Building Cyberinfrastructure for the Reuse and Reproducibility of Complex Hydrologic Modeling Studies

4

7

Iman Maghami^a, Ashley Van Beusekom^b, Lauren Hay^b, Zhiyu Li^c, Andrew Bennett^d, YoungDon
 Choi^a, Bart Nijssen^b, Shaowen Wang^c, David Tarboton^e, Jonathan L. Goodall^{a,*}

- 8 ^a Department of Engineering Systems & Environment, University of Virginia, Charlottesville,
- 9 Virginia, USA
- 10 ^b Department of Civil & Environmental Engineering, University of Washington, Seattle,
- 11 Washington
- 12 ° Department of Geography & Geographic Information Science, University of Illinois at Urbana-
- 13 Champaign, IL, USA
- ^d Department of Hydrology and Atmospheric Sciences, University of Arizona, Tucson, AZ, USA
- ^e Department of Civil and Environmental Engineering, Utah Water Research Laboratory, Utah
- 16 State University, Logan, Utah, USA
- 17 * To whom correspondence should be addressed (E-mail: goodall@virginia.edu; Address:
- 18 University of Virginia, Department of Engineering System and Environment, University of
- 19 Virginia, 151 Engineers Way, P.O. Box 400747, Charlottesville, VA, 22904, USA; Tel: (434)
- 20 243-5019)
- 21 This is the accepted version of the article published in final form at:
- 22 Maghami, I., A. Van Beusekom, L. Hay, Z. Li, A. Bennett, Y. Choi, B. Nijssen, S. Wang, D.
- 23 Tarboton and J. L. Goodall, (2023), "Building cyberinfrastructure for the reuse and
- 24 reproducibility of complex hydrologic modeling studies," Environmental Modelling & Software,
- 25 164: 105689, <u>https://doi.org/10.1016/j.envsoft.2023.105689</u>.
- 26

27 Highlights

- Presents novel cyberinfrastructure for complex hydrologic modeling studies
- Focuses on the challenges introduced by computationally and data intensive studies
- Uses Globus for large data transfers between scientific cloud services
- Leverages containerization for model portability across compute environments
- Combines model APIs and Jupyter notebooks to document modeling workflows

33 Abstract

34 Building cyberinfrastructure for the reuse and reproducibility of large-scale hydrologic modeling

- 35 studies requires overcoming a number of data management and software architecture challenges.
- 36 The objective of this research is to advance the cyberinfrastructure needed to overcome some of 27 these shallon gos to make such asymptotic and hudrals gis studies assists reveal where Wa
- 37 these challenges to make such computational hydrologic studies easier to reuse and reproduce. We 38 present novel cyberinfrastructure capable of integrating HydroShare (an online data repository).
- 38 present novel cyberinfrastructure capable of integrating HydroShare (an online data repository), 39 CyberGIS-Jupyter for Water and high performance computing (HPC) resources (computational
- 40 environments), and the Structure for Unifying Multiple Modeling Alternatives (SUMMA)
- 41 hydrologic modeling framework through its application programming interface for orchestrating
- 42 model runs. The cyberinfrastructure is demonstrated for a complex computational modeling study
- 43 on a contiguous United States dataset. We present and discuss key capabilities of the
- 44 cyberinfrastructure including 1) containerization for portability across compute environments, 2)
- 45 Globus for large data transfers, 3) a Jupyter gateway to HPC environments, and 4) Jupyter
- 46 notebooks for capturing the modeling workflows.

47 Keywords

48 Reproducibility; Computational Hydrology; Jupyter; HPC; Containerization

49 Software and Data Availability

50 The data and Jupyter notebooks used in this study were published on HydroShare with persistent

51 citable Digital Object Identifiers (DOIs). A collection resource in Hydroshare (Choi et al., 2022a)

52 holds the resources containing the data and Jupyter notebooks further described in the following

table. A Hydroshare account (http://hydroshare.org) and access to the CyberGIS-Jupyter for Water

54 computing gateway (accessed through HydroShare) are required to execute the Jupyter notebooks

55 in the second and third resources.

Resource Description	Reference
Original NLDAS forcings for the CAMELS basins can be obtained as a NetCDF file*	Mizukami and Wood, 2021

SUMMA Simulations using CAMELS Datasets on CyberGIS-Jupyter for Water ^{**}	Choi et al., 2022b
SUMMA Simulations using CAMELS Datasets for HPC use with CyberGIS-Jupyter for Water ^{**}	Choi et al., 2022c

⁵⁶ *The data from the CAMELS dataset (Newman et al., 2015a) was consolidated into one NetCDF

57 file taking advantage of OPeNDAP data services supported by the HydroShare THREDDS

58 server and web application connector (Tarboton and Calloway, 2021).

59 **The SUMMA setup for the CAMELS basins can be obtained from the summa_camels folder

60 of the HydroShare resources.

61 List of relevant URLs

- 62 CyberGIS-Jupyter for Water: https://go.illinois.edu//cybergis-jupyter-water
- 63 Docker: https://www.docker.com
- 64 HydroShare REST API: https://www.hydroshare.org/hsapi/
- 65 Numpy: https://www.numpy.org
- 66 Pandas: https://pandas.pydata.org
- 67 pySUMMA: https://github.com/UW-Hydro/pysumma/releases/tag/v3.0.3
- 68 Seaborn: https://seaborn.pydata.org
- 69 Singularity: https://sylabs.io
- 70 SUMMA: https://github.com/CH-Earth/summa/releases/tag/v3.0.3
- 71 xarray: http://xarray.pydata.org
- 72 XSEDE: https://www.xsede.org

73 **1 Introduction**

74 Reproducibility, the ability to duplicate and verify previous findings, is a foundational principle in 75 scientific research. In computational hydrology, Melsen et al., (2017) highlighted two contrasting definitions of model reproducibility: (1) "bit- reproducibility" which is defined as exact replication 76 77 of a study, including the exact same numbers forming the results, and (2) "conclusionreproducibility" which focuses on reproducibility of the conclusions of a study as the conclusions 78 79 are expected to hold if the same experimental approach is applied. They argue that "conclusion-80 reproducibility" (replicating a study's conclusions) may be more important than "bitreproducibility" (exactly replicating model runs) because hydrological theories need to be tested 81 beyond bit-reproducibility by investigating conditions under which theories can be confirmed or 82 83 falsified. Even so, conclusion-reproducibility itself goes beyond the simple sharing of code and 84 data as open-source and online resources typically touted for achieving reproducibility. The code 85 and data must be accompanied by well-documented workflows with readable and reusable code 86 (Chen et al., 2020; Mullendore et al., 2021; Simmonds et al., 2022). Reusable code requires 87 providing open-source computational environments in which the code can be executed. Ensuring 88 this reuse and reproducibility is a non-trivial task; it requires not only adopting new capabilities 89 for handling complex software and big data, it also requires careful software engineering practices 90 to integrate these new capabilities into well designed and built cyberinfrastructure (Merkel, 2014).

91 A growing body of researchers have been discussing and proposing guidelines and strategies for 92 reproducible computational modeling (e.g., Bush et al., 2021; Choi et al., 2021; Knoben et al., 93 2022; Mullendore et al., 2021; Simmonds et al., 2022). In recent work, Knoben et el. (2022) 94 presented a novel approach for creating a hydrologic model at any location or scale (local to global) 95 by separating model-agnostic and model-specific configuration steps within cyberinfrastructure 96 workflows. Choi et al. (2021) described a general strategy for creating modern cyberinfrastructure 97 to support open and reproducible hydrologic modeling as the integration of three components: (1) 98 online data repositories; (2) computational environments leveraging containerization and self-99 documented computational notebooks; and (3) Application Programming Interfaces (APIs) that 100 provide programmatic control of complex computational models. As an example of this general 101 approach, Choi et al. (2021) also presented an implementation that used (1) HydroShare as the 102 online repository, (2) two different Jupyter instances, one hosted by the Consortium of Universities 103 for the Advancement of Hydrologic Science, Inc. (CUAHSI) and a second hosted by CyberGIS-104 Jupyter for Water, as the computational environments, and (3) pySUMMA, a Python wrapper for 105 manipulating, running, managing, and analyzing of SUMMA (Structure for Unifying Multiple 106 Modeling Alternatives), as the model API.

107 While Choi et al. (2021) focused mainly on the system design and demonstrated their approach 108 with a fairly simple modeling use case, reproducibility in computational hydrology can present 109 some difficult challenges when dealing with large-scale hydrologic studies (Hutton et al., 2016). 110 These challenges mostly pertain to the use of "big data" and computationally expensive and time-111 consuming resources needed for reproducibility of complex hydrologic modeling studies. Hutton 112 et al., (2016) notes that in these cases, new techniques are needed to ensure scientific rigor. In this 113 paper, we provide an example of the overall system design outlined by Choi et al. (2021) as applied 114 to a complex hydrologic study by Van Beusekom et al. (2022) (hereafter referred to as the VB 115 study). We develop the necessary cyberinfrastructure to reproduce this study for selected sub116 domains and discuss the challenges and opportunities in ensuring conclusion-reproducibility for

117 complex hydrologic studies.

118 The VB study evaluated the effect of the temporal resolution of surface meteorological inputs (or 119 forcings) on modeled hydrological fluxes and states for 671 basins across the contiguous United 120 States (CONUS). It quantified the difference in hydrologic outcomes based on daily or sub-daily 121 forcings for multiple model configurations and parameter values. Reproducibility of the VB study if one was given only the input data and model code would be challenging because it requires the 122 123 installation and configuration of the modeling framework SUMMA (Clark et al., 2015b, 2015a), 124 the data volumes are very large, and the model runs require High Performance Computing (HPC) resources. The complete VB study consisted of 704 6-year model runs for each of the 671 basins 125 126 (or 2.8 million model years). SUMMA was implemented with a single hydrologic response unit for each basin, resulting in a single output time series for each basin for each model configuration. 127 128 For every model run, the output consisted of 14 hydrological variables, which required 6 MB per 129 model simulation, or 2.834 TB for the entire study. While few researchers may be interested in 130 reproducing the entire VB study, the more common use case and the focus of this study, would be 131 to repeat or extend the VB study for a subset of the basins. We want to enable others to reproduce 132 the VB study for subsets of the original domain as a basis for doing additional research enabling 133 conclusion-reproducibility rather than the bit-reproducibility. For such an approach to be effective, 134 it is not sufficient to provide the open-source SUMMA code and model input data; one must also 135 provide the additional components described by Choi et al. (2021), i.e., computational 136 environments, models exposed through APIs, and documented model workflows to create a 137 cyberinfrastructure that lowers the barrier to reuse and reproducibility.

138 This research contributes to the growing literature advancing cyberinfrastructure for hydrology 139 and other geoscience fields. Yang et al., (2010) illustrated the importance of using HPC in 140 computationally intensive geospatial sciences and hydrologic modeling. Essawy et al., (2016) 141 demonstrated server-side workflows for large-scale hydrologic data processing, although they did 142 not make use of HPC in their application. Lyu et al., (2019) used containerization and combined 143 computational environments including HPC and High Throughput Computing (HTC) 144 cyberinfrastructure to directly run the models using Jupyter notebooks. Gan et al. (2020), 145 integrated a hydrologic data and modeling web service with HydroShare as a data sharing system 146 to show how this integration leads to a findable and reproducible modeling framework. Gichamo 147 et al., (2020) used web-based data services to prepare input data for hydrologic models. Kurtz et 148 al. (2017) introduced a cloud-based real-time data assimilation and modeling framework and 149 showed how parallel processing can be used for complex hydrologic models in the cloud. 150 However, unlike the VB study, none of Lyu et al. (2019), Gan et al. (2020), Gichamo et al. (2020) 151 and Kurtz et al. (2017) applied their methods on a computationally extensive complex hydrologic 152 use case. Therefore, the challenges and opportunities of using cyberinfrastructure for 153 reproducibility of complex, large-scale hydrologic modeling. for which HPC and big data 154 approaches are required, remain largely unexplored.

To address this research gap, we designed and implemented cyberinfrastructure to enable intuitive access to HPC computational environments and to support data transfers into and out of the HPC environment. Additionally, we provide a workflow that allows users to replicate parts of the study

157 environment. Additionary, we provide a worknow that anows users to replicate parts of the study 158 within their own computing environments. We also perform a workflow run-time performance

analysis that compares different model scenarios by varying the size of simulations across different

160 computing environments, providing users with a guide towards selection of the computing 161 environment depending on the size of their simulations. The cyberinfrastructure provides a starting 162 point for users to modify the hydrologic model setups, thus going beyond reproducibility (i.e., the 163 ability to duplicate and verify previous findings) into replication where one modeling methodology 164 can be used to answer the same scientific research question but with new input data (as highlighted 165 by Essawy et al. (2020)). The cyberinfrastructure may also serve as an educational resource by

- 166 providing an intuitive way for students to perform complex hydrologic modeling studies. The data
- 167 and cyberinfrastructure are provided through HydroShare to run on any basin for which we provide
- 168 a SUMMA setup to assist the modeler in analyzing basins individually.

169 The remainder of this paper is organized as follows. In Section 2, we provide a brief overview of 170 the VB study, the cyberinfrastructure, the model workflows, and the model scenarios used for a 171 science use case subsetted from the VB study as well as the model workflows run-time 172 performance analysis. Section 3 provides results and discussion. The results focus on the modeling 173 use case and an analysis of the workflow run-time performance for different computing 174 environments. The discussion focuses on opportunities and challenges learned from our experience 175 designing and building the cyberinfrastructure to support our modeling workflows. Finally, our 176 conclusions and recommendations are provided in Section 4.

177 **2 Methods**

178 2.1 Overview of the VB study

179 The VB study used 671 basins to study the effects of the temporal resolution of the meteorological 180 forcings on hydrologic model simulations across the CONUS. The basins are part of the CAMELS 181 dataset (Catchment Attributes and MEteorology for Large-sample Studies; Newman et al., 2015b) 182 a large-sample hydrometeorological dataset across the CONUS consisting of input forcings, basin 183 attributes, and relevant historical streamflow records. The VB study used SUMMA (Clark et al., 184 2015b) to configure multiple model instances for each basin, representing eight different model 185 configurations and 11 different sets of model parameter values. In addition, eight forcing datasets 186 were constructed. In each of these forcing datasets one of the meteorological inputs was modified 187 so that the diurnal cycle was replaced by the mean value over that day. The VB study performed 188 704 ($8 \times 11 \times 8 = 704$) 6-year model runs for each CAMELS basin, consisting of one year of model 189 initialization and five years of actual simulation. Model outputs for 14 simulated variables were 190 stored to evaluate the sensitivity of the simulations to changes in model forcings, model 191 configurations, and model parameters (Figure A1 and Table A1). The VB study results 192 demonstrated that (1) the effect of each forcing input on each model output varies by model output 193 and model location, (2) the use of a particular parameter set may not be critical in determining the 194 most and least influential forcing variables, and (3) the choice of model physics (i.e., using 195 different model configurations) could change the relative effect of each forcing input on model 196 outputs.

197 The VB study was run with scripts on the Cheyenne supercomputer (a 5.34-petaflops, high-

- 198 performance computer built for the National Center for Atmospheric Research; Computational and 199 Information Systems Laboratory (2017)), and it took a few days to complete the runs. For each
- 199 Information Systems Laboratory (2017)), and it took a few days to complete the runs. For each 200 basin, the output size for a single 6-year run was 6 MB. Thus, reproducing the entire study is

201 computationally expensive and also requires large amounts of storage (704 runs \times 671 basins \times 6 202 MB = 2.834 TB). However, the cyberinfrastructure allows individual basins to be run 203 independently. Here, we focus on a use case in which a researcher wishes to reproduce a subset of 204 the VB study by analyzing one or a few basins within a cloud cyberinfrastructure environment to reach conclusion-reproducibility. The conclusion-reproducibility that we aimed in this study is 205 206 solely a qualitative one and if the presented cyberinfrastructure can be successfully applied to 207 studies differing from the original study, i.e., the VB study, the conclusion- reproducibility is 208 achieved.

209 2.2 Cyberinfrastructure design and implementation

210 Following the approach described in Choi et al. (2021), we designed and implemented

211 cyberinfrastructure (Figure 1) to replicate the VB study by integrating (1) the HydroShare online

212 data repository, (2) CyberGIS-Jupyter for Water Computing Gateway (CJW CG) and high-

213 performance computational environments, and (3) a model API that can be utilized in scripts using

214 Jupyter notebooks (here the pySUMMA API). Each of these three components is further explained

215 in the following subsections.



216

Figure 1: The three primary components of the general cyberinfrastructure (following Choi et al.

218 2021) with seamless data transfers for open and reproducible environmental modeling.

219 2.2.1 Online data repositories

We used HydroShare, an online collaboration environment, as the online data repository (Horsburgh et al., 2016; Tarboton et al., 2014). A collection resource in HydroShare, which can be found at Choi et al. (2022a), contains three resources holding the data, computational environment, and models (Figure 1).

224 The HydroShare resource holding the data (Mizukami and Wood, 2021) contains the forcing data 225 set for the 671 CAMELS basins. The forcings are based on the hourly NLDAS-2 (North American 226 Land Data Assimilation System; NLDAS-2, 2014; NLDAS-2 is hereafter referred to as NLDAS). 227 The original NLDAS hourly forcing data were created on a 0.125 x 0.125 degree grid. To create 228 hourly SUMMA model forcings, NLDAS outputs were spatially averaged over each of the 671 229 CAMELS basins and merged into one NetCDF file. With this format, an OPeNDAP server 230 (OPeNDAP, 2021) can extract data for selected basins on the server, so that the user does not have 231 to download the entire CONUS dataset to a local computer. HydroShare offers this capability via 232 its THREDDS Data Server (TDS).

233 2.2.2 Computational environments

234 The developed computational environments provide a consistent software environment that is 235 independent of each user's own operating system and software libraries, making it possible to 236 study a computationally expensive research problem. Figure 2 shows each computational 237 environments component, and the interoperability between the computational environments and 238 HydroShare. One computational environment was implemented on the CJW CG cloud service for 239 studies with limited computational demand, e.g., a study of only a few basins, or as an instructional 240 tool, or for model debugging. A second computational environment was developed on an HPC 241 resource to reproduce a problem more representative of challenges posed by the use of big-data in 242 the VB study. The HPC environment also allows the user to study a particular basin in greater 243 detail. In this study, the CJW CG computational environment is used to provide (1) the model 244 execution environments configured as Docker images to enable execution of the SUMMA model 245 for studies with limited computational demand (i.e., those need to use CJW CG Workflow), and 246 (2) cyberinfrastructure for preprocessing, postprocessing and data storage for both studies with 247 limited computational demand (need to use CJW CG Workflow) and with high computational 248 demand (i.e., those need to use HPC Workflow) (Figure 2). The HPC computational environment 249 is only used for providing model execution environments configured as Singularity containers to 250 enable execution of the SUMMA model for studies with higher computational demand. More 251 details on each computational environment are provided in the rest of this section.



Figure 2: CJW and HPC computational environments with model execution environments
 configured as Docker image or Singularity container to support concurrent model execution
 through Jupyter notebooks, and use of Globus to transfer model outputs from HPC.

CJW CG is a cloud computing environment interoperable with HydroShare. It is an instance of CyberGISX (Yin et al., 2017) that serves the data- and computation-intensive needs of the water and environmental communities. We used CJW CG because it is publicly available, is interoperable with advanced cyberinfrastructure resources (such as the HPC resource used in this study) and has been serving the water and environmental communities to support their modeling needs.

263 Reproducibility was facilitated by using containerization of the SUMMA model and the 264 pySUMMA API with Docker (Merkel, 2014) in the case of the CJW CG environment or Singularity in the case of the HPC environment (Kurtzer et al., 2017) along with a computational 265 gateway interface to Jupyter notebooks (pySUMMA and the notebooks are described in a later 266 267 section) (Figure 2). Although using Docker is a common approach to containerize the model dependencies, we used Singularity in the HPC environment because it is designed to work 268 269 seamlessly with existing batch job systems to support HPC applications. The containerization and 270 interface are hosted on the CJW scientific cloud service hosted on Jetstream cloud (Hancock et al., 271 2021; Stewart et al., 2015; Towns et al., 2014). The Dockerfile is hosted on a GitHub repository 272 (Li, 2021) with pre-built docker images being shared on a Docker Hub repository. Singularity 273 container used by the HPC environment is hosted on CyberGIS-Compute Service, a middleware platform allowing seamless access to HPC resources via Python-based Software Development Kit 274 275 and core middleware services (CyberGIS-Compute Service, 2021; Li et al., 2022). The singularity 276 container was created through docker images conversion. CyberGIS-Compute Service also 277 handles submitting jobs to HPC as well as large data transfer from HPC through Globus (will be 278 discussed in section 2.2.4).

The Conda software package was used to manage the project specific computational environment on CJW, allowing the user to build a Python environment with the SUMMA model, pySUMMA API, and other computational dependencies. This was done by providing a kernel version for the project (CyberGIS Center HydroShare Development Team, 2022). Using this stable kernel, which captures all the required dependencies with their specific versions, ensures careful software version control.

285

286 2.2.3 Model Application Programming Interface (API)

The model API pySUMMA was chosen to be part of the interactive tool. The pySUMMA API 287 288 (Choi et al., 2021) wraps the SUMMA hydrologic modeling framework (Clark et al., 2015a) and 289 allows the user to script the use of the SUMMA model using Python. It facilitates model 290 configuration and allows for local execution of the model by either using a Docker container or a 291 locally compiled SUMMA executable (Choi et al., 2021). With pySUMMA, a user can modify 292 SUMMA input files and run SUMMA inside a Python script, as well as automatically parallelize 293 runs and visualize output. In the simplest case the pySUMMA Simulation object wraps a single 294 instance of a SUMMA simulation.

295 For users who choose to analyze multiple basins at a time in the CJW CG environment instead of

the HPC environment, the notebook automatically will configure a pySUMMA Distributed object,

which provides an interface to spatially distributed simulations and handles parallelism and job

298 management under the hood. In this study, multiple SUMMA simulations are run in each basin,

so a pySUMMA Ensemble object is used to manage multiple runs with different configurations.

- 300 In the HPC computational environment a custom backend was written to handle parallelism using
- 301 Message Passing Interface (MPI), reducing the need for users to customize the configuration based

302 on the type of job that they are running. A high-level description of pySUMMA is presented in

- 303 Figure 1. The simulation.py enables the execution of the SUMMA model and, along with
- 304 file_manager.py, decisions.py, force-file_list, and output_control.py, allows for manipulating
- 305 SUMMA configuration files. The distributed.py enables the parallel execution of SUMMA.
- 306 2.2.4 Data management and transfer

The input data for this study consists of the SUMMA configuration files and the forcing data for the 671 CAMELS basins. The configuration files (e.g., geometries information for the 671 CAMELS basins along with their attributes such as hru_id) are shared within each of the two HydroShare resources holding the Jupyter notebooks. The forcing data are provided in a HydroShare resource (Mizukami and Wood, 2021).

312 The output files resulting from running the notebooks using the CJW CG and HPC computational

313 environments are: (1) NetCDF output files generated by the SUMMA simulations, (2) a NetCDF

314 file recording the model performance for each basin as measured by the Kling-Gupta Efficiency

315 (KGE) (Gupta et al., 2009), and (3) additional files created by the notebooks such as the figures

that visualize the model results.

317 In the case of the CJW CG environment, after running the notebooks, all files are saved in the CJW

318 CG and are directly accessible to the user. In the case of the HPC environment, the KGE results 319 and other files created by the notebooks (e.g., figures) are automatically transferred to the CJW

and other files created by the notebooks (e.g., figures) are automatically transferred to the CJW
 CG, but the NetCDF output files remain within the HPC environment to avoid transferring large

- volumes of model output (as a reminder, the size of the model output for the entire VB study was
- 322 2.834 TB).

323 However, if the user of the HPC environment wishes to transfer selected SUMMA NetCDF output 324 files from the HPC to be directly accessible for further analysis and long-term storage, then the 325 CyberGIS-Compute Service (Li et al., 2022) can be used for reliable high-performance large file 326 transfers through the Globus service (Chard et al., 2016; Foster, 2011). As shown in Figure 2, data 327 is transferred from HPC to the CJW using Globus without going through the job submission server. 328 Globus is a software as a service that enables the transfer of datasets of any size between different 329 storage options (personal computers, HPC, etc.) without users being required to be constantly 330 logged in and monitoring the data transfer (Chard et al., 2016). Technically, the CyberGIS-GIS Compute acts as a Globus app client holding a community Globus account that has access to both 331 332 data endpoints on the Jupyter and target HPC. When data transfer is needed, CyberGIS-Compute 333 initiates a Globus task between the two endpoints and monitors the progress. Users are updated

- 334 with data transfer status in the notebooks environment during the entire process.
- 335
- 336

337 2.3 Model Workflows as Jupyter Notebooks

338 As mentioned earlier, the model workflows allow the user to reproduce all or subsets of the VB 339 study using either the CJW CG computational resources (referred to later as CJW CG) or the HPC and CJW CG computational resources (referred to later as HPC). The CJW CG and HPC 340 341 HydroShare resources can be found at Choi et al. (2022b) and Choi et al. (2022c), respectively. 342 The model workflows are documented in three (for CJW CG) or four (for HPC) Jupyter notebooks. 343 Table 1 shows the summary of the steps taken in each notebook, while Figure A2 - A5 show more 344 detailed information for notebooks 1-4. The first three notebooks for both the CJW CG and HPC 345 environments focus on (1) selecting the study basins, simulation period, and model input forcings, 346 (2) running the SUMMA model, and (3) exploring outputs to analyze the effect of each forcing 347 variable in each basin. The HPC computational resource uses a fourth notebook to transfer large 348 unprocessed output data from the HPC to CJW using GLOBUS. Notebooks 1 and 3 are very 349 similar between the two HydroShare resources, and both CJW CG and HPC HydroShare resources 350 use CJW CG computational resources to run these two notebooks. The second notebook differs 351 for the two environments, and the difference is explained in Section 2.3.2. These notebooks assist 352 a modeler in analyzing CAMELS basins individually, providing information on forcings and 353 output variables that are the most/least sensitive in their basin. With some additional work, the 354 CJW CG computational environment can also be hosted on other (non CJW) cloud services, but 355 the HPC environment is more tailored to interact with the CJW cloud service used here.

356

357

Table 1. Overview of the notebook 1-4.

#	Notebook Name	Goal	CJW CG or HPC
1	Preprocessing	Prepares forcings, and sets study basins and simulation period	Very similar between HPC and CJW CG environment
2	SUMMA execution	Runs the SUMMA model	Different versions for HPC and CJW CG environment
3	Post-processing	Explores outputs to find out effect of each forcing variable in each basin	Very similar between HPC and CJW CG environment
4	Use Globus to transfer big data	Transfer raw output from HPC to CJW using Globus service	Only for HPC environment

358

359

360 To use the HPC computational resource, the user must obtain access to the HPC by issuing a request through HydroShare to use CJW. Once this access is granted, users are automatically given 361 362 free access to two alternative HPC resources: (1) the Virtual ROGER (Resourcing Open Geospatial 363 Education and Research) HPC administered by the School of Earth, Society, and Environment at University of Illinois Urbana-Champaign (UIUC) which is integrated with the Keeling compute 364 365 cluster at UIUC ("Virtual Roger User Guide," 2022) and (2) the Expanse HPC, a much larger NSF 366 XSEDE resource operated and managed by San Diego Supercomputer Center (SDSC) ("Expanse System Architecture," 2022). In theory, the CyberGIS-Compute Service can support other HPCs 367 as well, but we did not test other HPCs. In this study, among the provided HPC options, we only 368

369 used Expanse to demonstrate the cyberinfrastructure: in our initial experiments Expanse HPC 370 performed faster than Virtual ROGER and the goal here was to show how a HPC can scale up a 371 study by speeding up the modeling process compared to a non-HPC environment rather than an 372 inter-comparison between different HPCs. Users who do not wish to use HPC computational 373 resources can use CJW CG computational resources directly to run smaller modeling jobs.

374 The hardware specifications of the CJW CG and the Expanse HPC are compared in Table 2. The 375 CJW CG has only 3 compute nodes each of which has eight CPUs with 1.996 GHz Clock Speed and 30 GB DRAM. Each user can only use up to six CPUs and the CPUs can be shared among 376 377 users. This means the maximum degree of parallelism for simulations using this computational resource is six. Thus, in case of running one basin from the VB study (704 runs) and using all the 378 379 six available CPUs, each CPU will need to run 117.33 simulations (some of them 117 and others 380 118 simulations). The Expanse HPC has 728 AMD Rome standard compute nodes each of which 381 is equipped with 256 GB DRAM and 128 2.25 GHz CPUs ("Expanse User Guide," 2022). The 382 Expanse HPC allows the user to only use up to 2 nodes at a time, i.e., 256 CPUs or the maximum 383 degree of parallelism for simulations. Thus, if a user is running one basin from the VB study (704 runs) and using all the available 256 CPUs, then each CPU will need to run 2.75 simulations (some 384 385 of them 2 and others 3). This shows how the HPC resource can scale up the model runs offering a high-performance tool. More details about the run-time performance of the notebooks are 386 387 discussed in the results and discussion section.

388

- 389
- 390
- 391

Table 2. Hardware specifications of the computational environments.

Computational Environment	Node count	Number of CPU cores per node (for parallel runs only)	Clock Speed (GHz)	DRAM/node (GB)
CJW CG*	3	8	1.996	30
Expanse HPC**	728	128	2.25	256

^{392 *}AMD EPYC-Milan Processor. Each user can only up to 6 CPUs and the CPUs can be shared393 among users.

394 **AMD Rome Standard Compute Nodes. Each user can only use up to 2 nodes, which means

395 256 CPUs, the maximum number of parallelism for simulations.

396

397 The following subsections discuss the general purpose of each notebook used to reproduce parts

- of the VB study. For specific coding details, refer to the notebooks in the HydroShare resourcesat Choi et al. (2022b) and Choi et al. (2022c).
- 400 2.3.1 Data processing notebook

401 The first notebook (JN 1: Preprocessing) processes the original CAMELS SUMMA files and the

402 input forcing datasets (Table A2). The user can select one or more CAMELS basins (1-671 basins)

403 but by selecting a higher number of basins the computational time and expense increases. 404 Notebook 1 subsets the original CAMELS SUMMA files, producing SUMMA attributes, 405 parameters, initial conditions, and hourly NLDAS forcing files for the selected basin(s). Then, 406 additional forcing datasets for the hydrologic model sensitivity study are developed from the NLDAS data files (FORCINGS box in Figure A1) as discussed below. 407

408 For each SUMMA-model setup, variations in 14 SUMMA-generated outputs, described in Table 409 A1, are examined with respect to variations in seven input forcings (air pressure (prs), air 410 temperature (*tmp*), long wave radiation (*lwr*), precipitation rate (*ppt*), specific humidity (*hum*), shortwave radiation (swr), and wind speed (wnd)), under different model parameterizations and 411 412 configurations. The SUMMA outputs generated with the one-hour NLDAS forcing dataset are 413 considered the benchmark (NLDAS dataset 1; FORCINGS box in Figure A1). The rest of datasets 414 (*ppt* to *prs* datasets; FORCINGS box in Figure A1) are developed, holding each of the individual 415 forcing variables constant over a 24-hour period while the other six forcing variables contain the

416 original hourly NLDAS values.

417 Figure A2 shows the steps taken in the first notebook. This notebook is the same for the CJW CG

418 and HPC environments except that the simulation time period and basins to be explored are pre-

419 populated differently. The user can change these setups in the third step of this notebook (step

420 1 3). In the last step of this notebook, users can visualize the individual forcing variables held 421

- constant over a 24-hour period against the original hourly NLDAS values using hourly and
- 422 cumulative plots.

2.3.2 SUMMA execution notebook 423

424 The second notebook (JN 2: Running SUMMA) executes the SUMMA model using the input data 425 from the first notebook for four different sets of SUMMA basin runs, outlined in Figure A1 (RUNS 426 box) and described in detail in The VB study. The first set of basin runs (DEFAULT; 8 SUMMA 427 runs per basin; RUNS box) uses the eight forcing datasets (FORCINGS box) combined with 428 default parameters and a default SUMMA configuration. The SUMMA default configuration is 429 set in the resource model decision file.

430 The second set of basin runs (LHS; 88 SUMMA runs per basin; RUNS box in Figure A1) uses the 431 eight forcing datasets combined with 11 parameter sets and a default SUMMA configuration. The 432 11 parameter sets consist of the default parameter set and 10 additional parameter sets with 15 433 commonly calibrated parameters (Table A2). As detailed in the VB study, the parameters are 434 sampled using Latin Hypercube Sampling (LHS) over their defined range. The pyDOE LHS function (Lee, 2014) is used to create unique 10 x 15 LHS sampling matrices for the selected basin. 435 436 Then the LHS matrices are used to produce 10 parameter sets of the 15 parameters while 437 considering the parameter constraints listed in Table 2. The choice of a different seed value will 438 lead to different LHS sets (and these sets will be different from the ones used by the VB Study).

439 The third set of basin runs (CONFIG; 64 SUMMA runs per basin; RUNS box in Figure A1) uses 440 the eight forcing datasets combined with the default parameter set and eight SUMMA 441 configurations. The eight SUMMA configurations, outlined in the CONFIGURATIONS box in 442 Figure A1, test three model decisions (stomatal resistance (stomResist), choice of snow 443 interception parameterization (snowIncept), and choice of canopy wind profile (windPrfile) with two options for each decision. Note the default configuration for this study is shown in bold in the
 CONFIGURATIONS box in Figure A1:BallBerry, lightSnow, and logBelowCanopy.

The fourth set of basin runs (COMPREHENSIVE; 704 SUMMA runs per basin; RUNS box in Figure A1) includes the DEFAULT, LHS, and CONFIG basin runs, and is the only set that needs to be run to replicate a single basin sensitivity study following the VB study method (six years of simulation must be run for replication). For testing purposes, sets 1-3 can also be run by themselves. The 10 parameter set files for the basin from the LHS sampling plus the default parameters (11 parameter sets) are run each with eight SUMMA configurations (CONFIGURATIONS box in Figure A1).

453 Figure A3 shows the steps taken in the second notebook. The first two steps in this notebook are 454 the same for the CJW CG and HPC environments but the rest of the workflow differs. In the CJW 455 CG notebook, the user can define the simulations by selecting the simulation period, model 456 configuration, and/or parameter values. Depending on which run complexity choice (i.e., 457 DEFAULT, LHS, CONFIG, COMPREHENSIVE in the RUNS box in Figure A1) is selected the 458 notebook executes a specific set of code cells using a conditional statement logic (e.g., if user 459 selects *config* prob == 1, step 2 7 is run which leads to CONFIG runs as shown in the RUNS box 460 in Figure A1). Users need to carefully consider the number of basins and the length of the 461 simulation period as the CJW CG environment is not powerful enough to run large simulations in 462 a reasonable time. In the HPC notebook, we only provided the user with the option to run the most 463 complex problem, i.e., *lhs config prob*, as the HPC is powerful enough to run the full problem 464 making it unnecessary to allow for simpler problems. The user can still change the simulation period (in step 2 3 of the workflow in Figure A3). The other main difference between the CJW 465 466 CG and HPC notebooks is that the codes calculating KGE values for the HPC notebook are 467 executed on the HPC (Step 2 8 in HPC branch in Figure A3) while for the CJW CG environment, 468 the KGE values are calculated locally on CJW CG (Step 2 9 in CJW CG branch in Figure A3). In 469 the HPC environment, the KGE values are calculated on the HPC resource to prevent having to 470 transfer large data volumes from the HPC to the CJW CG with the sole purpose of calculating 471 performance metrics. Users can use Globus to transfer selected output files from HPC to the CJW 472 CG for additional analysis. Notebook 4, which exists only in the HPC environment, was developed

473 for this purpose and is discussed in section 2.3.4.

A modified and scaled (range between -1 and 1) version of the KGE was used as an indicator of 474 475 model output sensitivity to a change in input forcing based on the work of Clark et al. (2021) and Mathevet et al. (2006) and is described in the VB study. The KGE test compares hourly model 476 477 outputs generated with the benchmark forcing dataset (NLDAS dataset 1; Table A2) with outputs generated with the forcing datasets with one forcing held constant (CNST datasets 2-8; Table A2). 478 479 KGE values are ranked from low to high to determine relative order of forcing influence on model 480 outputs with highest rankings associated with least influence of change to 24-hour constant 481 forcing.

482 2.3.3 Post-processing notebook

The third notebook (JN 3: Post-processing) produces visualizations of the sensitivity of SUMMA model output to the temporal resolution of the model forcing. Figure A4 shows the steps taken in the third notebook. The notebooks for CJW CG and HPC environments are the same. For the selected basin(s), eight plots are generated with Notebook 3 that follow the analysis in the VB
study. The reader is referred to the supplementary materials and the VB study for a detailed
explanation of each of the eight plots. In this paper, we only present the second figure generated
by Notebook 3, i.e., KGE values for each output variable for all 8 DEFAULT model runs.

490 2.3.4 Model output transfer

The fourth notebook (JN 4: Use Globus) is only included in the HPC resource (Figure A5) to transfer SUMMA output files from HPC to CJW on HydroShare. To retrieve the data from the HPC, this notebook needs a job ID submitted to the HPC and created in Notebook 2. While this notebook is running users can see the live status of the file transfer managed by the CyberGIS-Compute Service. Once running of this notebook is successfully finished, the user will be able to see the location of the transferred file on CJW.

497 2.4 Performance analysis

498 We tested the performance of the cyberinfrastructure using a number of model scenarios, using six 499 years of simulation (to be consistent with the VB study) and varying the number of studied basins 500 for each computational environment, described in Table 3. For the CJW CG environment, we 501 tested the performance of notebooks 1-3 for three scenarios (Table 3, rows 1 - 3): (1) one basin (a 502 total of six years of simulations), (2) four basins (a total of 24 years of simulations), and (3) six 503 basins (a total of 36 years of simulations). We decided not to test the CJW CG environment for 504 more basins as the CJW CG runs were slow and the HPC resource was available for larger 505 simulations.

506 For the HPC environment, we used Expense HPC, and tested the performance of notebooks 1-3 507 for 12 scenarios (Table 3, rows 4 - 15). In these scenarios, we varied the number of allocated CPUs 508 (128 or 256) for parallelism and the total number of basins ranging from one basin (a total of six 509 years of simulations) to 20 basins (a total of 120 years of simulations, which equals about three percent of the total simulation years for the whole VB study). To test the performance of Notebook 510 511 4, transferring output files from HPC to the CJW, we only used scenarios HPC 256 1 to HPC 512 256 6 (rows 4 - 9 in Table 3) and repeated each transfer 5 times to obtain a range of run-time for 513 each of the scenarios.

514

Table 3. Model scenarios for notebooks run-time performance analysis.

Row	Model scenario name	Number of CPU cores allocated	Number of basins	Simulation years	Total simulation years
1	CJWVM_1	6	1	6	6
2	CJWVM_2	6	4	6	24
3	CJWVM_3	6	6	6	36
4	HPC_256_1	256	1	6	6
5	HPC_256_2	256	4	6	24

6	HPC_256_3	256	6	6	36
7	HPC_256_4	256	10	6	60
8	HPC_256_5	256	15	6	90
9	HPC_256_6	256	20	6	120
10	HPC_128_1	128	1	6	6
11	HPC_128_2	128	4	6	24
12	HPC_128_3	128	6	6	36
13	HPC_128_4	128	10	6	60
14	HPC_128_5	128	15	6	90
15	HPC_128_6	128	20	6	120

516 **3 Results and Discussion**

In this section, we first briefly present results of the modeling case study that served as a motivating use case for the cyberinfrastructure. Then, we present results of the performance analysis focusing on contrasting the CJW CG and HPC notebooks using a variety of model setups. Then, we summarize the resulting resources from this study that are shared on HydroShare. Finally, we discuss the resulting system including opportunities and challenges identified through this research that can be the focus of future research.

523 3.1 Results of the modeling case study

Four CAMELS basins with diverse characteristics (Table 4) were chosen as examples of the effect 524 525 of basin characteristics on model results. We specifically selected these four basins for this modeling case study because we found that they all show different patterns. For the four selected 526 527 basins, Figure 3 shows the KGE values for each SUMMA output variable using the DEFAULT 528 (BIL; CONFIGURATIONS box in Figure A1) model configuration runs. The runs consist of one 529 reference simulation in which all forcing variables vary on an hourly basis (NLDAS dataset 1; FORCINGS box in Figure A1) and seven simulations in which one forcing variable is held 530 531 constant at the mean daily value throughout each day (the seven datasets ppt to prs; FORCINGS 532 box in Figure A1). KGE values were calculated relative to the reference simulation for each of the 533 seven simulations using five years of hourly model output from 10/1/1991 - 9/30/1996.

534

Table 4. Basin descriptions for individual basin analysis.

			CAMELS Attributes					
USGS Station ID	Name	Drainage area (km ²)	Gage datum (m)	Mean daily precipitation (mm/day)	Fraction of precipitation falling as	Aridity	Mean daily discharge (mm/day)	Runoff ratio*

					snow			
01632900	Smith Creek Near New Market, VA	242	268	2.91	0.10	0.89	0.80	0.27
02212600	Falling Creek near Juliette, GA	187	1202	3.37	0.01	1.19	0.74	0.22
09378630	Recapture Creek Near Blanding, UT	10	2195	1.58	0.50	0.50	0.21	0.13
11264500	Merced River at Happy Isles Bridge near Yosemite, CA	469	1228	2.64	0.91	1.15	1.94	0.73

535 * Annual runoff / annual precipitation



541 Figure 3 demonstrates the variability in model output sensitivity to the temporal resolution of the

542 forcing variables. The first three basins (gages 01632900, 02212600, and 09378630) show a strong 543 *ppt* temporal aggregation influence using DEFAULT, whereas gage 11264500 is more influenced by *tmp*, *hum*, and *swr* temporal aggregation. In other words, a higher temporal resolution is necessary for the aforementioned forcing variables in the given basins to capture the sub-daily hydrologic response shown by the reference simulation. The weaker influence of *ppt* temporal aggregation on the gage 11264500 compared to other gages can be attributed to its high fraction of precipitation falling as snow, 0.91 as opposed to 0.1, 0.01, 0.5 (Table 4).

549 Also in Figure 3, we see varying ranges in KGE values for particular output variables. As an 550 example, SurfaceRunoff is affected by constant hourly values of ppt for gages 01632900 and 551 09378630; ppt and hum for gage 02212600; and tmp, hum, swr, wnd, ppt, and prs (most to least 552 dominant) for gage 11264500. This shows the forcing variables in each basin that need to have a 553 higher temporal resolution to reproduce the SurfaceRunoff output in the reference simulation. In 554 this section, we only presented one example of an inter-basin comparison to illustrate how different 555 the results can be across different basins. Researchers can further explore the differences between 556 individual basins using other plots that can be made using the interactive Jupyter notebooks, and 557 also reproduce the results from the original VB study.

558 3.2 Results from performance analysis

559 Figure 4 shows the run-time for the data processing notebook (Notebook 1) and the post-560 processing notebook (Notebook 3) for the 15 scenarios listed in Table 3. Notebooks 1 and 3 are very similar between CJW CG and HPC computational environments. Notebooks 1 and 3 do not 561 562 take a significant time to run because they are only preprocessing and output analysis notebooks, 563 and no simulations are run. For scenarios with fewer than 30 simulation years, Notebook 1 takes 564 longer than Notebook 3, but this changes for scenarios with more simulation years as the rate of 565 run-time increase with simulation years is much higher with Notebook 3 than with Notebook 1. 566 For the CJW CG environment, the average time to run Notebooks 1 and 3 across the tested 567 scenarios only takes 0.6% of the entire time needed to run all Notebooks 1, 2, and 3. This means 568 the time required to run data processing and post-processing notebooks is not a limiting factor for 569 running the simulations. For the HPC environment, this ratio increases to 8.5% and 11.3% when 570 using 128 and 256 CPUs, respectively. This dramatic increase in the ratio is due to the significant 571 decrease in run-time of Notebook 2 when using HPC.



573

Figure 4. Notebook 1 (JN1) and 3 (JN3) run-time performance analysis for different model
simulations (both JN1 and JN3 were run on CJW CG no matter whether the HPC or CJW CG
environment was used for the modeling; therefore, we do not distinguish between the
environments in this figure).

579 The run-time for the SUMMA execution notebook (Notebook 2) for the 15 model scenarios using 580 different computation environments is shown in Figure 5. The high rate of run-time increase with 581 increasing simulation years for the CJW CG environment emphasizes that while the CJW CG 582 environment is technically able to simulate smaller models, it might not be fast enough to run 583 larger simulations. In the case of running six basins for six years, the HPC was 3.6 and 2.6 times 584 faster than the CJW CG, when using 256 and 128 CPUs, respectively. HPC with 256 CPUs (scenario HPC 256 6) could finish the simulations for 120 years (3% percent of the VB study) in 585 586 2 hr and 10 min while HPC with 128 CPUs (scenario HPC 128 6) could run the same problem in 1.48 times of the time need by HPC 256 6. Using the HPC with 256 CPUs, assuming a 587 588 conservative linear extrapolation, the SUMMA simulations from Notebook 2 are expected to be 589 done in about 75 hours for the entire VB study. In summary, HPC provides considerably faster 590 simulations making them ideal to run for larger studies.

591 When using the HPC resource and in the case of 120 years of simulation, dividing the number of 592 the allocated CPUs by two led to about a 50% increase in the run-time and not 100% as one might 593 expect. This non-linear scaling can be mainly attributed to 1) communication overhead in the 595 Notebook 2 did not utilize parallelism. For example, KGE values were only calculated after they 596 were exported as NetCDF files instead of being calculated directly from the raw SUMMA output 597 files. The rate of run-time increase for HPC with 128 CPUs is higher compared to that for HPC 598 with 256 CPUs. This may be attributed to the communication overhead because each CPU in the 599 case of the HPC with 128 CPUs needs to run twice as many simulations compared to HPC with 500 256 CPUs.

601



Figure 5. Notebook 2 run-time performance analysis for different model simulations using the
 CJW CG, or HPC (Expanse with 256 or 128 CPUs) options.

605

602

606 The run-time for transferring the SUMMA output files from Expanse HPC to CJW on HydroShare using the Globus service integrated by CyberGIS-Compute Service is shown in Figure 6. Each 607 608 transfer was repeated 5 times to obtain a range of run-time for each of the model simulations with 609 a different total number of simulation years. The range of the transfer time for each total number 610 of simulation years is small, indicating a consistent data transfer. For 120 years of simulation, it 611 took 14.5 min on average to transfer 118 GB of data from HPC to CJW, highlighting that the data transfer approach from HPC to CJW is fast and stable. The transfer rate (GB/min) is independent 612 of data size (Figure 6). 613



Figure 6. Boxplots for Notebook 4 run-time performance analysis for five different simulation
years to transfer data from Expanse HPC to CJW on HydroShare. Each transfer was repeated
five times to obtain a range of run-time for each of the model simulations with a different total
number of simulation years.

620 3.3 Data organization in HydroShare

615

The data for this study was pre-processed and the output post-processed by using existing Python packages. The study demonstrates the potential for using the online repository of HydroShare to not only store data and modeling code, but to also store computational environments, API version documentation, and container installation. HydroShare, as a hydrology-based repository service, facilitated this by allowing all the parts of the problem to be stored together as one resource. Furthermore, parts of the resource can be extracted and made into a new version of the resource (updated, revised, or modified), to promote collaboration.

To this point, a HydroShare collection resource was created that contained three composite resources. These resources are published and have Digital Object Identifier (DOI) which makes them immutable and findable. Figure 7 shows the landing page for the HydroShare collection

- resource that groups the three composite resources. The three composite resources that are contained by this collection resource are shown in dialogue box 1, the "Related Resources" in box 2 refers to this paper, and box 3 shows the information on how to cite this resource. Figure 8 shows the landing page for the HydroShare composite resource holding the HPC notebooks. Box 1 shows the contents of the resource, most importantly the four Jupyter notebooks and the readme.md file. The readme.md file (box 2) provides the user with the instructions on how to run the notebooks. Box 3 shows the information on how to cite this HydroShare resource.



Figure 7 The HydroShare landing page for the collection resource developed by this study (Choi et al., 2022a).

	🤐 🤼 🖓 19 🙆 🔺	1
Authors:	Young-Don Choi Ashley Van Beusekom Zhiyu/Drew Li Bart Nijssen Sharing Status: Public Lauren Hay Andrew Bennett David Tarboton Iman Maghami Manuari 200	U
	Jonathan Goodall Martyn P. Clark Downloads: 17	
Owners:	Iman Maghami Jonathan Goodall Bart Nijssen Young-Don Choi Andrew Bennett Ashley Van Beusekom Zhiyu/Drew Li +1 Votes: Be the first one to 🛐 this.	
Resource type:	Composite Resource Comments: No comments (yet)	
Storage:	The size of this resource is 9.0 MB	
Created:	May 20, 2021 at 12:35 a.m.	
Last updated:	Mar 13, 2022 at 1:23 a.m. Iman Maghami	
Citation:	see how to cite this resource not shown due to space limit	
Abstract		
Subject K	eywords	
VSEDE	NAMES CALLES SILAMA	
(ABEDE) (C	JURIUIS CAMILES SUMMA	
Resource	Level Coverage	
Contont		
Content		
(+) $(+)$	1 II V Sort by - Q. Search current directory	0
	🗞 🕘 📩 🔀 🚔 Learn more about the Bagit download	
Ch contents		
E) contents		
1_camels_m	L 2_camels_py_ 3_camels_an_ 4_globus_do_ constantDay_ constantSW_	
Readme.mo	i summa_cam	
	A.	_
🗋 Readr	ne.md	
Harris	a num ale s since lations	
How	to run the simulations	
This Readm The steps, i	ie file provides the users with the step-by-step guide to successfully run the three developed notebooks. In the order they need to be taken, are explained in what follows.	
STEP_	u Preliminary step	
In this sten	the modellers make sure that they have access to the content files of the resource and required compute platform	-
Related I	Resources	-
Oradita	Related Resources and Credits not	
Credits	shown due to space limit	
	snown due to space mint	
How to (Dite	
Chol, Y., A. Vi for HPC use	an Beusekom, Z. Li, B. Nijssen, L. Hay, A. Bennett, D. Tarboton, I. Maghami, J. Goodall, M. P. Clark (2022). SUMMA Simulations using CAMELS Datasets with CyberGIS-Jupyter for Water, HydroShare, http://www.hydroshare.org/resource/9d73d61696ea4f6b9c9a11e21cd44e24	Co
		-

Figure 8 The HydroShare landing page for the HPC resource developed by this study (Choi et al., 2022c).

646 3.4 Opportunities and challenges

647 This study demonstrated a real-world working implementation application of strategies for 648 reproducible hydrologic modeling presented by Choi et al. (2021) to a large-scale hydrologic study 649 (the VB study). This section discusses the opportunities and challenges of this implementation. If 650 one needs to adopt this cyberinfrastructure for studies significantly differing from the VB study, 651 considerable changes or extra steps might be needed. For instance (1) if exploring non-CAMELS 652 basins, then extra steps to prepare the inputs might be needed, or (2) if using hydrologic models 653 other than SUMMA, then containerization of the model might be needed. Despite the plausible 654 challenges when making these non-trivial extra steps, the intended main opportunity here is that 655 the modeling community can learn from the presented open cyberinfrastructure considering the 656 commonalities among the hydrologic models with regard to the input data, preprocessing, 657 processing, and postprocessing steps needed by them (Knoben et al., 2022).

658 Minimal changes in the notebooks are required to use the presented cyberinfrastructure to rerun 659 parts or all of the VB study or to extend the experiments performed in that study for selected 660 CAMELS basins. With these minimal changes, a user could use (1) different CAMELS basins, (2) different parameters in the LHS set, (3) different simulation periods, e.g., a drought period, (4) 661 more than 10 LHS sets, e.g., a more thorough exploration of the parameter space, and (5) additional 662 663 SUMMA model configurations. The last two changes, i.e., using a larger number of LHS sets and 664 different model configuration/decisions, highlights a major challenge in reproducing a 665 computationally complex study. Here, the limit on manageable data size was pushed, even when 666 running a few basins. HPC computational power was required to run the full six years of 667 simulation; expanding the parameter exploration space or adding model decisions would 668 compound the data size. Thus, while this work is advancing cyberinfrastructure used for big data 669 in hydrology, challenges remain.

670 The second major challenge that is encountered is implementing version control. What if users 671 need to run the Jupyter notebooks presented in this study in their own computational environment 672 (not deployed on CJW), or they need to install a newer version of a model API? How can they make sure they have a reproducible framework that is robust enough to tackle the version control 673 674 problem? Because there are many individual pieces of software, it was challenging at times for the 675 study team to keep all the software versions synchronized. We propose that future research should 676 tackle the version control challenge by making the computational environment all documented and 677 installable via a Python environment file. The pySUMMA code, which is used for hydrology 678 modeling, was installed via conda just as the rest of the infrastructure. In the future, Python 679 package updates will break compatibility, but compatibility can be preserved by installing the older versions (as documented in the environment file), or the user understanding the updates in order 680 681 to manually work around the updated package incompatibility. If a researcher wants to use a newer 682 (future) version of pySUMMA, then they may need to debug some parts of the Jupyter notebooks 683 that are affected by the changes. While this is not an ideal way to handle version updates, at least 684 the researcher has options of a working, albeit older, computational environment, from which to 685 begin reproducing the study before updating to newer software.

686 The specifics of the environment can be placed in a Python environment.yml file that can be shared 687 as part of the online model and data repositories, and can be installed with an installation notebook 688 inside the repository. This can use best practice for transparency about what dependencies the 689 computational gateway interface notebooks need to run. The specifics of each dependency can be

described in the installation notebook, so that if in the future there are issues with the availability

691 of that dependency, then a suitable substitute can be found. Version control issues can be thus 692 addressed through this methodology, albeit an imperfect solution depending on possible user

693 troubleshooting.

694 In addition to the two major challenges described above, there are two additional challenges related 695 to the use of the HPC environments: 1) large data transfers between computational environments, 696 online data repositories, and a user's personal computer and 2) allowing users to execute their 697 workflows on different HPC environments based on their use case and access to HPC 698 environments. There may be cases, for example, where users does not want to utilize HPC 699 resources due to financial cost concerns and need to transfer a large amount of model outputs from 700 an HPC environment's temporary scratch directory to a Jupyter compute environment to further 701 analyze the data using the Jupyter compute environment. Transferring large datasets, e.g., the 702 entire output from VB study or even the four selected basins study explored in this paper, would 703 be slow and unreliable using standard data transfer approaches, i.e., compress data into a big 704 package and then transfer it. In this study, we used Globus to do this data transfer which can 705 transfer multiple individual files in parallel without a need to compress data a big package, and 706 other related cyberinfrastructures that do not currently use Globus or a related technology could 707 benefit from doing so. Globus is not limited to data transfers between the HPC environment and 708 the Jupyter compute environments (CJW in the case of this study), however. In fact, it is possible 709 that the full or a large portion of the model output can be stored on an online data repository or 710 even on a user's own personal computer. In either case, the online data repository or the user's 711 personal computer, the outputs could be downloaded using Globus if Globus is installed, and they 712 become a Globus server. Making a user's personal computer a Globus server may be the case that 713 the user prefers to back up a model run not in an online data repository but at some other location. 714 In this case, Globus could be used to connect directly with the HPC environment thereby bypassing 715 both any Jupyter compute environments (CJW in the case of this study) as well as online data 716 repositories (HydroShare in the case of this study) as an intermediate storage location. If the large 717 data takes much of the space in the user's personal computer, user may consider transferring it to 718 external hard drives that offer larger capacity. To allow users to execute their workflows on 719 different HPC environments, users would need to set up their own job submission service and 720 configure the Jupyter environment (e.g., CJW) to the specific HPC environment that they have 721 access to. Although the job submission software used in this study is open source, it is customized 722 for the UIUC HPC used in the study, so it cannot be directly used for other HPCs. Future work 723 could be for CJW to act as a connector to user supplied HPC environments. In this case, CJW 724 would ask users to provide their own credentials and to their own HPC, rather than only using the UIUC HPC service. While not a simple task, standardization of job submission approaches across 725 726 HPC environments makes this functionality possible. Generalizing the approach through future 727 research could benefit users to access their own institutional HPCs and other HPCs at the national 728 level that the user has access to.

729 4 Conclusions

730 The importance of reproducibility is broadly recognized across different scientific disciplines. 731 When it comes to computational hydrology, this can be a significant challenge. This research 732 shows how an architecture that integrates the (1) online data repositories, (2) computational 733 environments, and (3) model API can facilitate reproduction of the components of modern and 734 complex hydrologic studies. For this purpose, we used a recently published large-scale hydrologic 735 study (VB study) as an example. We designed and built cyberinfrastructure that utilized software 736 components to enable intuitive, and online access to computational environments. This approach 737 was used to remove the potential software inconsistencies from users' differing personal software 738 editions, as well as to make implementation easier with pre-compiled software, with the added 739 complication of a computationally expensive research problem instead of a case study. This 740 approach gave the user the option to use either the CJW CG or HPC computational environments, 741 depending on how much they need to reproduce a problem more representative of the big-data 742 problem. Using HydroShare as the data repository, and containerization of the pySUMMA API 743 (with Docker or Singularity in the case of the HPC environment) along with a computational 744 gateway interface of Jupyter notebooks both hosted on the CJW made this possible. Three Jupyter 745 notebooks for the CJW CG environment and four Jupyter notebooks for HPC environment were 746 developed. Notebooks 1-3 for both CJW CG and HPC environments enable, (1) preparing the 747 forcing data, simulation period, and study CAMELS basins, (2) executing SUMMA hydrologic 748 model, and (3) visualization of the results. Notebook 4, only developed for the HPC environment, 749 enables transferring large data from HPC to the scientific cloud service (i.e., CJW) using Globus 750 service integrated by CyberGIS-Compute in a reliable, high-performance and fast way.

We presented a modeling case study subset from the VB study that served as a motivating use case for the cyberinfrastructure. The case study showed how four individual basins with different characteristics can lead to different patterns of temporal aggregation for each of the forcing variables given the same model setup. The case study served to show that the developed cyberinfrastructure enables others to reproduce the VB study for subsets of the original domain as a basis for doing additional research enabling conclusion-reproducibility beyond bitreproducibility.

758 We analyzed performance of the notebooks focusing on contrasting HPC and CJW CG notebooks 759 using a variety of model scenarios. The HPC environments could perform significantly faster 760 simulations compared to CJW CG, enabling users to explore a large number of basins and 761 simulation periods. This clearly showed how the use of HPC from a Jupyter gateway could advance 762 the reproducibility of modern and complex hydrologic studies. The run-time performance analysis for the big data transfer notebook for the HPC environment showed that the method used was 763 764 stable, reliable and fast. Therefore, similar studies could easily benefit from the same approach for 765 transferring large data between scientific cloud services.

766 With the focus of this research was on conclusion-reproducibility over bit-reproducibility of the 767 VB study, users can easily modify the notebooks to test different situations by varying the study 768 basins and periods, parameterizations, and model configurations. These situations highlighted two 769 major challenges. First, the complexity of the big-data problem eventually became large enough 770 that it needed to be run using the HPC computation environment, which presented other smaller challenges of data transfer and portability of the HPC environment. Second, implementation of a 771 772 version control system was needed (e.g., when a user needs to install a newer version of a model 773 API or when a user needs to run these codes on their local machine rather than the used cloud-774 based computational environment). Sharing the dependencies of the computational environments

- as a Python environment yml file and an installation notebook that installs them was discussed as
- a future solution to tackle the version control issue.

777 Finally, as a broader impact, the VB study methodology replicated with interactive codes could 778 also serve as a valuable educational resource, allowing educators to present sophisticated modeling 779 experiments for use within classrooms through online Python notebooks. Likewise, the basic 780 approach could be extended to enable new water decision-support systems that take advantage of 781 the SUMMA framework and HPC yet remain easy to interact with through notebooks. This can 782 help to, for example, evaluate forcing sensitivity to a water resources management objective, or 783 explore the parameter and model uncertainties of SUMMA using different algorithms such as 784 Markov chain Monte Carlo (MCMC), and Bayesian model averaging (BMA) (Samadi et al., 2020) 785 in a systematic manner. With more work to harden and improve the usability of the system 786 presented here, these additional use cases can be possible.

787 Declaration of Competing Interest

788 The authors declare no conflicts of interest relevant to this study.

789 Acknowledgments

790 This work was supported by the National Science Foundation (NSF) under collaborative grants 1664061,

791 1664119 and 1664018 for the development of HydroShare (<u>http://www.hydroshare.org</u>), 1928369 for the

integration of Reproducibility methods into HydroShare. The work also was supported by the Institute for

793 Geospatial Understanding through an Integrative Discovery Environment (I-GUIDE) that is funded by NSF

under award No. 2118329. Any opinions, findings, and conclusions or recommendations expressed in this

795 material are those of the authors and do not necessarily reflect the views of the NSF.

796 Appendix.

797 This section provides supplemental material to support our methods and results. The figures and

tables are referred to in the main text.

799

	PARAMETERS	8 Model Configurati					
FORCINGS 8 Forcing Datasets (F) • NLDAS 1 2 3 4 NLDAS-2 hourly forcings • ppt 1 2 3 4 Precipitation rate • tmp 1 2 3 4 Air temperature • hum 1 2 3 4 Specific humidity • swr 1 2 3 4 Shortwave radiation	11 Parameter sets (P) Default 1 2 3 4 LHS1 2 4 LHS2 2 4 LHS3 2 4 LHS4 2 4 LHS5 2 4 LHS5 2 4 LHS5 2 4 LHS6 2 4 LHS7 2 4	 BIE 3 4 BIL 1 2 3 4 Default BSE 3 4 BSL 3 4 JIE 3 4 JIL 3 4 JSE 3 4 JSL 3 4 	RUNS 4 sets of SUMMA model runs (704 unique runs per basin) 1. DEFAULT 8F x 1P (Default) x 1C (BIL) = 8 r 2. LHS (includes DEFAULT run 8F x 11P (Default & LHS1-10) x 1 3. CONFIG (includes DEFAULT 8F x 1P (Default) x 8C = 64 runs 4. COMPREHENSIVE 8F x 11P (Default & LHS1-10) x 8	runs ns) I C (BIL) = 88 runs I T runs) 3 C = 704 runs			
 <i>Iwr</i> 1 2 3 4 Longwave radiation <i>wnd</i> 1 2 3 4 Wind speed <i>prs</i> 1 2 3 4 Air pressure All F datasets based on the NLD/ hourly values for the designated 	LHS8 2 4 LHS9 2 4 LHS9 2 4 LHS10 2 4 The Latin Hypercube Samp over their defined ranges (T AS-2 hourly data. For datase forcing (3-letter abbreviation	LHS8 2 4 The 3-letter abbreviations define the 2 model decision options for the 3 model decisions tested (Table 1): LHS9 2 4 stomResist: BallBerry (B) or Jarvis (J) snowIncept: lightSnow (I) or stickySnow (s) windPrfile: Exponential (E) or LogBelowCanopy (L) Latin Hypercube Sampling (LHS) sampled 15 parameters 10 times r their defined ranges (Table 3) 2 hourly data. For datasets 2-8, the ng (3-letter abbreviation) are set to					

800

Figure A1. An overview of the forcing datasets (FORCINGS; yellow box), parameter sets
(PARAMETERS; blue box), and model configurations (CONFIGURATIONS; green box) used
in the 704 SUMMA model runs (RUNS; pink box) performed for each of the 671 CAMELS
basins. Note the pink numbers that follow each forcing, parameter, and configuration refers to
the SUMMA model run set as numbered in the pink RUNS box (e.g., the Default parameter set
in the PARAMETERS box is used with SUMMA model runs 1, 2, 3 and 4 in the RUNS box)
(source: modified from Van Beusekom et al., 2022).

JN 1: Preprocessing

Prepares forcings, and sets study basin(s) and simulation period

1_1 Preliminary step

- 1_1_1 Check the environment
- 1_1_2 Import libraries

1_2 Set up paths to SUMMA configuration files for CAMELS basins

1_2_1 Unzip the folder contatining SUMMA CAMELS configurations

1_2_2 Set up paths to SUMMA configuration files

1_3 Select basins and simlaution period

1_3_1 Retrieve the meteorological forcings

- 1_3_2 Select basins and simulation period
- 1_3_3 Slice the forcings to selected basins and simulaiotn period

1_3_4 Slice the SUMMA CAMELS shapefile to selected basins

- 1_3_5 Show the selected basins in map
- 1_3_6 Slice the SUMMA CAMELS paramters abd

attributes files to selected basins

1_3_7 Make constant intiail conditions

1_4 Create forcing files with constant daily values at their daily means

1_4_1 Write and save the truth forcing

- 1_4_2 Shifting to local time zones using longitude values
- 1_4_3 Downsample hourly time-series data to daily data

1_4_4 Upsample back to hourly data and undo time zone changes

1 4 5 Scale constant SW radiation

1_4_6 Create files with only one variable constant at their mean daily values

1_5 Check processed forcing files through plotting

1_5_1 Hourly plots of the forcings

1_5_2 Cummulative plots



810

Figure A2. The preprocessing notebook (JN1) diagram.





Figure A3. Running SUMMA notebook (JN2) diagram.

JN 3: Postprocessing

Explore outputs to find out effect of each fording varible in each basin

3_1 Preliminary steps

3_1_1 Import libraries

3_2 Set up paths to SUMMA configuration files for CAMELS basins

V

3_2_1 Set up paths to SUMMA configuration files

3_3 Make problem complexity choices

- Suggested to choose the most complex problem ran in JN2. DO NOT choose one of these to be "1" here, if you did not choose it or a more complex option to equal "1" in JN2

¥

3_4 Summary statistics of KGE error on outputs

3 4 1 Divide the decision set

3_4_2 Get the forcing and output names, and find the

- HRUs and their locations
- 3_4_3 Summarize KGE error

3_5 Make the results plots

3_5_1 Setup plots

- 3_5_2 The first set of plots
- 3_5_3 The second set of plots

814

Figure A4. Post-processing notebook (JN3) diagram.

816



Figure A5. HPC Data transfer notebook (JN4) diagram.

Table A1. SUMMA output variables chosen for analysis (source: Van Beusekom et al., 2022).

#	Variable Type	SUMMA Variable Name	Description (units)
1		SurfaceRunoff	surface runoff (m s-1)
2		AquiferBaseflow	baseflow from the aquifer (m s-1)
3	liquid water fluxes	Infiltration	infiltration of water into the soil profile (m s-1)
4	for the son domain	RainPlusMelt	rain plus melt (m s-1)
5		SoilDrainage	drainage from the bottom of the soil profile (m s-1)
6		LatHeatTotal	latent heat from the canopy air space to the atmosphere (W m-2)
7	turbulent heat transfer	SenHeatTotal	sensible heat from the canopy air space to the atmosphere (W m-2)
8		SnowSublimation	snow sublimation/frost (below canopy or non-vegetated) (kg m-2 s-1)
9	snow	SWE	snow water equivalent (kg m-2)
10	vegetation	CanopyWat	mass of total water on the vegetation canopy (kg m-2)
11		NetRadiation	net radiation (W m-2)
12	1 1 1	TotalET	total evapotranspiration (kg m-2 s-1)
13	derived	TotalRunoff	total runoff (m s-1)
14		TotalSoilWat	total mass of water in the soil (kg m-2)

Parameter Name	Minimum	Maximum	Default	Constraints
k_macropore	1.0d-7	0.1	0.0001	
k_soil	1.0d-7	1.0d-5	variable	
theta_sat	0.3	0.6	variable	<pre>> critSoilTranspire; > fieldCapacity; > theta_res</pre>
aquiferBaseflowExp	1	10	2.0	
aquiferBaseflowRate	0	0.1	0.1	
qSurfScale	1	100	50	
summerLAI	0.01	10	3	
frozenPrecipMultip	0.5	1.5	1	
heightCanopyTop	0.05	100	variable	> heightCanopyBottom
heightCanopyBottom	0	5	variable	
routingGammaShape	2	3	2.5	
routingGammaScale	1	100000	20000	
albedoRefresh	1	10	1.0	
tempCritRain	272.16	274.16	273.16	
windReductionParam	0	1	0.28	

Table A2. Parameters chosen for Latin Hypercube Sampling (source Van Beusekom et al.,2022).

829

- 830 The eight plots generated by Notebook 3 are described as follows:
- 831 1. Location of the selected CAMELS basin.
- 832 2. KGE values for each CNST forcing dataset (datasets 2-8; Table A2) by output variable
 833 using the DEFAULT model runs. This is a subset of Figure 9A from Van Beusekom et al.
 834 (2022)*.
- 835
 3. Boxplots depicting the range in the KGE values for each set of model runs (DEFAULT, LHS, CONFIG, and COMPREHENSIVE; Table A1) by output variable. Note, boxplots only appear for the model runs selected in Notebook 2. This is a subset of Figure 9B from Van Beusekom et al. (2022).
- 839
 4. Boxplots depicting the range in the KGE values for each set of model runs (DEFAULT, LHS, CONFIG, COMPREHENSIVE; Table A1) by CNST forcing dataset (datasets 2-8; Table A2). Note, boxplots only appear for the model runs executed in Notebook 2. This is a subset of Figure 9C from Van Beusekom et al. (2022).
- Ranks 1 7 stacked barplots depicting the relative basin KGE rank counts by CNST forcing dataset (datasets 2-8; Table A2) for the 14 SUMMA output variables. Note, bars on this

plot will only appear if the COMPREHENSIVE basin runs are executed in Notebook 2.
This is a subset of Figure 8 from Van Beusekom et al. (2022).

- Ranks 1 7 stacked barplots depicting the relative basin KGE rank counts by CNST forcing dataset (datasets 2-8; Table A2) for the eight SUMMA configurations. Note, the complete figure will only appear if the COMPREHENSIVE basin runs are executed in Notebook 2.
 A stacked bar for the default configuration (BIL) will be plotted if the LHS basin runs are executed in Notebook 2. This is a subset of Figure 8 from Van Beusekom et al. (2022).
- 852
 7. Boxplots for each output variable depicting the range in the seven-summed KGE values
 853 (from CNST forcing datasets 2-8) for the eight SUMMA configurations, or for the default
 854 configuration if only the default configuration was run (DEFAULT or LHS basin runs in
 855 Notebook 2. This is a subset of Figure 6 from Van Beusekom et al. (2022).
- 8. Boxplots depicting the range in the summed SUMMA hourly output variables over the 856 857 period of record produced using the benchmark (NLDAS) forcing dataset for the eight 858 SUMMA configuration, or for the default configuration if only the default configuration 859 was run (DEFAULT or LHS basin runs in Notebook 2). Note, a point will appear instead 860 of a boxplot if only the default parameter set was run (DEFAULT or CONFIG basin runs 861 in Notebook 2). This analysis is not in Van Beusekom et al. (2022); it is included in the 862 interactive tool to supply users with potential SUMMA output variable ranges for their 863 selected basin.

* To reproduce the modeling case study presented in the current paper, the selected four
basins need to be specified in Notebook 1 (Figure A2, "Step 1_3_2 Select basins and
simulation period") and then Notebook 3 can be used to reproduce Figure 3 (KGE values
using the DEFAULT model runs for each CNST dataset (datasets 2-8; Table A2), grouped
by SUMMA output variable)

869 **References**

- Bush, R., Dutton, A., Evans, M., Loft, R., Schmidt, G.A., 2021. Perspectives on Data
 Reproducibility and Replicability in Paleoclimate and Climate Science. Harvard Data Sci.
 Rev. 2. https://doi.org/https://doi.org/10.1162/99608f92.00cd8f85
- 873 Chard, K., Tuecke, S., Foster, I., 2016. Globus: Recent enhancements and future plans, in:
 874 Proceedings of the XSEDE16 Conference on Diversity, Big Data, and Science at Scale.
 875 https://doi.org/10.1145/2949550.2949554
- Chen, M., Voinov, A., Ames, D.P., Kettner, A.J., Goodall, J., Jakeman, A.J., Barton, M.C.,
 Harpham, Q., Cuddy, S.M., DeLuca, C., Yue, S., Wang, J., Zhang, F., Wen, Y., Lü, G.,
 2020. Position paper: Open web-distributed integrated geographic modelling and simulation
 to enable broader participation and applications. Earth-Science Rev.
 https://doi.org/10.1016/j.earscirev.2020.103223
- Choi, Y.-D., Beusekom, A. Van, Li, Z., Nijssen, B., Hay, L., Bennett, A., Tarboton, D.,
 Maghami, I., Goodall, J., Clark, M.P., 2022a. Hydrologic Model Sensitivity to Temporal
 Disaggregation of Meteorological Forcing Data across CONUS [WWW Document].
 HydroShare. URL
- 885 https://www.hydroshare.org/resource/c0e8de47aee744d088db7019d78c2b3f/

886 Choi, Y.-D., Beusekom, A. Van, Li, Z., Nijssen, B., Hay, L., Bennett, A., Tarboton, D., 887 Maghami, I., Goodall, J., Clark, M.P., 2022b. SUMMA Simulations using CAMELS 888 Datasets on CyberGIS-Jupyter for Water [WWW Document]. HydroShare. URL 889 https://www.hydroshare.org/resource/50e9716922dc487981b71e2e11f3bb5d/ 890 Choi, Y.-D., Beusekom, A. Van, Li, Z., Nijssen, B., Hay, L., Bennett, A., Tarboton, D., 891 Maghami, I., Goodall, J., Clark, M.P., 2022c. SUMMA Simulations using CAMELS 892 Datasets for HPC use with CyberGIS-Jupyter for Water [WWW Document]. HydroShare. 893 URL https://www.hydroshare.org/resource/9d73d61696ee4f6b9c9a11e21cd44e24/ 894 Choi, Y.-D., Goodall, J., Sadler, J.M., Castronova, A.M., Bennett, A., Li, Z., Nijssen, B., Wang, 895 S., Clark, M.P., Ames, D.P., Horsburgh, J.S., Yi, H., Bandaragoda, C., Seul, M., Hooper, 896 R., Tarboton, D.G., 2021. Toward open and reproducible environmental modeling by 897 integrating online data repositories, computational environments, and model Application 898 Programming Interfaces. Environ. Model. Softw. 135. 899 https://doi.org/10.1016/j.envsoft.2020.104888 900 Clark, M.P., Nijssen, B., Lundquist, J.D., Kavetski, D., Rupp, D.E., Woods, R.A., Freer, J.E., 901 Gutmann, E.D., Wood, A.W., Brekke, L.D., Arnold, J.R., 2015a. The structure for unifying 902 multiple modeling alternatives (SUMMA), Version 1.0: Technical description. 903 Clark, M.P., Nijssen, B., Lundquist, J.D., Kavetski, D., Rupp, D.E., Woods, R.A., Freer, J.E., 904 Gutmann, E.D., Wood, A.W., Brekke, L.D., Arnold, J.R., Gochis, D.J., Rasmussen, R.M., 905 2015b. A unified approach for process-based hydrologic modeling: 1. Modeling concept. 906 Water Resour. Res. 51, 2498–2514. https://doi.org/10.1002/2015WR017198 907 Clark, M.P., Vogel, R.M., Lamontagne, J.R., Mizukami, N., Knoben, W.J., Tang, G., Gharari, S., 908 Freer, J.E., Whitfield, P.H., Shook, K.R., Papalexiou, S.M., 2021. The abuse of popular 909 performance metrics in hydrologic modeling. Water Resour. Res. 57, p.e2020WR029001. 910 https://doi.org/https://doi.org/10.1029/2020WR029001 911 Computational and Information Systems Laboratory, 2017. Chevenne: SGI ICE XA System 912 (NCAR Community Computing). Boulder, CO: National Center for Atmospheric Research; 913 [WWW Document]. https://doi.org/https://doi.org/10.5065/D6RX99HX 914 CyberGIS-Compute Service, 2021. What is CyberGIS-Compute Service? [WWW Document]. 915 URL https://cybergisxhub.cigi.illinois.edu/knowledge-base/components/cybergis-916 compute/what-is-cybergis-compute/ (accessed 2.14.23). 917 CyberGIS Center HydroShare Development Team, 2022. CyberGIS-Jupyter for Water (CJW) 918 Announcements [WWW Document]. URL 919 http://www.hydroshare.org/resource/a901d83a1281404fae58cc41c1cc9889 920 Essawy, B.T., Goodall, J., Voce, D., Morsy, M.M., Sadler, J.M., Choi, Y.-D., Tarboton, D.G., 921 Malik, T., 2020. A taxonomy for reproducible and replicable research in environmental 922 modelling. Environ. Model. Softw. 134. https://doi.org/10.1016/j.envsoft.2020.104753 Essawy, B.T., Goodall, J.L., Xu, H., Rajasekar, A., Myers, J.D., Kugler, T.A., Billah, M.M., 923 924 Whitton, M.C., Moore, R.W., 2016. Server-side workflow execution using data grid

- technology for reproducible analyses of data-intensive hydrologic systems. Earth Sp. Sci. 3,
 163–175. https://doi.org/https://doi.org/10.1002/2015EA000139
- Expanse System Architecture [WWW Document], 2022. . San Diego Super Comput. Cent. URL
 https://www.sdsc.edu/services/hpc/expanse/expanse_architecture.html (accessed 5.6.22).
- Expanse User Guide [WWW Document], 2022. . San Diego Super Comput. Cent. URL
 https://www.sdsc.edu/support/user_guides/expanse.html (accessed 5.6.22).
- Foster, I., 2011. Globus Online: Accelerating and democratizing science through cloud-based
 services. IEEE Internet Comput. 15, 70–73. https://doi.org/10.1109/MIC.2011.64
- Gan, T., Tarboton, D.G., Dash, P., Gichamo, T.Z., Horsburgh, J.S., 2020. Integrating hydrologic
 modeling web services with online data sharing to prepare, store, and execute hydrologic
 models. Environ. Model. Softw. 130.
- 936 https://doi.org/https://doi.org/10.1016/j.envsoft.2020.104731
- Gichamo, T.Z., Sazib, N.S., Tarboton, D.G., Dash, P., 2020. HydroDS: data services in support
 of physically based, distributed hydrological models. Environ. Model. Softw. 125, 104623.
 https://doi.org/https://doi.org/10.1016/j.envsoft.2020.104623
- Gupta, H.V., Kling, H., Yilmaz, K.K., Martinez, G.F., 2009. Decomposition of the mean squared
 error and NSE performance criteria: Implications for improving hydrological modelling. J.
 Hydrol. 377, 80–91. https://doi.org/https://doi.org/10.5194/hess-23-4323-2019.
- Hancock, D.Y., Fischer, J., Lowe, J.M., Snapp-Childs, W., Pierce, M., Marru, S., Coulter, J.E.,
 Vaughn, M., Beck, B., Merchant, N., Skidmore, E., Jacobs, G., 2021. Jetstream2:
 Accelerating cloud computing via Jetstream, in: Practice and Experience in Advanced
 Research Computing (PEARC '21). Association for Computing Machinery, New York, NY,
- 947 USA, p. Article 11, 1–8. https://doi.org/https://doi.org/10.1145/3437359.3465565
- Horsburgh, J.S., Morsy, M.M., Castronova, A.M., Goodall, J., Gan, T., Yi, H., Stealey, M.J.,
 Tarboton, D.G., 2016. HydroShare: Sharing Diverse Environmental Data Types and Models
 as Social Objects with Application to the Hydrology Domain. J. Am. Water Resour. Assoc.
 52. https://doi.org/10.1111/1752-1688.12363
- Hutton, C., Wagener, T., Freer, J., Han, D., Duffy, C., Arheimer, B., 2016. Most computational
 hydrology is not reproducible, so is it really science? Water Resour. Res. 52, 7548–7555.
 https://doi.org/10.1002/2016WR019285
- Knoben, W.J.M., Clark, M.P., Bales, J., Bennett, A., Gharari, S., Marsh, C.B., Nijssen, B.,
 Pietroniro, A., Spiteri, R.J., Tang, G., Tarboton, D.G., Wood, A.W., 2022. Community
 Workflows to Advance Reproducibility in Hydrologic Modeling: Separating modelagnostic and model-specific configuration steps in applications of large-domain hydrologic
 Water Resour. Res. https://doi.org/https://doi.org/10.1029/2021WR031753
- Kurtz, W., Lapin, A., Schilling, O.S., Tang, Q., Schiller, E., Braun, T., Hunkeler, D., Vereecken,
 H., Sudicky, E., Kropf, P., Franssen, H.J.H., 2017. Integrating hydrological modelling, data
 assimilation and cloud computing for real-time management of water resources. Environ.

963 Model. Softw. 93, 418–435. https://doi.org/https://doi.org/10.1016/j.envsoft.2017.03.011 964 Kurtzer, G.M., Sochat, V., Bauer, M.W., 2017. Singularity: Scientific containers for mobility of 965 compute. PLoS One 12, e0177459. https://doi.org/10.1371/journal.pone.0177459 966 Lee, A., 2014. pyDOE [WWW Document]. URL https://pythonhosted.org/pyDOE/ 967 Li, Z., 2021. Dockerfile to create Singularity image for HPC resource [WWW Document]. URL 968 https://github.com/cybergis/cybergis-compute-v2-summa/tree/main/images (accessed 969 5.11.22). 970 Li, Z., Michels, A., Lu, F., Padmanabhan, A., Wang, S., 2022. CyberGIS-Jupyter for Water 971 [WWW Document]. HydroShare. URL 972 http://www.hydroshare.org/resource/4cfd280e8eb747169b293aec2862d4f5 973 Lyu, F., Yin, D., Padmanabhan, A., Choi, Y.-D., Goodall, J.L., Castronova., A.M., Tarboton, 974 D.G., Wang, S., 2019. Reproducible hydrological modeling with CyberGIS-Jupyter: a case 975 study on SUMMA, in: Proceedings of the Practice and Experience in Advanced Research 976 Computing on Rise of the Machines (Learning). pp. 1–6. 977 https://doi.org/https://doi.org/10.1145/3332186.3333052 978 Mathevet, T., Michel, C., Andréassian, V., Perrin, C., 2006. A bounded version of the Nash-979 Sutcliffe criterion for better model assessment on large sets of basins, in: Large Sample 980 Basin Experiments for Hydrological Model Parameterization: Results of the Model 981 Parameter Experiment-MOPEX. IAHS PUBLICATION, Wallingford, UK, pp. 211-219. 982 Melsen, L. A., Torfs, P.J.J.F., Uijlenhoet, R., Teuling, A.J., 2017. Comment on "Most 983 computational hydrology is not reproducible, so is it really science?" by Christopher Hutton 984 et al. Water Resour. Res. 53, 2568–2569. https://doi.org/10.1002/2016WR020208 985 Merkel, D., 2014. Docker: lightweight Linux containers for consistent development and 986 deployment. Linux j 239, no. 2. Mizukami, N., Wood, A., 2021. NLDAS Forcing NetCDF using CAMELS datasets from 1980 to 987 988 2018 [WWW Document]. HydroShare. URL 989 http://www.hydroshare.org/resource/a28685d2dd584fe5885fc368cb76ff2a 990 Mullendore, G.L., Mayernik, M.S., Schuster, D.C., 2021. Open Science Expectations for 991 Simulation-Based Research. Front. Clim. 3. 992 https://doi.org/https://doi.org/10.3389/fclim.2021.763420 993 Newman, A.J., Clark, M.P., Craig, J., Nijssen, B., Wood, A.W., Gutmann, E.D., Mizukami, N., 994 Brekke, L., Arnold, J.R., 2015a. Gridded ensemble precipitation and temperature estimates 995 for the contiguous United States. J. Hydrometeorol. 16, 2481–2500. 996 https://doi.org/https://doi.org/10.1175/JHM-D-15-0026.1 997 Newman, A.J., Clark, M.P., Sampson, K., Wood, A., Hay, L.E., Bock, A., Viger, R.J., Blodgett, 998 D., Brekke, L., Arnold, J.R., Hopson, T., Duan, Q., 2015b. Development of a large-sample 999 watershed-scale hydrometeorological data set for the contiguous USA: Data set 1000 characteristics and assessment of regional variability in hydrologic model performance.

- 1001Hydrol. Earth Syst. Sci. 19, 209–223. https://doi.org/https://doi.org/10.5194/hess-19-209-10022015, 2015
- 1003 NLDAS-2, 2014. NLDAS-2 Forcing Dataset Information [WWW Document]. URL
 1004 https://ldas.gsfc.nasa.gov/nldas/v2/forcing (accessed 4.9.21).
- 1005 OPeNDAP, 2021. OPeNDAP User Guide [WWW Document]. URL
 1006 https://www.earthdata.nasa.gov/opendap-user-guide (accessed 8.12.21).
- Samadi, S., Pourreza-Bilondi, M., Wilson, C.A.M.E., Hitchcock, D.B., 2020. Bayesian Model
 Averaging With Fixed and Flexible Priors: Theory, Concepts, and Calibration Experiments
 for Rainfall-Runoff Modeling. J. Adv. Model. Earth Syst. 12.
 https://doi.org/10.1029/2019MS001924
- Simmonds, M., Riley, W.J., Agarwal, D., Chen, X., Cholia, S., Crystal-Ornelas, R., Coon, E.,
 Dwivedi, D., Hendrix, V., Huang, M., Jan, A., 2022. Guidelines for Publicly Archiving
 Terrestrial Model Data to Enhance Usability, Intercomparison, and Synthesis. Data Sci. J.
 21.
- Stewart, C.A., Cockerill, T.M., Foster, I., Hancock, D., Merchant, N., Skidmore, E., Stanzione,
 D., Taylor, J., Tuecke, S., Turner, G., Vaughn, M., Gaffney, N.I., 2015. Jetstream: a selfprovisioned, scalable science and engineering cloud environment, in: Proceedings of the
 2015 XSEDE Conference: Scientific Advancements Enabled by Enhanced
 Cyberinfrastructure. St. Louis, Missouri, p. ACM: 2792774. p. 1-8.
- 1020 https://doi.org/https://dx.doi.org/10.1145/2792745.2792774
- Tarboton, D., Calloway, C., 2021. THREDDS DAP2 [WWW Document]. URL
 http://www.hydroshare.org/resource/70070fa1b382496e85ca44894683b15d
- Tarboton, D.G., Idaszak, R., Horsburgh, J.S., Ames, D.P., Goodall, J.L., Band, L.E., Merwade,
 V., Couch, A., Hooper, R.P., Maidment, D.R., Dash, P.K., 2014. HydroShare: Advancing
 collaboration through hydrologic data and model sharing, in: Proceedings of the 7th
 International Congress on Environmental Modelling and Software. Int. Environ. Modell.
 and Software Soc, San Diego, Calif., pp. 978–988.
- Towns, J., Cockerill, T., Dahan, M., Foster, I., Gaither, K., Grimshaw, A., Hazlewood, V.,
 Lathrop, S., Lifka, D., Peterson, G.D., Roskies, R., Scott, J.R., Wilkins-Diehr, N., 2014.
 XSEDE: Accelerating Scientific Discovery. Comput. Sci. Eng. 16, 62–74.
 https://doi.org/https://dx.doi.org/10.1109/MCSE.2014.80
- 1032 Van Beusekom, A.E., Hay, L.E., Bennett, A.R., Choi, Y.-D., Clark, M.P., Goodall, J., Li, Z.,
 1033 Maghami, I., Nijssen, B., Wood, A.W., 2022. Hydrologic Model Sensitivity to Temporal
 1034 Aggregation of Meteorological Forcing Data: A Case Study for the Contiguous United
 1035 States. J. Hydrometeorol. 23, 167–183. https://doi.org/10.1175/JHM-D-21-0111.1
- 1036 Virtual Roger User Guide [WWW Document], 2022. . CyberGIS Cent. Adv. Digit. Spat. Stud. 1037 Univ. Illinois Urbana Champaign. URL https://cybergis.illinois.edu/infrastructure/hpc-user 1038 guide/ (accessed 6.5.22).

Yang, C., Raskin, R., Goodchild, M., Gahegan, M., 2010. Geospatial cyberinfrastructure: past, present and future. Computers, Environment and Urban Systems. Comput. Environ. Urban Syst. 34, 264–277. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2010.04.001
Yin, D., Liu, Y., Padmanabhan, A., Terstriep, J., Rush, J., Wang, S., 2017. A CyberGIS-Jupyter

- 1043 Framework for Geospatial Analytics at Scale, in: Proceedings of the Practice and
- 1044 Experience in Advanced Research Computing 2017 on Sustainability, Success and Impact.
- 1045 pp. 1–8. https://doi.org/https://doi.org/10.1145/3093338.3093378