydrology: I. Formulation  
2 of the VIC-2L model coupled to a routing model. *Hydrol. Sci. J.* 43, 131–141.  
3 doi:10.1080/02626669809492107

4 Ludäscher, B., Altintas, I., Berkley, C., Higgins, D., Jaeger, E., Jones, M., Lee, E.A., Tao, J., Zhao, Y.,  
5 2006. Scientific workflow management and the Kepler system. *Concurr. Comput. Pract. Exp.* 18,  
6 1039–1065. doi:10.1002/cpe.994

7 Miles, B., Band, L., 2013. EcohydroWorkflowLib. <[http://pythonhosted.](http://pythonhosted.org/ecohydroworkflowlib/index.html)  
8 [org/ecohydroworkflowlib/index.html](http://pythonhosted.org/ecohydroworkflowlib/index.html)> (verified 05.29.2013).

9 Moriasi, D.N., Arnold, J.G., Liew, M.W. Van, Bingner, R.L., Harmel, R.D., Veith, T.L., 2007. *M e g s q a*  
10 *w s* 50, 885–900.

11 Oinn, T., Greenwood, M., Addis, M., Alpdemir, M.N., Ferris, J., Glover, K., Goble, C., Goderis, A., Hull,  
12 D., Marvin, D., Li, P., Lord, P., Pocock, M.R., Senger, M., Stevens, R., Wipat, A., Wroe, C., 2006.  
13 Taverna: lessons in creating a workflow environment for the life sciences. *Concurr. Comput. Pract.*  
14 *Exp.* 18, 1067–1100. doi:10.1002/cpe.993

15 Perraud, J., Fitch, P.G., Bai, Q., 2010. Challenges and Solutions in Implementing Hydrological Models  
16 within Scientific Workflow Software. *AGU Fall Meet. Abstr.* -1, 06.

17 Piasecki, M., Lu, B., 2010. Development of a Hydrologic Modeling Platform Using a Workflow Engine, in:  
18 *AGU Fall Meeting Abstracts.* p. 1239.

19 Rajasekar, A., Moore, R., Hou, C.-Y., Lee, C. a., Marciano, R., de Torcy, A., Wan, M., Schroeder, W.,  
20 Chen, S.-Y., Gilbert, L., Tooby, P., Zhu, B., 2010a. iRODS Primer: Integrated Rule-Oriented Data  
21 System, Synthesis Lectures on Information Concepts, Retrieval, and Services.  
22 doi:10.2200/S00233ED1V01Y200912ICR012

- 1 Rajasekar, A., Moore, R., Wan, M., Schroeder, W., 2010b. Policy-based Distributed Data Management  
2 Systems. JODI J. Digit. Inf. 11, 1–16.
- 3 Sheffield, J., Goteti, G., Wen, F., Wood, E.F., 2004. A simulated soil moisture based drought analysis for  
4 the United States. J. Geophys. Res. D Atmos. 109, 1–19. doi:10.1029/2004JD005182
- 5 Sheffield, J., Wood, E.F., 2007. Characteristics of global and regional drought, 1950–2000: Analysis of soil  
6 moisture data from off-line simulation of the terrestrial hydrologic cycle. J. Geophys. Res. Atmos.  
7 112, 1–21. doi:10.1029/2006JD008288
- 8 Troy, T.J., Wood, E.F., Sheffield, J., 2008. An efficient calibration method for continental-scale land  
9 surface modeling. Water Resour. Res. 44, 1–13. doi:10.1029/2007WR006513

10