1 **Integrating Scientific Cyberinfrastructures to Improve Reproducibility in**

2 **Computational Hydrology: Example for HydroShare and GeoTrust**

3 Bakinam T. Essawy [a], Jonathan L. Goodall [a*], Wesley Zell[b], Daniel Voce[c], Mohamed M. Morsy[a,d],

4 Jeffrey Sadler[a], Zhihao Yuan[e], and Tanu Malik[e]

5 [a] Department of Civil and Environmental Engineering, University of Virginia, 351 McCormick

6 Road, PO Box 400742, Charlottesville, VA, 22904, USA

7 [b] Earth Systems Modeling Branch, US Geological Survey, 12201 Sunrise Valley Dr., Reston, VA,

8 USA

9 [c] Department of Electrical and Computer Engineering, University of Virginia, 351 McCormick

10 Road, PO Box 400743, Charlottesville, VA, 22904, USA

11 [d] Irrigation and Hydraulics Engineering Department, Faculty of Engineering, Cairo University,

12 P.O. Box 12211, Giza 12613, Egypt

13 [e] College of Computing and Digital Media, DePaul University, Chicago, IL 60604, USA

14 [*] To whom correspondence should be addressed (E-mail: *goodall@virginia.edu*; Address:

15 University of Virginia, Department of Civil and Environmental Engineering, PO Box 400742,

16 Charlottesville, Virginia 22904; Tel: (434) 243-5019)

17

18 *Highlights:*

19 • Method for packaging and publishing scientific workflows

20 • Integration between GeoTrust and HydroShare projects

21 • GeoTrust is used to easily package environmental models as containers

22 • HydroShare is used to document and share packaged workflows

23 • An example application is provided for using a MODFLOW-NWT model

**Abstract**

24

25      The reproducibility of computational environmental models is an important challenge that

26    calls for open and reusable code and data, well-documented workflows, and controlled

27    environments that allow others to verify published findings. This requires an ability to document

28    and share raw datasets, data preprocessing scripts, model inputs, outputs, and the specific model

29    code with all associated dependencies. HydroShare and GeoTrust, two scientific

30    cyberinfrastructures under development, can be used to improve reproducibility in computational

31    hydrology. HydroShare is a web-based system for sharing hydrologic data and models as digital

32    resources including detailed, hydrologic-specific resource metadata. GeoTrust provides tools for

33    scientists to efficiently reproduce and share geoscience applications. This paper outlines a use case

34    example, which focuses on a workflow that uses the MODFLOW model, to demonstrate how the

35    cyberinfrastructures HydroShare and GeoTrust can be integrated in a way that easily and

36    efficiently reproduces computational workflows.

37    **Keywords:**

38     Computational reproducibility; hydrologic modeling; MODFLOW; metadata

39

40 **1. Software availability**

41 The software created in this research is free and open source. The software information and

42 availability are as follows:

43 Developers: Bakinam T. Essawy, Daniel Voce, and Wesley Zell

44 Programming language: Python, Bash

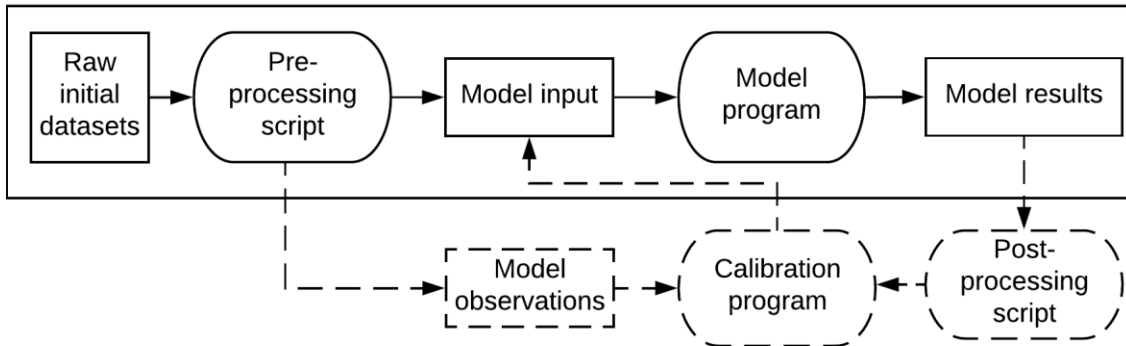45 GitHub link: https://github.com/uva-hydroinformatics-lab/AWS_MODFLOW.

46

## 2. Introduction

The challenge of creating more open and reusable code, data, and formal workflows that allow others to verify published findings is gaining attention in the scientific community (Borgman, 2012; David et al., 2016; Gorgolewski and Poldrack, 2016; Meng et al., 2015; Peng, 2011; Qin et al., 2016). Reproducibility is important for both verifying previous results as well as building upon the prior computational research of other scientists. Although we can achieve standard reproducibility for most computational research, there are certain cases in which reproducibility remains difficult to achieve. This challenge is not caused only by technical barriers but also by limited documentation of the research to be replicated and the potentially complex requirements for how the software is packaged, installed, and executed (Piccolo and Frampton, 2016). Recent papers have argued the need and have proposed approaches to improve reproducibility, both within geosciences generally and the hydrologic sciences specifically (David et al., 2016; Essawy et al., 2016; Gil et al., 2016; Hutton et al., 2016). Reproducibility of research is said to be achieved if the scientist was able to preserve sufficient computational artifacts in a way that can be replicated in the future (Meng et al., 2015).

Here we consider reproducibility to be the ability to repeat in the same exact form and then document and share digital resources previously used to complete an analysis. These digital resources include (1) initial raw, unprocessed datasets; (2) data preprocessing scripts used to clean and organize the data; (3) model inputs; (4) model results; and (5) the specific model code along with all of its dependencies. Figure 1 shows a typical conceptual workflow that needs to be repeated for computational reproducibility. These data, software, and environments are often integrated into workflows (as computational experiments) that allow scientists to re-run an analysis from raw initial datasets and obtain the same model results.

4

70       There are different requirements for reproducibility depending on the nature of the

71    research. For example, laboratory experiments require capturing descriptive information about

72    protocols and methods, leading to empirical reproducibility. Computational reproducibility, on the

73    other hand, requires descriptive information about the software and workflow details of model-

74    based research (Todden, 2013). Any workflow that is computationally reproducible must be

75    general and able to address the heterogeneous landscape of tools and approaches used within the

76    target scientific community. In hydrology, scientists use a large variety of computational models,

77    many of which have decades of development effort behind them (Singh et al., 2002).

78    Computational modeling can often require a significant amount of effort and time to prepare model

79    inputs and to calibrate and validate model parameters. Depending on the complexity of the system

80    being modeled and the experience of the modeler, these aspects can make reproducing

81    computational hydrologic experiments particularly challenging.

82    Addressing the challenges for achieving reproducibility in computational workflow has been

83    the topic of many studies. Until now, most approaches have either focused on the logical

84    preservation (i.e., sufficient documentation of a workflow and its components to allow for

85    reproduction later on) or physical preservation (i.e., workflow conservation by packaging all of its

86    components allowing identical replication) (Santana-Perez et al., 2017). It is hard to achieve a high

87    level of reproducibility while using one of these approaches in isolation; rather, the integration of

88    both physical and logical preservation is required to achieve a high level of reproducibility. Some

89    efforts have been made to integrate both logical and physical preservation for computational

90    workflows, such as the Topology and Orchestration Specification for Cloud Applications

91    (TOSCA). The TOSCA framework supports documentation for both the top-level structure of the

92    abstract workflow and the execution environment details (logical). TOSCA also provides

93    packaging functionality for the workflow (physical) (Qasha et al., 2016). In a similar way, our

94    approach provides both logical and physical preservation. However, the functionality is extended

95    to allow for automated creation, documentation, publication, and cloud-based execution of

96    scientific workflow packages.



97

98          **Figure 1** A typical conceptual workflow that needs to be repeated for computational

99      reproducibility. Dashed lines indicate processes for model calibration that are not discussed in

100                                            this study.

101          This research presents a solution for achieving a higher level of reproducibility by using

102    GeoTrust's *Sciunit-CLI* tool and HydroShare. HydroShare (http://www.hydroshare.org) and

103    GeoTrust (http://geotrusthub.org) are two new cyberinfrastructures under active development that

104    aim to improve reproducibility in computational hydrology. The methods described in this  paper

105    can be used to assist scientists to more easily repeat, reproduce, and verify a computational

106    experiment (Malik, 2017). This method goes beyond open source and simply shared by allowing

107    portability in different hardware and software environments and reproducible analyses with

108    different datasets. This level of reproducibility is not easily achieved by using HydroShare or

109    GeoTrust in isolation.  For example, GeoTrust does not provide a community of users who can

110    verify analyses or the variety of datasets that are required for verification; HydroShare, however,

111    does provide these. Similarly, while HydroShare simplifies the process of sharing code, data, and

112    descriptive metadata, it does not address the challenge of sharing the computational environment

113    required for the workflow and then repeating the computational workflow with different datasets.

114    This paper presents the design and implementation of a workflow that takes advantage of the

115    complementary strengths of the two systems. HydroShare is used to share key digital resources in

116    the workflow, while GeoTrust is used to capture, encapsulate, and make portable model execution.

117    An example application of the approach is presented using MODFLOW-NWT, a version of the

118    United States Geological Survey's groundwater model, MODFLOW (Niswonger et al., 2011).

119        The remainder of the paper is organized as follows. First, additional background on the

120    HydroShare and GeoTrust projects is provided. This background section is meant to orient readers

121    on key aspects of these projects. Next, the methodology section shows the system design and the

122    use case application for the MODFLOW-NWT model.   In the results section, the system

123    implementation of the HydroShare and GeoTrust integration approach is presented and

124    demonstrated by using the use case results as an example application. Finally, a discussion and

125    conclusions section summarizes the key aspects of the approach and outlines opportunities for

126    future research to advance on known limitations of the approach.

127    **3.  Background**

128    *3.1. HydroShare*

129        HydroShare is an open source web-based system developed for hydrologic scientists to

130    easily share, collaborate around, and publish all types of scientific data and models including

131    detailed, hydrologic-specific resource metadata (Tarboton et al., 2014a, 2014b). HydroShare has

132    been developed with the support of the United States National Science Foundation (NSF).

133    Following the completion of the original NSF grant, the Consortium of Universities for the

7

134    Advancement of Hydrologic Sciences Incorporated (CUAHSI) (also funded by the NSF) assumed

135    long-term support for HydroShare's operation and maintenance. In HydroShare, digital content is

136    stored and referred to as a "resource." Each resource is a unit used for management and access

137    control within HydroShare. Every resource has a resource type (Horsburgh et al., 2015).

138    HydroShare assigns a unique identifier for each newly created resource; this identifier is known as

139    the Resource ID. The "generic" resource type supports the Dublin Core metadata standard (Weibel

140    et al., 1998) and more specific resource types expand on this metadata standard for well-defined

141    data types. For example, "Model Operating System" is one of the extended metadata terms for the

142    "Model Program" resource type, which is used for sharing a computational model programs in

143    HydroShare (Morsy et al., 2017).

144         HydroShare provides a Representational State Transfer (REST) Application Program

145    Interface (API) that allows third-party applications to interact with HydroShare resources.

146    (https://github.com/hydroshare/hydroshare/wiki/HydroShare-REST-API#design-document).

147    Developers can create web-apps that use HydroShare's REST API to interact with HydroShare

148    resources. Web-app developers can catalogue their apps in HydroShare via the "Web-app"

149    resource type (Swain et al., 2016). When a developer creates a web-app resource in HydroShare,

150    the developer specifies which resource types are relevant to the web-app and the URL that will be

151    called when the web-app is executed from the landing page of the resource that the web-app is

152    acting on. After a developer adds a web-app as a resource in HydroShare, HydroShare users can

153    execute that app through HydroShare's web interface to act on relevant resources that they have

154    access to.

155         Although there are several different resource types supported by HydroShare, two of the main

156    resource types relevant to this paper deal with computational models. HydroShare divides

157    computational models into two separate but linked resource types: a) the model program and b)

158    the model instance. The model program includes the software for executing a specific instance of

159    the model and the model instance are the input files required for executing the model and,

160    optionally, the output files after a model instance has been executed by a model program

161    (Horsburgh et al., 2015; Morsy et al., 2017, 2014). Additionally, a Model Instance Resource type

162    can be linked to a model program resource type using the "ExecutedBy" term, assisting with

163    reproducibility of the model instance (Morsy et al., 2017). Other HydroShare resource types used

164    in this paper include the Composite resource type, which allows uploading metadata files at both

165    file and resource level; the collections resource type, which stores any number of individual

166    resources within HydroShare as a single, aggregate resource; and the web-app resource type, which

167    is the Digital content stored in HydroShare and referred to it as a "resource."

168    *3.2. GeoTrust*

169    The GeoTrust project, also funded by the NSF through their EarthCube program, aims to

170    create cyberinfrastructure that assists scientists to efficiently reproduce and share geoscience

171    applications used in research (Malik et al., 2017). The project has done this primarily by

172    developing the concept of a "sciunit" (https://sciunit.run/), an efficient, lightweight, self-contained

173    digital package of an ad-hoc computational workflow that can be repeated in other environments.

174    The sciunit advances the concept of a research object, an aggregation of digital artifacts such as

175    code, data, scripts, and temporary experiment results associated with a research paper. The sciunit

176    provides an authoritative and far more complete record of a piece of research (Hai et al., 2017).

177    To create, maintain, and publish sciunits, the GeoTrust project has developed a software tool for

178    Linux environments called *Sciunit-CLI*.

179    One of the main advantages of a sciunit is its portability, which allows it to be easily run on

180    various computing environments. To accomplish this, *Sciunit-CLI* creates sciunits using Docker,

181    a widely used containerization software. Docker wraps a piece of software in a complete filesystem

182    that contains everything needed to run the software, including code, software runtime, system

183    tools, and system libraries in a Docker container (Owsiak et al., 2017). By leveraging Docker,

184    sciunits are packaged with all of their dependencies. In this way, any sciunit can be executed in

185    any environment in which both Docker and the *Sciunit-CLI* tool are installed regardless of other

186    computer configurations (Hai et al., 2017). This capability eliminates the burden of configuring a

187    running environment with all software dependencies, which can be complex, in order to reuse a

188    scientific workflow and reproduce its results.

189        In addition to ensuring the portability of sciunits, *Sciunit-CLI* automates some documentation

190    of the workflow packaged into a sciunit, including environment dependencies. The automation of

191    documenting all code, data, and environment dependencies alleviates what is typically a

192    burdensome task for scientists. Importantly, *Sciunit-CLI* also records retrospective provenance of

193    the workflow execution, which can be used for re-running containers (Pham et al., 2014). Because

194    it contains all of the required dependencies, the sciunit can be rerun, and the outputs reproduced,

195    using any other deployment configuration that also has *Sciunit-CLI* installed. When *Sciunit-CLI*

196    creates a sciunit, it includes three types of metadata: annotation metadata (populated by the user)

197    and provenance and version metadata (generated automatically by *Sciunit-CLI*).

198        Figure 2 shows an example user interaction with the *Sciunit-CLI* tool. The user runs the

199    *create* command and provides a name, *"Model"* in the example. To create a container or a package

200    within the sciunit, the user runs the *package* command and provides the workflow name (e.g.,

201    "workflow.sh") along with any inputs for the workflow (e.g., "data"). The user application can be

202    written in any combination of programming languages, and many containers can be created within

203    the same sciunit.

204        *Sciunit-CLI* works in a distributed fashion, similar to the Git version control philosophy,

205    such that the sciunits are stored only locally until explicitly shared with a remote repository.  This

206    method of operation allows distributed collaborators to work offline on the same sciunit. When a

207    user is ready to share, they can publish the sciunit container to any remote web-repository using

208    the *publish* command. To use the publish command, the remote repository must be configured

209    within the *Sciunit-CLI* tool. This command line prompts first-time users to provide their remote

210    web-repository credentials. The remote repository reads the container's contents, stores the

211    container's digital artifacts in the appropriate remote sciunit, and associates the container with an

212    appropriate cloud execution server on which it can potentially re-execute. In our case, we used

213    HydroShare as the remote repository to publish our packaged sciunit in order to use HydroShare's

214    support for rich metadata and its ability to integrate third-party applications. The latter allowed us

215    to automate the cloud-based execution of this packaged sciunit.

```
1. > create  Model
2. > annotate Model author: Bakinam Essawy
3. > exec workflow.sh 1 /data
4. > show
     id: e1
     sciunit: Model
     command: ./workflow.sh Data
     size: 1.18 GB
     started: 2017-11-30 21:23
5. > push my_new_article --setup hs
6. >  repeat e1
7. > stop
```

216

217                    **Figure 2** A example user interaction with sciunit client.

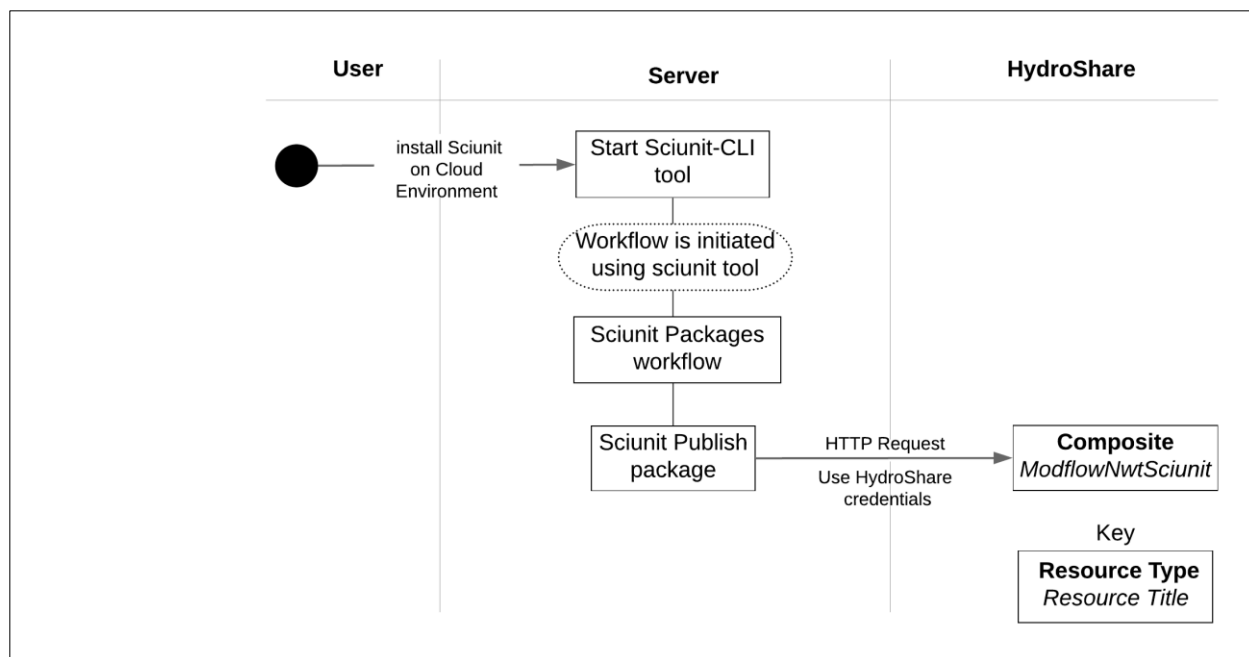218    **4.  Methodology**

219    *4.1. System Design*

220        The combined GeoTrust and Hydroshare system is designed to connect a repeatable

221    computational workflow with its input data in a reproducible way. As such, both the computational

11

222    workflow and the data must be stored in a public repository that has extensive metadata support.

223    In addition to public accessibility of the data and the computational workflow, the execution of the

224    workflow must also be made publicly available to ensure reproducibility and transparency.  The

225    technology for producing a repeatable computational workflow is provided by the GeoTrust

226    *Sciunit-CLI,* while the technology for public storage and metadata support is provided by

227    CUAHSI's HydroShare. Therefore, the main design aspect of this work consisted of designing a

228    publicly accessible method of execution in which sciunits built with the *Sciunit-CLI* and stored in

229    HydroShare could be executed using input data also stored in HydroShare. This was done in two

230    parts. The first was to build in functionality for publishing a sciunit through HydroShare. The

231    second part was to automate the execution of a sciunit from HydroShare using HydroShare web-

232    apps.

233    *4.1.1.  Integrating Sciunit-CLI with HydroShare*

234         Figure 3 shows an activity diagram of the system design for integrating GeoTrust *Sciunit-CLI*

235    and HydroShare. To achieve this integration, *Sciunit-CLI* was extended to support sharing of

236    sciunits through HydroShare. This functionality was implemented using HydroShare's REST API.

237    To publish their sciunit on HydroShare, the user must provide valid HydroShare credentials.  In

238    the current implementation, the sciunit resource is published on HydroShare as a Composite

239    Resource Type. Once the resource for the sciunit is created within HydroShare, the user can log

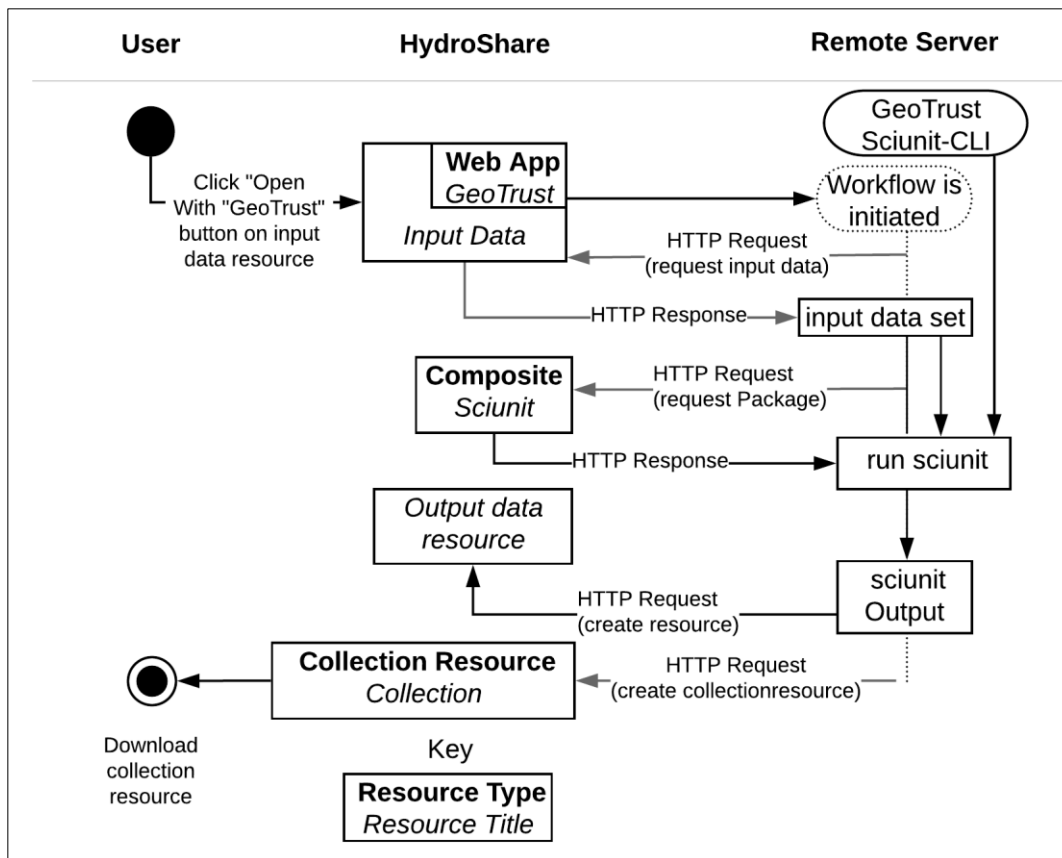240    into HydroShare and edit the metadata fields to more fully describe the sciunit resource.

**Figure 3** Activity diagram showing creating a sciunit using GeoTrust and publishing that sciunit on HydroShare.

*4.1.2. Automating sciunit execution through HydroShare*

Integrating the cloud-based sciunit execution from the HydroShare user interface was done using a HydroShare web-app. This web-app directs Hyper Text Transfer Protocol (HTTP) request to a web server where sciunits can be executed. The web-app configured to run a particular sciunit can be accessed through the "Open with" button on the landing page for the resource that stores the raw input data. When the scientist clicks on the web-app button from the "Open with" menu, an HTTP request containing the raw input data's resource ID will be sent to the server. With the resource ID, the HydroShare REST API can be used to download the raw input data and the sciunit to the server. The server can then execute the sciunit using the raw data, and return the output to the scientist as a new HydroShare resource.

Figure 4 shows the steps done in a generic form for the integration between the two cyberinfrastructures, GeoTrust and HydroShare, to improve reproducibility by automating the

13

256  execution of the published sciunit. The figure shows how the "Open with" app will perform a

257  HTTP GET request to a remote server, which has already been configured with the *Sciunit-CLI*.

258  This automation process is done using a Python script created on the web server machine. This

259  Python script uses the flask library to act as a web server with NGINX (https://www.nginx.com/)

260  used as a proxy to forward all HTTP requests from the user browser to the Python script, which

261  can handle multiple users simultaneously. The Python script is using the POST request to create a

262  new resource and upload the output generated from running the sciunit on this resource.

263  Simultaneously, a webserver is running on the remote machine, which handles the HTTP request

264  and automatically executes a Python script. This script uses the HydroShare user authentication to

265  download the input data from the resource and downloads the Composite resource that includes

266  the sciunit container. Once both resources are downloaded, the resources are unzipped and moved

267  to the working directory for the analysis. The *Sciunit-CLI* executes the downloaded sciunit

268  package.  After the sciunit is executed, a new resource is created in HydroShare and the output

269  from the *Sciunit-CLI* execution is uploaded into this new resource. A new collection resource is

270  also created on HydroShare to group all resources that were included during this execution.  In this

271  paper we used HydroShare API. Our Python script uses the Python Client Library for the REST

272  API (http://hs-restclient.readthedocs.io/en/latest/).

**Figure 4** The generic implementation for automating the execution of the published sciunit from the HydroShare web-app

*4.2. Use Case Application*

A use case application was designed to demonstrate the integration of GeoTrust *Sciunit-CLI* and HydroShare. This integration allows GeoTrust to package and publish a sciunit through HydroShare, after which HydroShare automates the execution of this sciunit. Execution of the packaged sciunit through HydroShare was demonstrated using EC2 instances from Amazon Web Services (AWS). A Linux-based, micro-sized machine (t2) was used for prototyping and demonstration purposes; this machine had 1 Gb of memory, 1 vCPU, 32 Gb of Solid State Drive (SSD)-based local instance storage, and a 64-bit platform ("Amazon EC2 Instances," 2015). This use case consisted of a workflow used for preprocessing model input data, running a computational

15

285    model, and handling the model outputs. The computational model used for the use case was

286    MODFLOW-NWT.

*4.2.1.  MODFLOW-NWT Use Case*

288    MODFLOW-NWT is a standalone version of MODFLOW, a commonly used groundwater

289    model (Niswonger et al., 2011). The concept of "packages" is key to the modularity of the different

290    versions of MODFLOW (including MODFLOW-NWT); packages are input files that define some

291    individual component of the groundwater-flow conceptual model or specify the solution method

292    used for the flow equation that is collectively formulated from the individual components.  For

293    example, the basic (BAS) and discretization (DIS) packages define the spatial and temporal

294    framework of the model, including the grid dimensions and the location of active and inactive grid

295    cells, while the recharge (RCH) package defines the spatial-distribution and rate of recharge to the

296    water-table.  For our use case using MODFLOW-NWT, the Newton-Raphson (NWT) package

297    defines the variables required to implement the Newton-Raphson solution method.

298    For this study, MODFLOW-NWT was used to simulate the shallow groundwater flow in the

299    James River watershed upstream of Richmond, VA, USA. The model includes recharge to the

300    water table, subsurface flow through the saturated zone, and base-flow discharge to surface water

301    bodies including the James, Rivanna, and Hardware Rivers and several smaller-order streams.

302    Depth-integrated effective transmissivity was assumed to be constant throughout the active model

303    area and spatially-distributed recharge was derived from the national recharge dataset developed

304    by Reitz et al. (2017). Base-flow discharge was simulated using the MODFLOW drain (DRN)

305    package with all drain elevations (i.e., the water-table elevation required to discharge base-flow to

306    a receiving stream) extracted from the National Elevation Dataset. The model runs to completion

307    and is unconstrained by calibration; as such it is to be only used as an example for the workflow

308    processes described in this paper (i.e., no hydrologic or management conclusions were drawn from

309    the results of the model). This workflow could be extended to include calibration (Figure 1). For

310    example, a HydroShare resource for a parameter estimation program such as PEST (Doherty and

311    Hunt, 2010) could be created and included in the sciunit container. Similarly, the pre-processing

312    script could include data retrieval from web services such as the USGS water services API

313    (https://waterservices.usgs.gov/) and the automated generation of PEST input files.

314    The FloPy library was used to create the MODFLOW-NWT model from raw input datasets

315    (Bakker et al., 2016). FloPy is a library of Python modules that allows scripting of the various

316    steps in MODFLOW model development, execution, and analysis. By combining FloPy with

317    GeoTrust and HydroShare, the workflow used to create and execute MODFLOW model (e.g., the

318    steps shown in Figure 1) can be stored within a reproducible container with descriptive metadata

319    in HydroShare.

320    **5. Results**

321    *5.1. System Implementation*

322    The system was implemented using the following steps. First, the script downloads raw input

323    data and the sciunit resources from HydroShare. Second, the script will unzip both the data and

324    sciunit, pass the data to the sciunit as an argument (this is how the sciunit accepts the input data),

325    and then run the sciunit with the downloaded data. Last, after the execution is completed, the

326    Python script will upload the results to HydroShare by using a POST request to create two new

327    resources: one for the sciunit output, which has the MODFLOW-NWT Model Instance Resource

328    type, and the other the collection resource that will include all the resources used within the study.

329    The script then returns the command status (including any errors) to the user.

17

330    *5.2. Use Case Results*

331    A digital workflow (bash script) was packaged into a sciunit using the *Sciunit-CLI* tool.

332    The digital workflow runs a Python script to prepare the MODFLOW-NWT input data files and

333    then executes a single run of the model. Figure 5 shows the component of the packaged digital

334    workflow.

```
#!/bin/bash
cp -a /home/$1/$1/data/contents/.  /home/Data/
(cd /home/; python build_modflow.py)
(cd /home/MODFLOW;  ./mfnwt *.nam)
```

335

336                    **Figure 5** component of the packaged digital workflow.

337    Figure 6 outlines the first steps taken in the process to start and create a new sciunit through

338    the GeoTrust *Sciunit-CLI* tool for the example workflow while Figure 7 shows the execution and

339    packaging of the digital workflow into a sciunit package. This package command traces all

340    dependencies for the workflow and includes them in a single Docker file. Figure 8 shows how the

341    *publish* command is used to publish a sciunit package on HydroShare. If this is the user's first time

342    connecting to HydroShare, *Sciunit-CLI* will ask for HydroShare user credentials, otherwise the

343    credentials stored will be used. Once the package is published, metadata can be provided by the

344    user via the HydroShare Graphical User Interface (GUI). Future implementations of the *Sciunit-*

345    *CLI* may expand this functionality by automatically populating more detailed metadata for

346    describing resources.

```
ubuntu@ip-172-31-25-113:~/test$sciunit create Model
Opened empty sciunit at /home/ubuntu/sciunit/Model
```

347

348    **Figure 6** The creation of a new sciunit through the GeoTrust *Sciunit-CLI* tool for the use case

349

```
ubuntu@ip-172-31-25-113:~/test$sciunit exec ./workflow.sh Data
Rasterizing shapefile: Data/James_Rivanna_5070.shp
Writing output raster to: Framework/James_Rivanna_IBOUND.tif
0...10...20...30...40...50...60...70...80...90...100 - done.
Executinggdalwarppath:gdalwarp
Clipping the raster to the model domain.
```

350 **Figure 7** Execution of the use case workflow through sciunit to create a package

351

```
ubuntu@ip-172-31-25-113:~/test$ sciunit push my_new_article --setup hs
Logged in as "Essawy, bakinam <btaessawy@gmail.com>"
Title for the new article: Model
my_new_article: 596MB [00:31, 18.8MB/s]
ubuntu@ip-172-31-25-113:~/test$
```

352 **Figure 8** Publishing the use case sciunit to HydroShare

353     The newly created resource on HydroShare is a Composite Resource Type. This resource

354 type allows the resource to include multiple files without file format limitations and with metadata

355 associated at a file level within the resource. The Composite resource contains two files. The first

356 is the provenance metadata file created while packaging the workflow; this metadata file contains

357 information concerning the creation and version history of the managed data. The second file is

358 the zipped package for the sciunit itself.

359     Once the sciunit is available as a HydroShare resource, HydroShare's integration with

360 third-party web apps is used to execute the sciunit.  In order to store data and make it accessible to

361 be used as the input required by the sciunit, we made a new model instance-type resource titled

362 "ModflowNwtRawData" (Essawy, 2018b). We also created a web-app resource titled "GeoTrust"

363 (Essawy, 2018a). This web-app pointed to the AWS-EC2 instance where the *Sciunit-CLI* tool and

364 our Python script were installed. The connection between the HydroShare resource and the web

365 server  was  made  by  providing  the  web  server's  URL  as  the  "App-launching  URL  Pattern"

366 metadata   term   in   the   resource.   The   GeoTrust   web-app   resource   is   linked   to   the

367 ModflowNwtRawData resource by the SupportedResourceType metadata property. This metadata
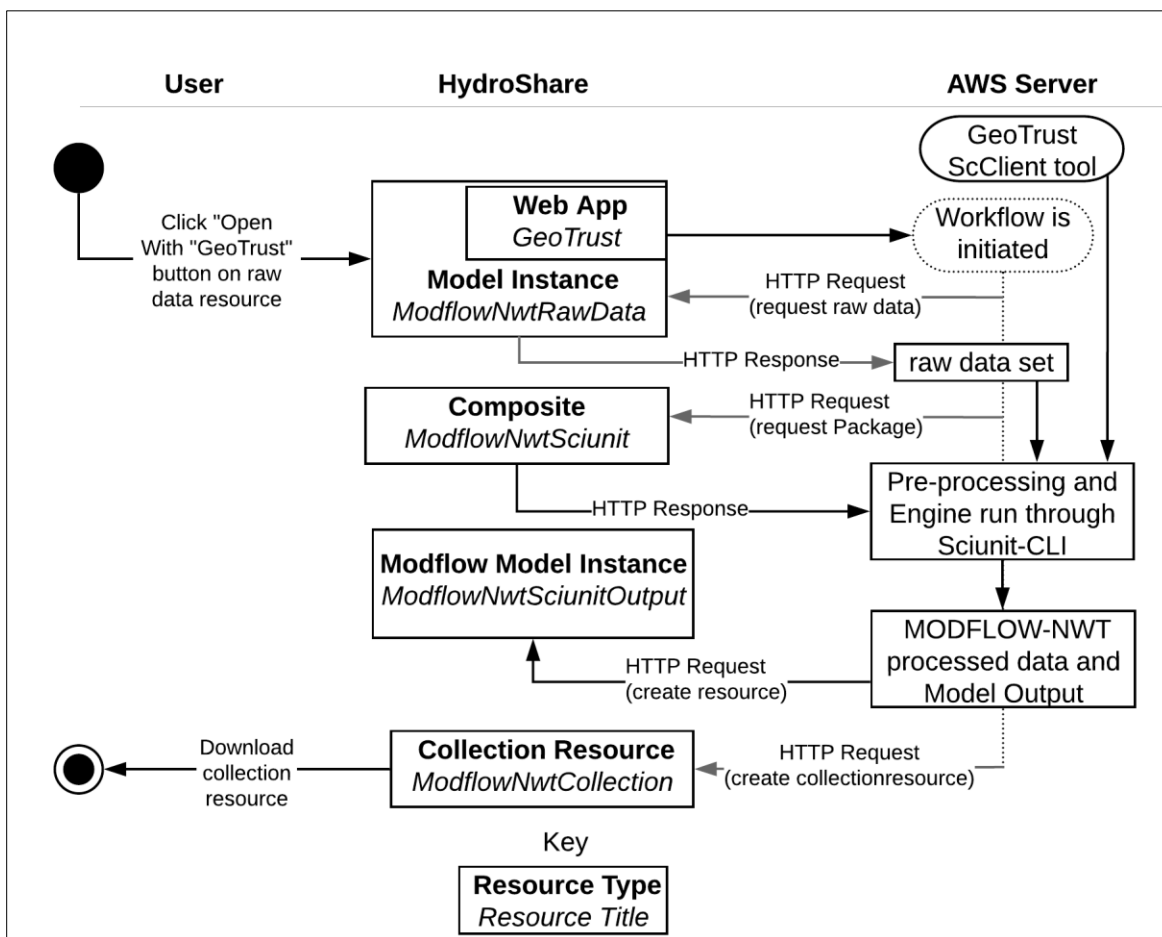
19

368  property was set to include the Composite Resource Type, which allowed the web-app to appear

369  in a drop-down list in the "Open with" menu on the ModflowNwtRawData resource landing page.

370  Figure 9 shows the Model Instance Resource type that includes the raw data, and the web apps

371  linked to this resource type to automate the sciunit execution. When the GeoTrust web-app on this

372  page is selected, the HTTP request is sent to server and the workflow is executed. The output is

373  written back to HydroShare as a new resource with the MODFLOW Model Instance Resource

374  type. This resource type is used because the resource can be executed by a MODFLOW model

375  program and it allows for adding extended metadata specific to MODFLOW (Morsy et al., 2017).



376

377  **Figure 9** The raw data within the Model Instance Resource type, and the web apps linked to this

378                    resource type to automate the sciunit execution.

379       Figure 10 presents the activity diagram for the steps that occur when the "Open with" button

380  is clicked and the "GeoTrust" app is selected on the ModflowNwtRawData resource landing page.

381  The "GeoTrust" app will perform an HTTP GET request to the AWS-EC2 machine, which has

382  already been configured with the *Sciunit-CLI*. The webserver running on the AWS-EC2 machine

383  handles the HTTP request and automatically executes a Python script. The script uses the

384    HydroShare user authentication to download both the raw data of the ModflowNwtRawData

385    resource and the sciunit container included within the ModflowNwtSciunit resource (Essawy,

386    2018c). Once the ModflowNwtSciunit and the ModflowNwtRawData resources are downloaded,

387    the script unzips the resources and moves them to the working directory for the analysis. The

388    *Sciunit-CLI* tool executes the downloaded sciunit package, which pre-processes the raw input data

389    for the model and executes the MODFLOW-NWT model. After the model is executed, a new

390    resource is created in HydroShare with the MODFLOW Model Instance Resource type named

391    ModflowNwtSciunitOutput (Essawy, 2018d) and the output from the *Sciunit-CLI* execution is

392    uploaded into this new resource. A new collection resource is also created on HydroShare to group

393    all the resources: the ModflowNwtRawData generic Model Instance Resource (the resource type

394    is a generic model instance because the data uploaded have no specific metadata or format that

395    could be tied to a specific resource type), the web-app GeoTrust resource, the ModflowNwtSciunit,

396    MODFLOW Model Instance Resource, the ModflowNwtSciunit Composite resource, and the

397    ModflowNwtSciunitOutput resource that includes the output resulting from executing the sciunit

398    package.

399

400

**Figure 10** Activity diagram showing the steps for the online execution of the sciunit through

HydroShare.

Figure 11 shows HydroShare user "My Resources page" after using the "Open with" action

button on the GeoTrust web-app on the ModflowNwtRawData resource for the online execution.

Two new resources are created. The first resource in the workflow is the

ModflowNwtSciunitOutput resource, which includes the input files for the MODFLOW-NWT

model program that are prepared through the preprocessing script and the output from the model

run. This resource is given the MODFLOW Model Instance Resource type, because the resource

has the inputs that are required by the MODFLOW-NWT model. This resource type allows for

extended metadata specific to a MODFLOW model instance. The second resource created is the

22

411    ModflowNwtCollection resource (Essawy, 2018e), which includes all the resources used in the

412    online execution for the MODFLOW-NWT. This provides a grouping of resources used for an

413    analysis and allows the user to share or download this collection of resources more easily.



414

415    **Figure 11** HydroShare user My Resources page after using the GeoTrust web app for the online

416    execution.

417        Figure 12 shows the output files within ModflowNwtSciunitOutput resource as viewed on

418    this resource's HydroShare landing page. The resource contains the output generated from running

419    the sciunit that prepares the model input for MODFLOW-NWT and the output from running the

420    MODFLOW-NWT model program itself. The MODFLOW Model Instance Resource type

421    includes extended metadata terms specific for MODFLOW.  In this use case the model has eight

422    packages.  In addition to the packages already described, this model instance includes: the output

423    control (OC) package, which specifies how the model output is written; the upstream-weighting

424    (UPW) groundwater flow package, which describes the system properties (e.g.,

425    transmissivity/conductivity); and the one output listing file (LIST), which contains all the

426    information about the current run (e.g., stress period, time step and the number of active and

427    inactive cells, the recharge, drains, and any errors). The name file (NAM) specifies the name of

428    the input and output files for the model instance.

429



430

431    **Figure 12** The output files within the ModflowNwtSciunitOutput resource landing page in

432               HydroShare.

433          Additional metadata associated with the MODFLOW output resource is divided into four

434   categories: 1) Authorship, 2) Related resources, 3) Resource Specific, and 4) Web Apps. Figure

435   13 shows the "Related Resources" metadata. Here all resources linked to the MODFLOW output

436   resource through formal relationships are listed. In this case, the MODFLOW output resource is

437   linked to the ModflowNwtRawData resource through the "Derived From" relationship and to the

438   MODFLOW-NWT resource through the "isExecutedBy" relationship. Figure 14 shows the

439   "Resource Specific" metadata. These are non-null metadata terms that apply only to the

440   MODFLOW Model Instances' such as grid attributes, solver, and boundary condition package

441   choices. Additional metadata terms not previously populated by the user can be populated later

442   within the edit mode and will appear in this section once populated.

443



444   **Figure 13** The ModflowNwtSciunitOutput Related Resources metadata tracking the resource's

445   provenance within HydroShare.

**Figure 14** ModflowNwtSciunitOutput specific metadata capturing key MODFLOW model properties.

Figure 15 shows details for the resulting ModflowNwtCollection resource as viewed on this resource's landing page. The collection resource contains four sub-resources: 1) the ModflowNwtRawData resource with the raw input data ready to be prepared for the MODFLOW-NWT model engine; 2) the ModflowNwtSciunit resource with the sciunit pre-processing workflow, which also includes running the MODFLOW-NWT model; 3) the ModflowNwtSciunitOutput resource, which stores the output generated from running the sciunit workflow; and 4) the GeoTrust web app used to perform the online model execution using AWS-EC2. By organizing all these resources into a single collection, it is possible to have one landing page where users can, referring back to the stated goals in the introduction of this paper, view, obtain, and execute (1) raw initial datasets, (2) data preprocessing scripts used to clean and

459    organize the data, (3) model inputs, (4) model results, and (5) the specific model code along with

460    of all its dependencies used for a computational analysis.

461

**Figure 15** The collection resource that includes all resources used within the study.

**6. Discussion and Conclusions**

In this paper, we demonstrated how HydroShare and GeoTrust can be integrated to easily and efficiently package, share, and publish model workflows. MODLFOW-NWT was used as an example application to demonstrate the functionality provided by these cyberinfrastructures for creating open, reusable data analysis and cloud-based model execution services. The approach showed how containers built using GeoTrust tools can be shared as HydroShare resources. A cloud-based service was created to automatically retrieve raw input data from HydroShare, execute a sciunit container that both prepares and runs a MODFLOW-NWT model, and share the results on HydroShare using a MODFLOW Model Instance Resource type. All the resources are aggregated in HydroShare into one collection resource with domain-specific metadata.

The integration of scientific cyberinfrastructures such as the HydroShare and GeoTrust projects can improve reproducibility in computational hydrology. New MODFLOW models can be directly built from unprocessed input data (e.g., land-surface DEMs or stream-network shapefiles) by running a sciunit container that includes automated data preparation steps implemented using the FloPy Python package. The container is run online using AWS resources initiated directly through the HydroShare user interface. A particular advantage of this approach is that the GeoTrust *Sciunit-CLI* tool provides scientists a method for efficiently creating containers for script-driven modeling workflows. Thus, the general approach demonstrated here for the MODFLOW-NWT use case could be applied for any workflow that can be automated and that is compatible with Docker requirements. For example, in prior work we have constructed pre- and post-processing workflows for the Variable Infiltration Capacity (VIC) hydrologic model (Liang et al., 1996) that could directly benefit from this method for packaging, sharing, and publishing resources (Billah et al., 2016; Essawy et al., 2016). These containers are efficient, lightweight,

486   self-contained packages of computational experiments that can be repeated or reproduced

487   regardless of deployment configurations.

488       In addition to integration with HydroShare for storing and publishing a sciunit, cloud resources

489   were used to execute sciunits directly through the HydroShare user interface. While only AWS

490   was presented, we evaluated as part of this work three different cloud computing services:

491   EarthCube Integration and Testing Environment (ECITE), CyVerse, and Amazon Web Services

492   (AWS). ECITE and CyVerse are funded by NSF and both are under active development. One main

493   advantage for using ECITE or CyVerse is that they are free of charge for scientific studies. AWS,

494   though not free, does offer a competitive grant program for researchers. From our experience, the

495   AWS platform made the process of obtaining computer resources the simplest when compared to

496   ECITE and CyVerse. The AWS user simply logs in to the console, selects the type of the machine

497   needed, and launches it. When using ECITE, we had to contact the developer and ask for an

498   instance with the required specifications and a short paragraph summarizing the project we are

499   working on to justify the allocation of compute resources. We also needed to contact the developer

500   each time we wanted to open a port (e.g., port 22 to SSH or port 80 for HTTP). The service did

501   not support Elastic IPs like AWS, so each time we restarted an instance and wanted to use SSH to

502   access to the machine, we needed to report the IP address used to access the machine to the

503   developer to add this address to the security rules. CyVerse is a more mature service, but allows

504   each user only a certain allocation of computational time. Once the user exceeds this allocation the

505   instance is suspended and the user needs to request more time from the administrators. This feature

506   was problematic for our use case of a continually available cloud-based resource for online model

507   execution. For these reasons, we used AWS-EC2 for much of the testing work described in this

508    paper, but ECITE and CyVerse are in active development and will likely be good options for this

509    use case in the future.

510        While this approach shows great promise, it is not without limitations: (1) the *Sciunit-CLI* tool

511    must be installed in order to re-execute a sciunit container and (2) HydroShare lacks methods for

512    uniquely identifying and managing web-app resources that will be needed as the number of these

513    resources continues to increase. Regarding the latter limitation, without a more organized structure,

514    naming conflicts could cause confusion when using the "Open with" button over which app is to

515    be requested. Also, this work does not fully explore computational challenges associated with the

516    proposed methodology. Using cloud services like AWS provides the opportunity for scalability as

517    more users are added. For example, this solution used small EC2 instances for prototyping. Future

518    work could explore AWS EC2 Container Service (ECS) as an alternative for a more scalable

519    solution to support multiple concurrent users. Data movement between HydroShare and AWS is

520    another potential issue as data volumes increase, which is not uncommon for hydrologic modeling.

521    HydroShare is built on iRODS (Integrated Rule-Oriented Data System), which includes the ability

522    to interface with AWS S3 storage resources. Future work could explore using this functionality to

523    automate the movement of large files between HydroShare and AWS to support computation

524    within AWS and still maintain access through the HydroShare user interface. iRODS is

525    specifically designed to handle such data federation needs and should provide a robust solution for

526    managing the large data flows common in hydrologic modeling.  Lastly, future work should

527    explore scaling of the general approach presented here to use cases in which multiple sciunits are

528    available for execution within a remote, cloud-based resource. In this case, a user could select from

529    available sciunits to process input data stored with HydroShare, making for a potentially very

530 powerful general approach applicable to many different modeling and analysis use cases that

531 require remote data processing.

**7. Acknowledgements and Disclaimer**

**8. References**

537 Amazon EC2 Instances [WWW Document], 2015. URL http://aws.amazon.com/ec2/instance-

538     types/ (accessed 6.7.15).

539 Bakker, M., Post, V., Langevin, C.D., Hughes, J.D., White, J.T., Starn, J.J., Fienen, M.N., 2016.

540     Scripting MODFLOW Model Development Using Python and FloPy. Groundwater 54, 733–

541     739. https://doi.org/10.1111/gwat.12413

542 Billah, M.M., Goodall, J.L., Narayan, U., Essawy, B.T., Lakshmi, V., Rajasekar, A., Moore, R.W.,

543     2016. Using a data grid to automate data preparation pipelines required for regional-scale

544     hydrologic     modeling.     Environ.     Model.     Softw.     78,     31–39.

545     https://doi.org/10.1016/j.envsoft.2015.12.010

546 Borgman, C.L., 2012. The conundrum of sharing research data. J. Am. Soc. Inf. Sci. Technol. 63,

547     1059–1078.

548 David, C.H., Famiglietti, J.S., Yang, Z.-L., Habets, F., Maidment, D.R., 2016. A decade of

549     RAPID—Reflections on the development of an open source geoscience code. Earth Sp. Sci.

550     226–244. https://doi.org/10.1002/2014EA000014.Received

551 Doherty, J.E., Hunt, R.J., 2010. Approaches to highly parameterized inversion-A guide to using

552     PEST for groundwater-model calibration. US Geol. Surv. Sci. Investig. Rep.

553 Essawy, B., 2018a. GeoTrust,. HydroShare,

554 http://www.hydroshare.org/resource/126701df868e4da9872d9b533db34ae6.

555 Essawy, B., 2018b. ModflowNwtRawData, HydroShare,

556 http://www.hydroshare.org/resource/4c9f9daa09e745a5b285481c7903c759.

557 Essawy, B., 2018c. ModflowNwtSciunit. HydroShare,

558 http://www.hydroshare.org/resource/995479b35b62486783e0da63e937ca89.

559 Essawy, B., 2018d. ModflowNwtSciunitOutput, HydroShare,

560 http://www.hydroshare.org/resource/19605cf6e91e415fb98b7a28cad263d6.

561 Essawy, B., 2018e. ModflowNwtCollection, HydroShare,

562 http://www.hydroshare.org/resource/bf598099ed384540aaa9284b7343a717.

563 Essawy, B.T., Goodall, J.L., Xu, H., Rajasekar, A., Myers, J.D., Kugler, T.A., Billah, M.M.,

564 Whitton, M.C., Moore, R.W., 2016. Server-side workflow execution using data grid

565 technology for reproducible analyses of data-intensive hydrologic systems. Earth Sp. Sci. 3,

566 163–175. https://doi.org/10.1002/2015EA000139

567 Gil, Y., David, C.H., Demir, I., Essawy, B.T., Fulweiler, R.W., Goodall, J.L., Karlstrom, L., Lee,

568 H., Mills, H.J., Oh, J.-H., Pierce, S.A., Pope, A., Tzeng, M.W., Villamizar, S.R., Yu, X.,

569 2016. Towards the Geoscience Paper of the Future : Best Practices for Documenting and

570 Sharing Research from Data to Software to Provenance. Earth Sp. Sci. 1–75.

571 https://doi.org/10.1002/2015EA000136

572 Gorgolewski, K.J., Poldrack, R.A., 2016. A Practical Guide for Improving Transparency and

573 Reproducibility in Neuroimaging Research. PLoS Biol. 14, 1–13.

574 https://doi.org/10.1371/journal.pbio.1002506

575 Hai, D., That, T., Fils, G., Yuan, Z., Malik, T., 2017. Sciunits : Reusable Research Objects. Tech.

576        Report, DBGroup, Sch. Comput. DePaul Univ.

577 Horsburgh, J.S., Morsy, M.M., Castronova, A.M., Goodall, J.L., Gan, T., Yi, H., Stealey, M.J.,

578        Tarboton, D.G., 2015. Hydroshare: Sharing diverse environmental data types and models as

579        social objects with application to the hydrology domain. JAWRA J. Am. Water Resour.

580        Assoc. 52. https://doi.org/10.1111/1752-1688.12363

581 Hutton, C., Wagener, T., Freer, J., Han, D., Duffy, C., Arheimer, B., 2016. Most computational

582        hydrology is not reproducible, so is it really science? Water Resour. Res. 50.

583        https://doi.org/10.1002/ 2016WR019285

584 Liang, X., Lettenmaier, D.P., Wood, E.F., 1996. One-dimensional statistical dynamic

585        representation of subgrid spatial variability of precipitation in the two-layer variable

586        infiltration capacity model. J. Geophys. Res. Atmos. 101(D16), 21403–21422.

587 Malik, T., 2017. GeoTrust: Improving Sharing and Reproducibility of Geoscience Applications

588        [WWW Document]. EOL Semin. Ser. URL https://www2.ucar.edu/for-

589        staff/daily/announcement-calendar-event/eol-seminar-series-dr-tanu-malik (accessed

590        6.6.17).

591 Malik, T., Valescu, C., Pham, Q., 2017. Sciunit, a system for creating, sharing, and running light-

592        weight containers. [WWW Document]. URL http://www.geotrusthub.org/ (accessed 1.1.17).

593 Meng, H., Kommineni, R., Pham, Q., Gardner, R., Malik, T., Thain, D., 2015. An invariant

594        framework for conducting reproducible computational science. J. Comput. Sci. 9, 137–142.

595        https://doi.org/10.1016/j.jocs.2015.04.012

596 Morsy, M.M., Goodall, J.L., Castronova, A.M., Bandaragoda, C., Greenberg, J., 2014. Metadata

597        for Describing Water Models, in: In Proceedings of the 7th International Congress on

598        Environmental Modelling and Software, DP Ames, NWT QuinnMorsy, M.M., Goodall, J.L.,

599       Castronova, A.M., Bandaragoda, C., Greenberg, J., 2014. Metadata for Describing Water

600           Models, in: In Proceedings of the. pp. 978–988.

601       Morsy, M.M., Goodall, J.L., Castronova, A.M., Dash, P., Merwade, V., Sadler, J.M., Rajib, M.A.,

602           Horsburgh, J.S., Tarboton, D.G., 2017. Design of a metadata framework for environmental

603           models with an example hydrologic application in HydroShare. Environ. Model. Softw. 93,

604           13–28. https://doi.org/10.1016/j.envsoft.2017.02.028

605       Niswonger, R.G., Panday, S., Motomu, I., 2011. MODFLOW-NWT , A Newton Formulation for

606           MODFLOW-2005. U.S. Geol. Surv. Tech. Methods 6, 44.

607       Peng, R.D., 2011. Reproducible research in computational science. Science. 334, 1226–1227.

608       Pham, Q., Malik, T., Foster, I., 2014. Auditing and maintaining provenance in software packages.

609           Int. Proven. Annot. Work. 97–109.

610       Piccolo, S.R., Frampton, M.B., 2016. Tools and techniques for computational reproducibility.

611           Gigascience 5, 1–13. https://doi.org/10.1186/s13742-016-0135-4

612       Qasha, R., Cala, Jacek, Watson, P., 2016. A Framework for Scientific Workflow Reproducibility

613           in the Cloud. IEEE 12th Int. Conf. e-Science 81–90.

614       Qin, J., Dobreski, B., Brown, D., 2016. Metadata and Reproducibility : A Case Study of

615           Gravitational Wave Research Data. J. Digit. Curation 11, 218–231.

616           https://doi.org/10.2218/ijdc.v11i1.399

617       Reitz, M., Sanford, W.E., Senay, Gabriel B., and Cazenas, J., 2017. Annual estimates of recharge,

618           quick-flow runoff, and ET for the contiguous US using empirical regression equations, 2000-

619           2013. U.S. Geol. Surv. data release. https://doi.org/https://doi.org/10.5066/F7PN93P0.

620       Reproducibility Guide - The rOpenSci Project [WWW Document], n.d. . 2017. URL

621           http://ropensci.github.io/reproducibility-guide/sections/introduction/ (accessed 6.16.17).

622     Santana-Perez, I., Ferreira da Silva, R., Rynge, M., Deelman, E., Pérez-Hernández, M.S., Corcho,

623           O., 2017. Reproducibility of execution environments in computational science using

624           Semantics and Clouds. Futur. Gener. Comput. Syst. 67, 354–367.

625           https://doi.org/10.1016/j.future.2015.12.017

626     Singh, V.P., Asce, F., Woolhiser, D.A., Asce, M., 2002. Mathematical Modeling of Watershed

627           Hydrology. J. Hydrol. Eng. 7, 270–292.

628     Stodden, V., 2013. Resolving Irreproducibility in Empirical and Computational Research. IMS

629           Bull. Online.

630     Swain, N.R., Christensen, S.D., Snow, A.D., Dolder, H., Espinoza-Dávalos, G., Goharian, E.,

631           Jones, N.L., Nelson, E.J., Ames, D.P., Burian, S.J., 2016. A new open source platform for

632           lowering the barrier for environmental web app development. Environ. Model. Softw. 85,

633           11–26. https://doi.org/10.1016/j.envsoft.2016.08.003

634     Tarboton, D.G., Horsburgh, J.S., Idaszak, R., Heard, J., Valentine, D., Couch, A., Ames, D.,

635           Goodall, J.L., Band, L., Merwade, V., Arrigo, J., Hooper, R., Maidment, D., 2014a. a

636           Resource Centric Approach for Advancing Collaboration Through Hydrologic Data and

637           Model Sharing. 11th Int. Conf. Hydroinformatics, HIC 2014, New York City, USA.

638     Tarboton, D.G., Idaszak, R., Horsburgh, J.S., Heard, J., Ames, D., Goodall, J.L., Band, L.,

639           Merwade, V., Couch, A., Arrigo, J., Hooper, R., Valentine, D., Maidment, D., 2014b.

640           HydroShare: Advancing Collaboration through Hydrologic Data and Model Sharing. Int.

641           Environ. Model. Softw. Soc. 7th Int. Congr. Environ. Model. Software, San Diego,

642           California, USA. www. iemss. org/society/index/php/iemss-2014-proceedings.

643           https://doi.org/10.13140/2.1.4431.6801

644     Weibel, S., Kunze, J., Carl Lagoze, A., Wolf, M., 1998. Dublin core metadata for resource

645        discovery. No. RFC 2413.

646