



1 **WaterSoftHack Cybertraining: Reproducible Data Science, Machine Learning, and Cloud and**
2 **Edge Computing Training for Collaborative Water Science Research**

3 Krishna Panthi¹, Vidya Samadi^{2,3*}, Nima Zafarmomen², Mostafa Saberian⁴, Adarsha Neupane², Carlos
4 Erazo Ramirez⁵, Bijaya Adhikari⁶, Anthony Castronova⁷, Ibrahim Demir^{5,8}

- 5 1. School of Computing, Clemson University, Clemson, SC, USA
6 2. Department of Agricultural Sciences, Clemson University, Clemson, SC, USA
7 3. Artificial Intelligence Research Institute for Science and Engineering (AIRISE), School of Computing,
8 Clemson University, SC, USA
9 4. The Glenn Department of Civil Engineering, Clemson University, Clemson, SC, USA
10 5. ByWater Institute, Tulane University, New Orleans, USA
11 6. Department of Computer Science, University of Iowa, IA, USA
12 7. Consortium of Universities for the Advancement of Hydrologic Science Inc, MA, USA
13 8. River-Coastal Science and Engineering Department, Tulane University, New Orleans, USA

14 *corresponding author: samadi@clemson.edu

17 **Highlights**

- 18 1. Hands-on cybertraining in data science, machine learning for water science.
19 2. Open-source WaterSoft Python package enables reproducible workflow development.
20 3. Participants show skill gains, collaboration, and applied hydrological research impact.

22 **Abstract**

23 The growing complexity and volume of data in water science demand advanced computational skills among
24 researchers, yet significant barriers limit rapid skill acquisition. We present WaterSoftHack, a two-week
25 cybertraining program designed to equip students, early-career professionals, and researchers with
26 reproducible data science, machine learning, and cloud/edge computing skills for water science applications.
27 The program combines open-access training resources, including the WaterSoft Python package, with
28 cohort-based capstone projects. A rigorous selection process of interested candidates ensures diverse
29 participation across academic levels, institutions, and backgrounds. The training integrates immersive,
30 hands-on instruction with formal science communication, emphasizing reproducibility, scalability, and
31 teamwork. Drawing on surveys and qualitative interviews from the first two years, we demonstrate notable
32 skill advancement, collaborative synergy, and career advancement outcomes. WaterSoftHack highlights the
33 importance of project-based, integrated cybertraining in building computational capacity and preparing a
34 diverse, capable workforce for the data-driven future of water science and engineering.

36 **Keywords:** WaterSoftHack; Cybertraining; Project-based Learning; Workforce Development.

38 **Short Summary**

39 WaterSoftHack is a cybertraining program that bridges water science and modern computational methods
40 through hands-on training in data analytics, machine learning, and cloud and edge computing. Using open-
41 source tools and project-based learning, the program builds reproducible workflows, strengthens
42 interdisciplinary collaboration, and enhances participants' research capacity and career readiness.

43



44 1. Introduction

45 Traditional hydrological modeling has long relied on conceptual and physically based models that represent
46 dominant watershed processes, such as rainfall–runoff dynamics and storage mechanisms (Seibert and
47 Bergström, 2022).. Early pioneers such as Bjerknes (1904) envisioned weather prediction as solving
48 differential equations for fluid dynamics and thermodynamics, while Richardson (1922) articulated a
49 numerical step-by-step procedure for weather forecasting, which remained impractical for decades. It was
50 Charney et al. (1950) who proved the feasibility of digital numerical weather prediction (NWP) experiments
51 via experiments on the ENIAC machine. In parallel, hydrological process modeling advanced with the
52 Saint-Venant equations (Saint-Venant, 1871) for open channel flow and later the kinematic wave theory of
53 Lighthill & Whitham (1955), which provided efficient formulations for flood and storm runoff routing.
54 Empirical rainfall-runoff models like unit hydrographs (Sherman, 1932) and more general instantaneous
55 unit hydrographs (Nash, 1957) further shaped the numerical modeling of early operational hydrology. By
56 the late twentieth century, integrated physically based frameworks such as Modelling Integrated Kinetic
57 Environment – Système Hydrologique Européen (SHE/MIKE-SHE; Abbott, et al., 1986), land-surface
58 hydrology models like Variable Infiltration Capacity (VIC; Liang et al., 1994), and terrain-driven
59 conceptualizations like TOPography-based hydrological MODEL (TOPMODEL; Beven & Kirkby, 1979)
60 broadened the modeling toolkit from catchment to continental scales.

61 Alongside modeling advances, the acquisition of hydrologically relevant data transformed radically.
62 Gauges provided the first temporally resolved precipitation and discharge records, while the radar networks
63 provide near-real-time rainfall estimates at kilometer and minute scales (Seo et al., 2019). Satellite missions
64 such as Tropical Rainfall Measuring Mission (TRMM; Kummerow et al., 1998) and Global Precipitation
65 Measurement (GPM; Hou et al., 2014) offer near-global precipitation at a few hourly resolutions. In addition,
66 Soil Moisture Active Passive (SMAP; Entekhabi et al., 2010) retrieves global soil moisture and Gravity
67 Recovery and Climate Experiment (GRACE; Tapley et al., 2004) reveals total water storage changes as an
68 effect to earth’s gravity. Reanalysis products such as the National Centers for Environmental Prediction
69 (NCEP)/ National Center for Atmospheric Research (NCAR) (Kalnay et al., 2018) and Fifth Generation of
70 ECMWF Atmospheric Reanalysis of the Global Climate (ERA5; Hersbach et al., 2020) provide consistent,
71 gridded meteorological forcing fields at hourly resolution. The data landscape available to hydrologists
72 today is richer than ever, spanning ground-based, airborne, and spaceborne observations (e.g., Demir et al.,
73 2015).

74 However, this abundance of data introduces a new challenge. Classical hydrology curricula remain
75 primarily focused on physical process theory and mathematical formulations, while training in the
76 computational cyberinfrastructure required exploiting modern computational skills which are limited in
77 hydrology. Students in hydrology learn the governing equations of infiltration, evapotranspiration, or
78 groundwater flow, but rarely acquire skills in managing terabyte-scale datasets, developing reproducible
79 workflows, or deploying scalable models on cloud or edge platforms. This mismatch is becoming
80 increasingly critical as hydrology transitions into data- and computation-rich science.

81 The gap can be described as threefold. First, although hydrological data from satellites, radars, and
82 reanalysis are openly available, researchers often lack the appropriate skills to acquire, preprocess, and
83 integrate these heterogeneous datasets. Second, while data science and machine learning have exploded in
84 Earth sciences (Krajewski et al., 2021), they are often encountered by hydrology students only superficially,



85 without guidance on their strengths, limitations, and integration with physical process understanding. Third,
86 modern hydrological prediction and water resource management increasingly require high-performance or
87 cloud/edge-based computational infrastructure; for example, collecting sensor monitoring data for running
88 distributed hydrology models such as Weather Research and Forecasting – Hydrological modeling system
89 (WRF-Hydro; Gochis et al., 2015) or producing ensemble forecasts (e.g., Demargne et al., 2014) for
90 understanding flood risk. However, most hydrology curricula lack advanced computing and do not train
91 students to use these models.

92 This gap has tangible consequences. Operational agencies such as the U.S. National Weather Service rely
93 on ensemble forecasting systems running on supercomputers and assimilating multiple streams of satellite
94 and in situ data (Brown et al., 2014). Research repositories such as NASA’s Earth data distribute petabyte-
95 scale archives that are impractical to download, requiring direct cloud/edge-based analysis instead. Without
96 training in data acquisition frameworks, data science methods, and modern computing platforms, new
97 generation of hydrologists risk being unprepared to fully exploit available resources or to contribute
98 effectively to operational and policy-relevant decision-making.

99 To address this gap, targeted training initiatives are needed to explicitly integrate computational methods
100 into water science and engineering education. Such programs should cover the modern computational stack
101 including (1) data acquisition and management for hydrological datasets, (2) data science principles
102 including statistical learning and visualization, (3) machine learning algorithms for both image and time
103 series pattern discovery and hybrid modeling, and (4) deployment on high-performance, cloud/edge
104 infrastructures. Most importantly, these approaches should be taught in a way that connects directly to
105 hydrological applications and emphasizes reproducibility, scalability, and collaborative practice.

106 Training workshops and hackathons have emerged as a way of bridging this gap. Hackathons provide a
107 project-based learning opportunity. They present challenges to participants that mirror real life scenarios,
108 often requiring skills not taught in academia (Blumenfeld et al., 1991). Hackathons have recently emerged
109 as innovative platforms for skill development, fostering collaboration, problem-solving, and rapid
110 prototyping (Haw and Crawford, 2025). Research by Briscoe (2014) and Lara & Lockwood (2016) show
111 that hackathons are effective in promoting interdisciplinary teamwork and accelerated learning. Similarly,
112 training workshops have been shown to be pivotal in fostering continuous learning and skill development
113 in hydrology, as they provide immersive, hands-on experiences with emerging methods, tools, and
114 interdisciplinary practices that help professionals and students stay current in this rapidly evolving field
115 (Huppenkothen et al., 2017). Indeed, hackathons with dedicated data-focused tracks provide participants
116 with valuable experience in handling real-world datasets and addressing common issues such as data
117 reliability, error assessment, and reproducibility (Anslow et al., 2016). Henceforth, participants are
118 challenged to apply newly acquired skills to water-related projects, ensuring that the training directly
119 translates to practical applications.

120

121 The WaterSoftHack program is an innovative hackathon-style cybertraining initiative that seamlessly
122 integrates data science, machine learning, and cyberinfrastructure principles. By combining hands-on
123 exercises with project-based learning, the program equips researchers with practical skills in data-driven
124 analytics and a deep understanding of the supporting cyberinfrastructure necessary for modern water
125 research. This study presents the outcomes of WaterSoftHack cybertraining program, combining one week



126 of theoretical and hands-on workshops with one week of capstone projects, during which participants
127 tackled real hydrological problems using state-of-the-art data and computational tools.

128

129 Participants included undergraduate and graduate students and early-career researchers with diverse
130 backgrounds in hydrology, many of whom had limited prior exposure to computational methods. By
131 documenting the program's structure, delivery, and evaluation, we aim to contribute to the broader
132 discussion on training the next generation of water research scientists. Specifically, we: (a) identify key
133 areas of computational knowledge currently underrepresented in hydrology education; (b) demonstrate how
134 intensive short-term training can enhance participants' competencies in data acquisition, data science, and
135 cloud/edge-based computation; and (c) highlight the challenges encountered and lessons learned in
136 designing interdisciplinary training that bridges hydrological science with modern computational practice.

137

138 This paper is organized as follows: Section 2 describes the WaterSoftHack program, including participant
139 recruitment and curriculum development. Section 3 outlines the project design and evaluation methodology.
140 Section 4 presents the results and discussion, and Section 5 concludes the paper with key findings and
141 insights.

142

143 **2. WaterSoftHack Program Description**

144 WaterSoftHack was conceived as a cybertraining hackathon aimed at accelerating the adoption of data
145 analytics, machine learning, and cloud/edge-based workflows in water science. Its underlying theory of
146 change posits that structured, mentored, hands-on exposure to reproducible computational methods would
147 (i) lower barriers to adoption, (ii) catalyze open-source contributions, and (iii) foster cross-disciplinary
148 collaboration among water science researchers. The program achieves this by providing accessible
149 instruction (educational empowerment), enabling open-source workflow development (research
150 innovation), cultivating a workforce proficient in advanced cyberinfrastructure (collaborative synergy), and
151 promoting broad community adoption of these approaches (community engagement).

152

153 **2.1. Program Organization and Recruitment**

154 In WaterSoftHack cybertraining, instruction was delivered by faculty and researchers with expertise in
155 hydroinformatics, data science, and computer science, while graduate students mentored and facilitated
156 small-group activities and provided technical support. Eligible participants included undergraduates,
157 graduate students, postdoctoral researchers, and early-career faculty, selected from a large and diverse
158 applicant pool. Recruitment opened in February via social media (e.g., LinkedIn) and professional channels
159 such as the Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI)
160 website and newsletters. Two months window was provided for applicants to fill in the application form.
161 After the application deadline, applicants were collectively evaluated based on technical experience,
162 training needs, alignment with program goals, and overall potential. The program selected a coherent of
163 WaterSoftHack Fellows each year. We intentionally sought diversity across different criteria when selecting
164 the Fellows. The criteria included institution types (very high research activity (R1), high research activity
165 (R2) , Primarily Undergraduate Institution (PUI), Hispanic-Serving Institution (HIS), and Historically
166 Black College or University (HBCU)), academic roles, gender, disability status, and disciplinary
167 backgrounds. In Year 1, 15 Fellows were selected to participate and compete in the hackathon. This number
168 was increased to 18 in Year 2 (2025), with several Year-1 Fellows continuing. The program attracted a
169 broad set of audiences, with over 150 general participants in 2024 and over 200 in 2025. The selected



170 fellows conducted capstone projects and competed for a prize, forming the core of the WaterSoftHack
171 Fellows program.

172

173 **2.2. Program Structure and Curriculum Development**

174 The training program was designed as an intensive, two-week online event hosted in late July and early
175 August. It was delivered synchronously via Zoom, with a daily schedule from 11:00 to 15:00 U.S. Eastern
176 Time, incorporating a one-hour lunch break. This format provided approximately 30 contact hours,
177 emphasized hands-on practical training (60%) combined with theoretical lectures (40%). To ensure
178 accessibility and long-term utility, all teaching materials, datasets, and code repositories were made openly
179 available on GitHub. A dedicated Discord server was established to facilitate asynchronous technical
180 support and peer-to-peer communication. All training sessions were recorded and publicly archived on the
181 project's YouTube channel to serve as a persistent educational resource.

182 In each year of the project, educational content was themed to introduce emerging technologies, for example
183 2024 training involved utilizing web-based tools for hydrologic analysis. Similarly, 2025 involved running
184 machine learning models via Google's Collab, and leveraging dedicated high performance computing
185 resources provided by the US National Science Foundation's (NSF) National AI Research Resource
186 (NAIRR) program.

187 **Week 1 (Open Training):** The start of each event consisted of theoretical and practical training that was
188 open to all the participants who applied to the program. An orientation was conducted in the first hour of
189 day 1 for the WaterSoftHack Fellows with general introduction of the program as well as code of conduct
190 and expectations. The first half of each day focused on scientific theory in formal lectures, while the second
191 half was structured as hands-on practical training with code implementation demonstration, model
192 development and data processing with real datasets. The final day of the week was restricted to the Fellows
193 and featured an interactive science communication seminar from Alan Alda Center at Stony Brook
194 University. This provided essential training to improve messaging and facilitate team cohesion, serving as
195 a preparatory exercise for team formation ahead of the capstone projects.

196

197 **Week 2 (Capstone Projects):** This week was limited to U.S.-based researchers and included a subset of
198 week 1 participants selected as WaterSoftHack Fellows. Fellows formed interdisciplinary research teams
199 by aligning their expertise and research interests around a central project idea. Program organizers ensured
200 that each capstone team included a balance of junior and senior researchers and maintained a 50/50 gender
201 balance. Teams then developed projects that applied the techniques covered during the training to extract
202 actionable insights from real-world datasets. All activities were conducted under the supervision of faculty
203 mentors and technical trainers. The week concluded with team presentations, a moderated Q&A session,
204 and submission of a short position paper. Projects were evaluated by program judges using a standardized
205 rubric emphasizing technical rigor, reproducibility, and potential scientific and societal impact, with the
206 top-scoring project receiving an award.

207

208 **3. Evaluation Design and Methods**

209 The development and evaluation of WaterSoftHack cybertraining adhered to a design-based research
210 (DBR) methodology (Design-Based Research Collective, 2003), emphasizing iterative refinement based on
211 empirical evidence and participant feedback. A preliminary need assessment survey was conducted through
212 CUAHSI within the water science community to identify specific training needs and priorities, which



213 directly informed curriculum design. Instructional materials were collaboratively peer-reviewed and made
214 publicly available on GitHub prior to the cybertraining, ensuring transparency and promoting community
215 engagement.

216 All procedures conformed to ethical research standards. Participants provided informed consent. All data
217 were de-identified and securely stored, and the study obtained Institutional Review Board (IRB)
218 determination prior to analysis. Recognizing limitations such as small cohort sizes and potential self-
219 selection bias, the analysis incorporated triangulation of multiple data sources and emphasized effect sizes,
220 and non-parametric methods to ensure robustness and interpretive rigor. This evaluation framework
221 provides a systematic and ethically grounded approach for assessing program effectiveness, emphasizing
222 both accountability and continuous improvement. The following section presents the results derived from
223 this evaluation, including outcomes of the needs assessment, participant learning gains, and thematic
224 insights from post-cybertraining feedback.

225

226 **4. Results**

227 **4.1. WaterSoftHack Fellow Selections**

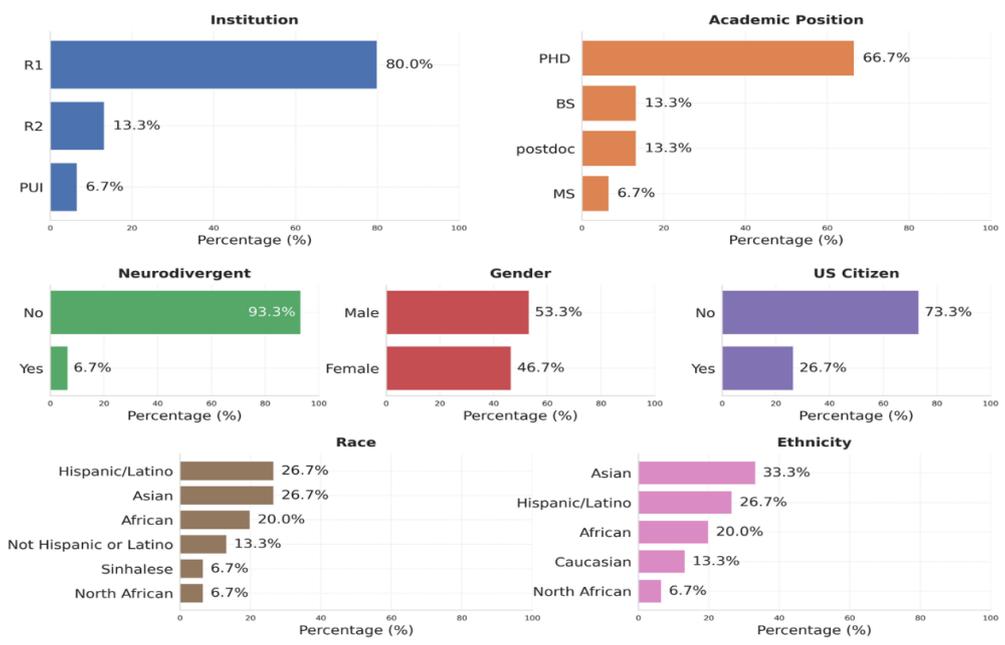
228 The selection process for the fellowship program was highly competitive. In its inaugural year, 2024, a
229 cohort of 15 fellows was selected from a pool of 109 US-based applicants. The program continued in 2025
230 with 130 US-based applications, from which 12 new fellows were chosen. Six of the Fellows from the 2024
231 cohort continued the WaterSoftHack program in 2025, resulting in a total of 18 fellows for the 2025
232 program. WaterSoftHack attracted a multidisciplinary group of participants with academic backgrounds
233 spanning multiple disciplines including civil engineering-water resources, environmental engineering,
234 geology, earth science, and computer science. The cohort was intentionally composed of individuals with
235 varying levels of expertise, encompassing both newcomers to the field and those with prior exposure to
236 data science methodologies.

237

238 Financial support for the selected Fellows was provided through the NSF CyberTraining program. This
239 funding was instrumental in ensuring the program's accessibility and facilitating broad participation from a
240 diverse pool of candidates. The selection criteria placed a strategic emphasis on fostering a diverse and
241 inclusive environment. Priority was given to applicants from institutions with low cyberinfrastructure
242 adaptation, with a particular focus on those from R2, HBCU, and PUI institutions. The final cohort was
243 deliberately structured to include a representative mix of academic positions, including undergraduate
244 students, master's graduate students, postdoctoral scholars, and early-career faculty. Furthermore, the
245 selection process was committed to enhancing diversity across multiple dimensions, including gender,
246 ethnicity, race, and disability status. A comprehensive demographic breakdown of the participants detailing
247 their institutional affiliations, academic positions, ethnicity, disability, gender, race, and U.S. citizenship
248 status is presented in Figures 1 and 2 for years 2024 and 2025, respectively. As illustrated in Figure 1, the
249 program attracted 33.3% Asian, 26.7% of Hispanic, 20% African, 13.3% Caucasian, and 6.7% North
250 African. The ethnicities of participants in 2025 include 27.8% Caucasian, 22.2% Asian, 16.7% African,
251 11.1% North African, 11.1% Asian American, 5.6% Hispanic, and 5.6% Native American.

252

253

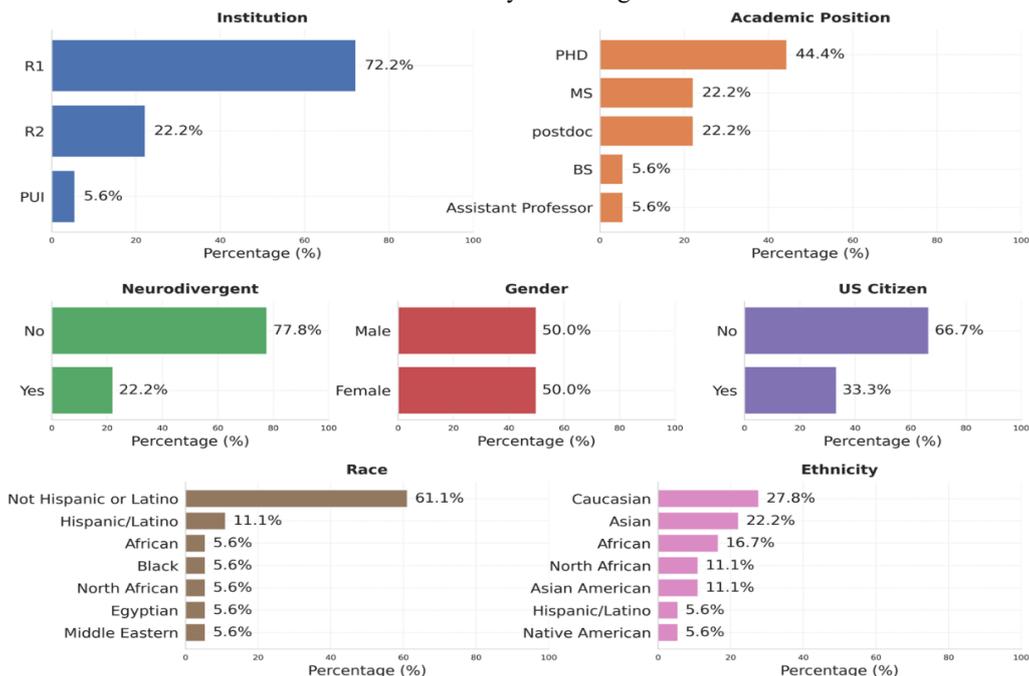


254

255

256

Figure 1. Breakdown of participant details for the WaterSoftHack Fellows that were selected (n=15) for the 2024 cybertraining.



257

258

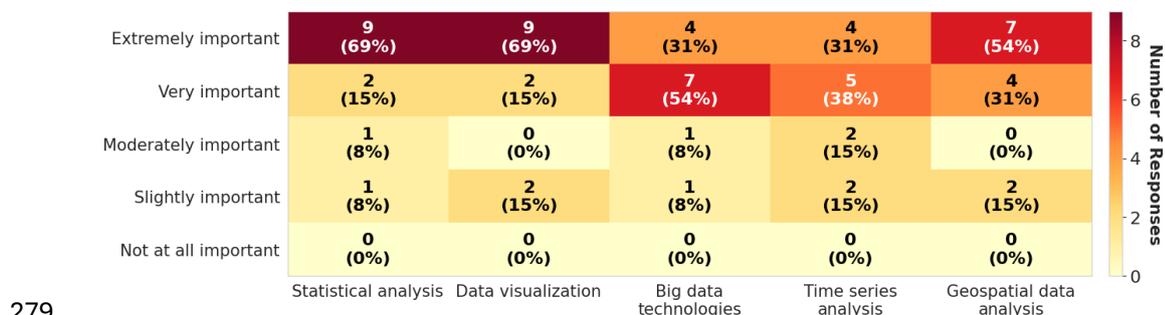
259

Figure 2. Breakdown of participation details for the WaterSoftHack Fellows selected (n=18) for the 2025 cybertraining.



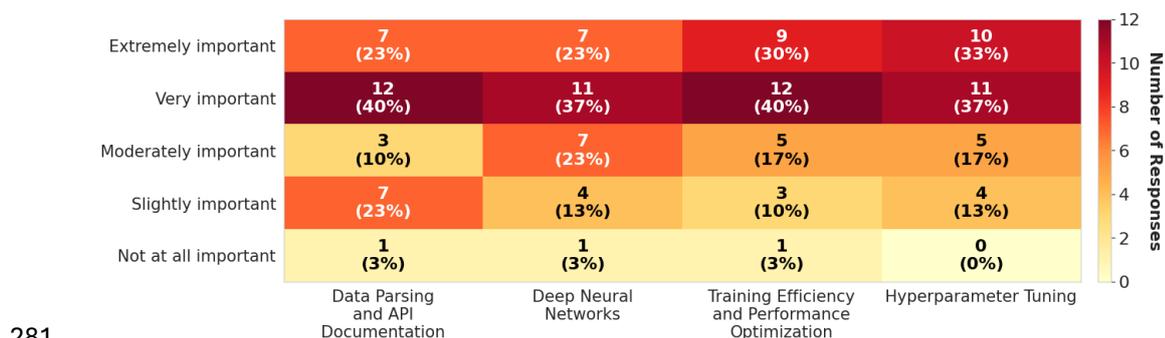
260 **4.2. Need Assessment Survey**

261 The WaterSoftHack program was directly informed by a pre-program need assessment survey designed to
 262 ensure alignment with workforce and community needs. The survey was distributed through CUAHSI to
 263 the broader water science community to systematically identify existing knowledge gaps in data science,
 264 machine learning, and cyberinfrastructure applications within hydrology. Findings from the survey guided
 265 the selection of training topics, hands-on activities, and capstone project themes, ensuring that the program
 266 addressed high-priority skills required for modern hydrologic research and practice. The survey results
 267 indicated a significant need for training in fundamental cyberinfrastructure and machine learning literacy
 268 and the development of reproducible computational workflows. The findings from 2024 and 2025 need
 269 assessments are summarized in Figures 3 and 4, respectively. The 2024 survey (n = 13) focused on
 270 identifying skill gaps in data analysis and visualization, revealing strong demand for training in statistical
 271 analysis, data visualization, big data technologies, time series analysis, and geospatial data processing
 272 across both academic and industry sectors. Similarly, the 2025 survey (n = 30) examined gaps in machine
 273 learning competencies, highlighting high demand for skills in deep neural networks, hyperparameter tuning,
 274 training efficiency, performance optimization, data parsing, and web-based scraping tool design and
 275 implementation. Statistical analysis and data visualization were identified as the most important training
 276 skills in 2024 while data parsing and machine learning training efficiency and performance optimization
 277 were viewed as the most critical skills to learn in 2025.
 278



279

280 Figure 3. Training need survey results for 2024 with 13 participants completing the survey.



281

282 Figure 4. Training need survey results for 2025 with 30 participants completing pre-event survey.



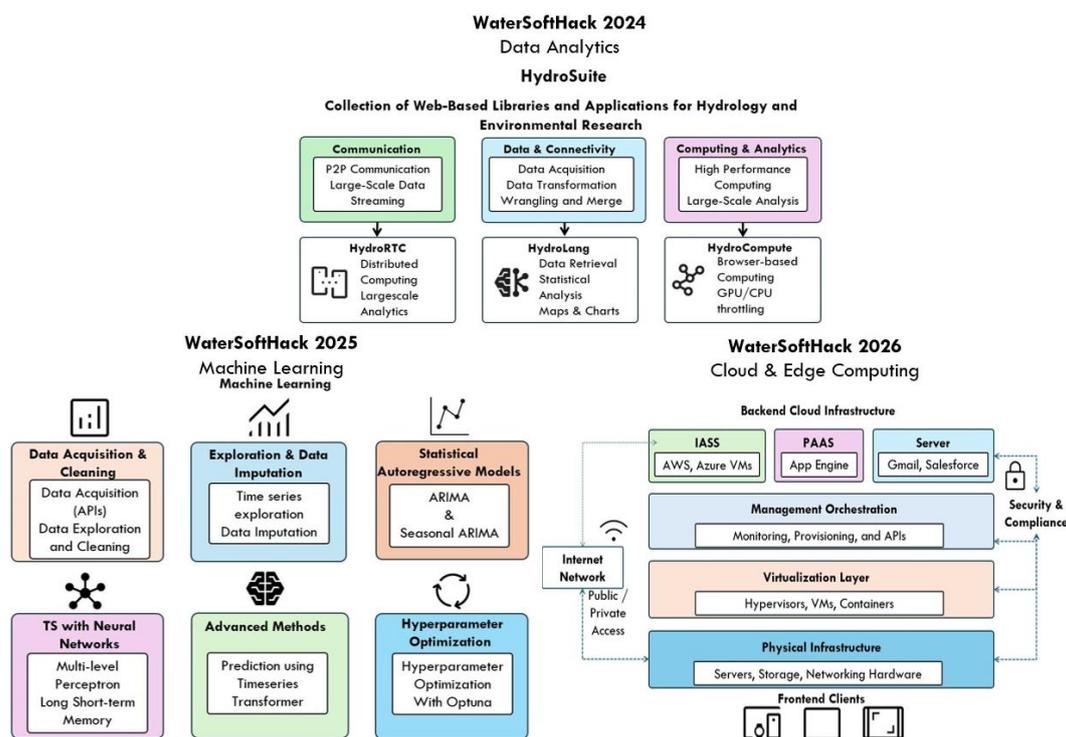
283 The WaterSoftHack cybertraining materials and curriculum evolved around three core cybertraining
284 components, data analytics and visualization, machine learning, and cloud/edge computing which were
285 implemented as modular WaterSoftHack training packages and assembled into the WaterSoft Python
286 package. WaterSoft package is hosted as a public monorepository on GitHub, containing training modules,
287 source codes, related datasets, and capstone projects developed by WaterSoftHack Fellows. The details of
288 these training components and the structure of the monorepository are discussed in the following sections.

289

290 **4.3. WaterSoft Python Package**

291 The preparatory materials provided the foundational knowledge necessary for participants to engage
292 effectively in cybertraining, covering core programming concepts and data science techniques. As of
293 February 2026, the WaterSoft monorepository on GitHub hosts a comprehensive suite of resources,
294 including HydroSuite (Hydrology Software Suite) training modules, data visualization tools, machine
295 learning workflows, hyperparameter tuning models, and capstone projects developed by the WaterSoftHack
296 Fellows. These components collectively form the foundation of the WaterSoft training framework, designed
297 to promote reproducible, data-driven hydrologic research. By the conclusion of the next WaterSoftHack
298 cycle (summer 2026), a dedicated cloud/edge computing training module will be integrated into the
299 WaterSoft package, further expanding its scope. Figure 5 presents a block diagram illustrating the key
300 components and interconnections within the WaterSoft framework. This includes three different types of
301 workflows: data analytics and visualization, machine learning, and cloud and edge computing.

302 Beyond its computational training components, the WaterSoft Python package also incorporates modules
303 that directly support hydrological science applications. The package provides tools and workflows for
304 accessing and processing hydrologic and meteorological datasets, performing time-series analysis of
305 hydrological observations, and developing simulation models for hydrologic variables such as streamflow,
306 river stage, and water-quality indicators. For example, training modules within the HydroSuite framework
307 demonstrate how machine learning and statistical methods can be applied to hydrological forecasting
308 problems using observational datasets from sources such as stream gauges and meteorological stations.
309 These workflows include data acquisition through APIs, preprocessing and quality control of hydrological
310 time-series data, feature engineering for watershed variables, hyperparameter tuning, and the development
311 of predictive models for hydrologic processes.



312

313 Figure 5. The workflow of the WaterSoft Python package which integrates three primary components -
 314 data analytics and visualization, machine learning, and cloud and edge computing to enable scalable,
 315 data-driven water science research and training.

316 4.3.1. HydroSuite Training Module

317 Year 1 (2024) of WaterSoftHack involved a comprehensive introduction to HydroSuite and HydroLang
 318 (Hydrology + Language; Ramirez et al., 2022), HydroCompute (Hydrology + Compute; Ramirez et al.,
 319 2024a), HydroRTC (Hydrology+Real-Time Communication; Ramirez et al., 2024b), and HydroSuite AI
 320 Helper (Pursnani et al., 2024). The pedagogical approach for the module was predicated on the principles
 321 of accessibility and immediate applicability. To remove the need for complex local software configuration,
 322 the instruction was delivered via online, open-access sandbox environments. This methodology ensured a
 323 standardized and universally accessible platform, allowing participants to execute and modify code directly
 324 within a web browser, thereby focusing on data analytics concepts rather than environment setup.
 325 HydroSuite contained three different modules: The initial phase of the training centered on HydroSuite, a
 326 JavaScript framework for client-side data acquisition, analysis, and visualization. The curriculum was then
 327 continued around computing four other modules including HydroLang, HydroCompute, HydroRTC, and
 328 HydroSuite AI Helper. Data module component instructed participants on programmatic interaction with
 329 environmental data. The tutorials covered methods for data retrieval from diverse sources with public
 330 application programming interfaces (APIs) and model repositories. Participants learned techniques for data
 331 manipulation and transformation, as well as client-side data input/output operations.



332 Additionally, the training introduced multiple computational capabilities of HydroLang including a hydro
333 component containing functions for fundamental hydrological analyses like rainfall-runoff modeling; a
334 stats component for statistical characterization of datasets; and a neural network component, which
335 demonstrated the implementation of feed-forward neural networks using the TensorFlow.js library for in-
336 browser machine learning applications.

337 Visualization and maps modules facilitated the interpretation and dissemination of results, with subsequent
338 training focused on data presentation. The visualization module provided instructions on rendering dynamic
339 charts and tables from analytical outputs using the Google Charts library. The maps module equipped
340 participants with the skills to generate interactive geospatial visualizations, including the rendering of
341 GeoJSON and Keyhole Markup Language (KML) data layers, utilizing either the Leaflet or Google Maps
342 platforms.

343 The second section of the training module addressed the challenge of computationally intensive tasks
344 through the HydroCompute library (Ramirez et al., 2024a). This library is engineered to execute high-
345 performance simulations on the client side, thereby leveraging the computational resources of the end-user
346 device. The instructions detailed HydroCompute's multi-engine architecture, which supports computations
347 via native JavaScript, WebAssembly (WASM), and WebGPU. The inclusion of WASM is particularly
348 significant, as it allows the execution of code compiled from high-performance languages like C and C++
349 at near-native speeds. The WebGPU engine allows for the utilization of a device's Graphics Processing Unit
350 (GPU) for parallelizable scientific computations. The training workflow involved data loading, defining
351 computational sequence with specified functions and dependencies, and the execution of the simulation,
352 demonstrating the feasibility of conducting sophisticated modeling within a standard browser environment.

353 The final component of the training introduced HydroRTC, a library designed to facilitate data sharing and
354 collaborative analysis. Built upon WebRTC and Socket.IO, HydroRTC provides a framework for both
355 server-to-peer (S2P) and peer-to-peer (P2P) communication. The curriculum explored several advanced use
356 cases. Participants learned to implement S2P data streaming for the efficient distribution of large datasets,
357 such as remote sensing imagery. A key concept presented was "smart data transmission," an approach where
358 an initial analysis on a data subset informs the prioritization of subsequent data delivery.

359 Furthermore, the training provided in-depth instruction on P2P architecture. This paradigm enables the
360 direct exchange of data and model results between clients, which can significantly reduce server load and
361 create decentralized, collaborative environments. The practical applications demonstrated included
362 volunteer computing for distributed modeling tasks and real-time data sharing among research collaborators
363 which highlight a shift towards more interactive web-native scientific workflows.

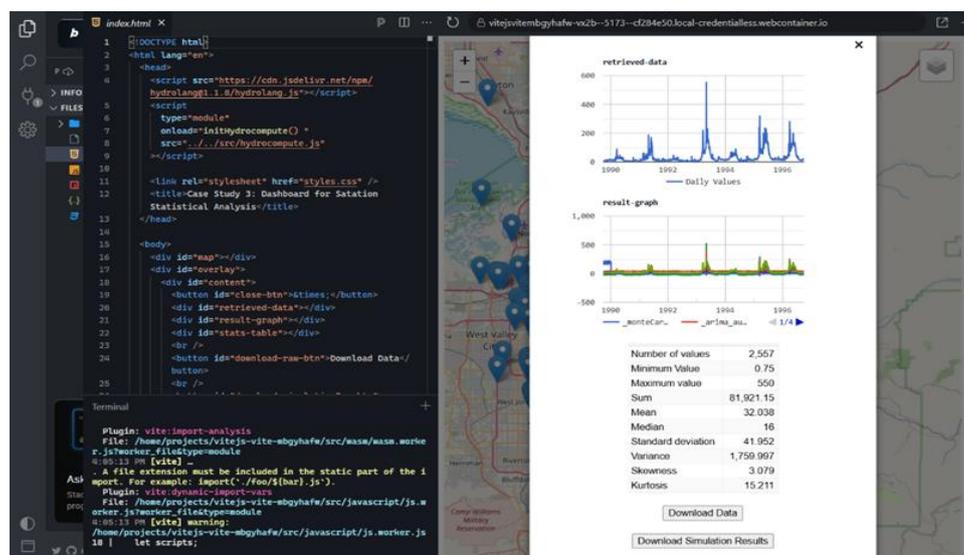
364 4.3.1.1. Case Study

365 During the training session, several case studies were presented as part of the accompanying session
366 materials, with the primary focus centered on demonstrating the interoperability and connectivity within
367 the HydroSuite collection of libraries spanning data acquisition, wrangling, and sorting through advanced
368 visualization for extracting insights from hydrological datasets. A representative case study examined the
369 Salt Lake City, Utah, USA using U.S. Geological Survey (USGS) gaging data, within which ARIMA
370 (Autoregressive Integrated Moving Average), Seasonal ARIMA (SARIMA), and Markov Chain Monte
371 Carlo modeling processes were implemented to quantify predictive trends derived from raw or minimally
372 processed data. The computations were performed dynamically on the user's local machine through a web-



373 based containerized environment, enabling participants to explore the underlying mechanisms inherent in
374 time-series hydrological data and to assess model responsiveness to real-world variability. By employing
375 web technologies and client-side hardware, the session further exposed participants to contemporary,
376 browser-based computational approaches that facilitate rapid prototyping and model deployment beyond
377 conventional development environments. These technologies supported the creation of interactive
378 dashboards and analytical utilities, allowing end users to conduct independent analyses across selected
379 geographic regions. This integration is exemplified in Figure 6, which presents the predictive outputs and
380 spatially contextualized information overlaid on a geospatial map illustrating the associated datasets.

381



382

383 Figure 6. Containerized app code shared with the participants in the 2024 WaterSoftHack cybertraining
384 for analyzing time series datasets obtained from multiple USGS stations.

385

386 4.3.2. Machine Learning Training Modules

387 Year 2 of WaterSoftHack cybertraining was dedicated to machine learning applications which provided
388 participants with a structured curriculum and hand-on training that progressed from foundational data
389 handling to advanced deep learning applications. The modules were designed to equip attendees with the
390 theoretical knowledge and practical skills necessary to implement machine learning solutions for complex
391 hydrological challenges. The initial phase of the training focused on the complete workflow development
392 of time series analysis, beginning with data acquisition, preparation, and normalization. Participants were
393 instructed on the programmatic retrieval of hydrological and meteorological datasets using APIs.
394 Following acquisition, the curriculum addressed critical data preprocessing techniques, including
395 methodologies for the imputation of missing values and the identification and treatment of outliers that are
396 frequently encountered in water science data. A significant component was dedicated to feature engineering,
397 wherein participants learned how to derive and select influential predictor variables from raw time series
398 data to enhance model performance. Upon establishing a foundation in data preparation, the module
399 introduced classical statistical methods for regression prediction. Instruction covered the theoretical



400 underpinnings and practical application of ARIMA and SARIMA models. These sessions provided
401 participants with a robust baseline for understanding and modeling temporal dependencies and seasonality
402 in hydrological data.

403 Building upon the statistical methods, the training program progressed to the application of neural networks
404 for water science time series prediction. This section began with an introduction to the Multi-Layer
405 Perceptron (MLP), a foundational class of feedforward artificial neural networks. The curriculum then
406 advanced to more complex recurrent architectures, with a specific focus on Long Short-Term Memory
407 (LSTM) networks. The training emphasized the advantages of LSTM in capturing long-term temporal
408 dependencies, a critical characteristic for many hydrological forecasting tasks.

409 The final phase of the WaterSoftHack week 1 training exposed participants to state-of-the-art machine
410 learning techniques and their real-world applications. Instructions were provided on the use of Transformer
411 models for time series forecasting which showed the architecture's powerful attention mechanism for
412 modeling complex sequences in time series data. To ensure the development of high-performing and
413 generalizable models, the curriculum included a dedicated segment on systematic hyperparameter
414 optimization methods as well.

415

416 **4.3.2.1. Case Study**

417 To consolidate the theoretical instruction, a parallel hands-on computation case study was demonstrated to
418 the participants. This case study was designed to provide practical experience in applying statistical and
419 machine learning models to real-world hydrological prediction tasks. The objective was to predict river
420 gauge height in a USGS gauging station one hour in advance. The models utilized a multivariate time series
421 dataset from the USGS gauging station 02336490. Hourly data from 2008 to 2024 was used for training,
422 testing and validation. This station is located at the outlet of the Proctor Creek-Chattahoochee River
423 watershed (HUC 031300020101) in Fulton County, Georgia (33°49'02.7" N, 84°28'49.2" W, NAD83).
424 Complementary meteorological data which included precipitation, air temperature, wind speed, relative
425 humidity and air pressure were obtained from the proximate Fulton County Airport monitoring station
426 (station id 03888). All data, including gauge height and meteorological variables, were aggregated to an
427 hourly temporal resolution.

428 Participants were guided through a complete data science workflow using computational notebooks
429 provided via Google Colaboratory. The standardized workflow included scripts for data acquisition, quality
430 control, imputation of missing values, and feature engineering. The modeling task was formulated as
431 predicting the gauge height at time $t + 1$ using a lookback window of the preceding 24 hours of
432 observations (both gauge height and meteorological covariates). The complete dataset was chronologically
433 partitioned into training (70%), validation (10%), and testing (20%) subsets to prevent temporal data
434 leakage. Five distinct time series forecasting models were developed and compared: ARIMA, SARIMA,
435 MLP, LSTM, a Transformer architecture. For the three machine learning models (MLP, LSTM, and
436 Transformer), hyperparameter optimization was demonstrated using the Optuna optimization framework
437 (Akiba et al., 2019). For the ARIMA model, the autoregressive order (p), differencing order (d), and
438 moving-average order (q) were set to 5, 1, and 5, respectively. These values were selected empirically
439 through trial-and-error, considering both predictive accuracy and computational cost. For the SARIMA
440 model, the same non-seasonal orders $(p, d, q) = (5, 1, 5)$ were used, together with seasonal orders



441 $(P, D, Q) = (2, 0, 2)$ and a seasonal period of 7. Relatively small seasonal orders were chosen to keep model
 442 training computationally feasible within the available time budget. The optimization process was
 443 configured to execute 100 trials, with the objective of minimizing the validation loss. Parameters were
 444 sampled from the specified intervals using tree-structured Parzen Estimator (TPE) (Bergstra et al.,
 445 2011). The final optimized hyperparameters obtained for each architecture are reported in Table 1.

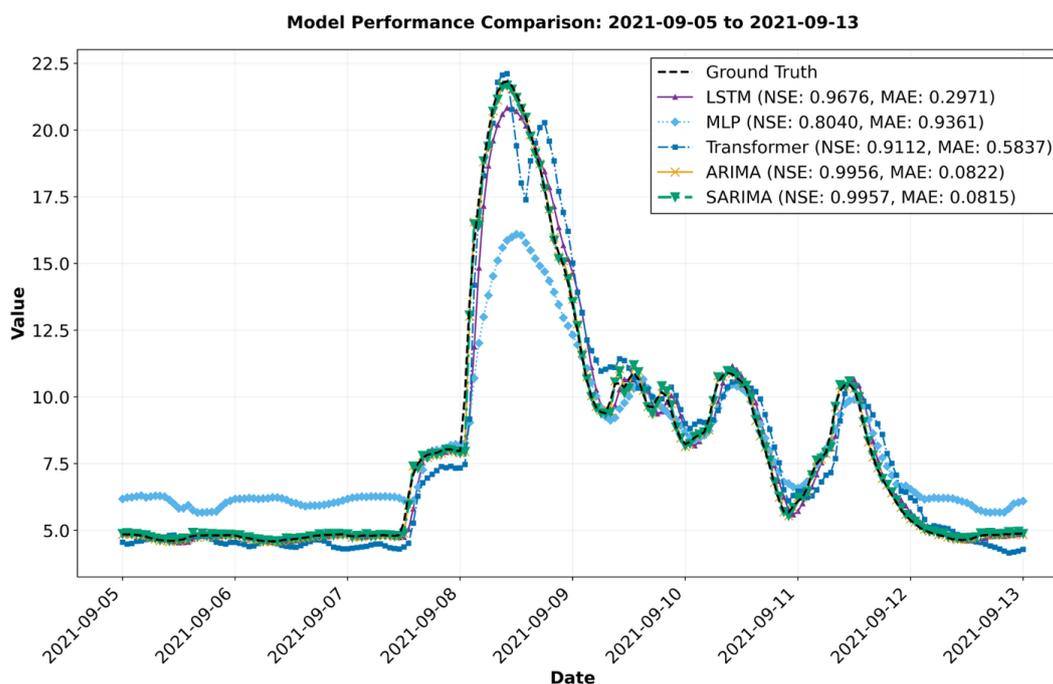
446 As shown in Table 1, the learning rate, which controls the step size during training, is highest for the LSTM
 447 (0.00103), lowest for the MLP (0.00022), and intermediate for the Transformer (0.00059). This reflects
 448 differences in how quickly each model adjusts its weight. Weight decay, which is a regularization parameter
 449 that penalizes large weights to reduce overfitting, is strongest for the LSTM (0.00025) while it is weakest
 450 for the MLP (3.22×10^{-5}) and moderate for the Transformer (0.00016). Dropout, another form of
 451 regularization where neurons are randomly deactivated during training, is not used in the LSTM, but it was
 452 set at 0.19 for the MLP, and slightly higher at 0.195 for the Transformer. This reflects the fact that the
 453 Transformer applies a smaller dropout of 0.110 to positional encoding to regularize the attention mechanism.
 454 Overall, these hyperparameter choices show that the LSTM prioritizes faster learning with strong weight
 455 decay while the MLP relies on careful weight updates with dropout for regularization, and the Transformer
 456 balances learning rate and weight decay while incorporating multiple layers of regularization to handle
 457 complex dependencies.

458 Table 1. Optimized hyperparameters for three different machine learning models demonstrated in the case
 459 study (Only a subset of hyperparameters were selected for demonstration).

Hyperparameters	LSTM	MLP	Transformer
Learning Rate	0.00103	0.00022	0.00059
Weight Decay	0.00025	3.22×10^{-5}	0.00016
Dropout	-	0.19	0.195
Dropout (Positional encoder)	-	-	0.110

460

461 Model performance was quantitatively evaluated on the unseen test dataset using two standard hydrological
 462 metrics: the Nash-Sutcliffe Efficiency (NSE) and the Mean Absolute Error (MAE). This practical
 463 application enabled a direct comparison of the models' respective abilities to capture complex temporal
 464 dynamics. Figure 7 provides a visual comparison of the predicted versus observed gauge heights for all five
 465 models over a representative one-week period for USGS02336490 gauging station. The USGS gauging
 466 station 02336490 is situated on the main stem of the Chattahoochee River near Atlanta, Georgia,
 467 downstream of several major reservoirs and receiving inflows from urbanized tributaries such as Proctor
 468 Creek. As illustrated, the ARIMA and SARIMA models achieved strong predictive performance which
 469 reflect the relatively stable and autocorrelated flow regime of this reservoir-regulated river segment. The
 470 LSTM model also demonstrated competitive skill. In contrast, the Transformer model exhibited lower
 471 predictive accuracy and tended to underestimate peak flows. This suggests that attention-based architecture
 472 may require larger datasets, inclusion of additional hydrological and meteorological covariates, or tailored
 473 training strategies to effectively capture the complex dynamics of gauge height at this urban-influenced and
 474 regulated site.



475

476 Figure 7. Comparison of predicted and observed gauge height data for September 5-13, 2021 in the USGS
477 gauging station 02336490. We used five models including ARIMA, SARIMA, LSTM, MLP and
478 Transformer to calibrate streamflow. The plot also displays the NSE and MAE values calculated for each
479 model over the test dataset.

480

481 4.3.3. Cloud and Edge Computing Module

482 The planned Year 3 module will introduce participants to cloud and edge computing frameworks and their
483 applications in water science and engineering. Building on knowledge acquired during the first two years,
484 the module will aim to familiarize participants with various cloud and edge computing paradigms,
485 commercial platforms, and water science-specific tasks that can be addressed using these technologies.

486 The training will begin with the fundamentals of cloud and edge computing. It will introduce the
487 participants to the concepts, architecture, and technologies that enable modern distributed systems. The
488 training will cover core cloud computing models such as Infrastructure as a Service (IaaS), Platform as a
489 Service, and Software as a Service, along with key characteristics like scalability, virtualization, and on-
490 demand resource provisioning. It will also cover edge computing fundamentals, focusing on processing
491 data closer to the source to reduce latency, to improve real-time performance, and to optimize bandwidth
492 usage. Finally, participants will see how to leverage cloud and edge computing for water science
493 applications.

494 Participants will learn to access and navigate the cloud and edge platforms via Cloudbank, with IBM Cloud
495 and its Environmental Intelligence Suite serving as the primary training platform. This suite provides ready-
496 to-use APIs and software development kits for a wide range of applications, including geospatial queries,



497 weather analytics, greenhouse gas emissions calculations, climate risk assessment, high-resolution imagery,
 498 and large curated datasets—making it particularly well-suited for water science projects. The curriculum
 499 will also cover optimal storage strategies for large hydrological datasets. The module will conclude with a
 500 comprehensive case study, where participants will apply these tools to a real-world project.

501 **4.4. WaterSoftHack Training Evaluation**

502 Pre- and post-surveys were conducted with WaterSoftHack Fellows to assess the effectiveness of the
 503 training. The evaluation followed a design-based research methodology, emphasizing iterative refinement
 504 of WaterSoftHack materials based on participant feedback. This approach enabled continuous
 505 improvements to both instructional content and delivery methods, informed by data collected during the
 506 two years of implementation. Data were collected by the educational research team through pre- and post-
 507 cybertraining surveys as well as follow-up online interviews with Fellows. The surveys included Likert-
 508 scale (Likert, 1932) and open-ended questions designed to evaluate satisfaction, skill development, and
 509 confidence in applying the training to participants’ research projects. Follow-up interviews were semi-
 510 structured, focusing on participants’ overall experiences, perceived challenges, and recommendations for
 511 enhancing the program.

512 The pre- and post-WaterSoftHack survey responses varied between years. In 2024, 13 of 15 Fellows
 513 completed the pre-cybertraining survey, and 12 responded to the post-cybertraining survey (Figure 8). In
 514 2025, 11 of 18 Fellows completed the pre-survey, while only 8 responded to the post-survey (Figure 9).
 515 Analysis of the raw survey data indicates a reduction in the proportion of participants who self-identified
 516 as “novices” across nearly all skill categories (see Tables 2 and 3), apart from software development in
 517 2025, which is planned to be explicitly addressed in the 2026 training. Figures 8 and 9 illustrate a clear
 518 decline in novice-level self-assessments for both years. For instance, in 2025, only 45.5% of participants
 519 rated themselves as competent or proficient in machine learning before the training, whereas this proportion
 520 increased to 87.5% afterward. Similarly, in 2024, only 23.1% of participants considered themselves above
 521 novice level in data analytics and prototyping before the WaterSoftHack cybertraining, but it increased to
 522 >50% after the completion. These trends suggest notable self-reported improvements across multiple skill
 523 domains. Formal statistical analysis of these changes was not possible because the post-training survey data
 524 were anonymous, unmatched, and limited in number. As a result, the dataset does not allow for a
 525 quantitative estimate of the program’s impact on skill development. The evaluation of WaterSoftHack’s
 526 effectiveness therefore draws primarily on qualitative evidence.

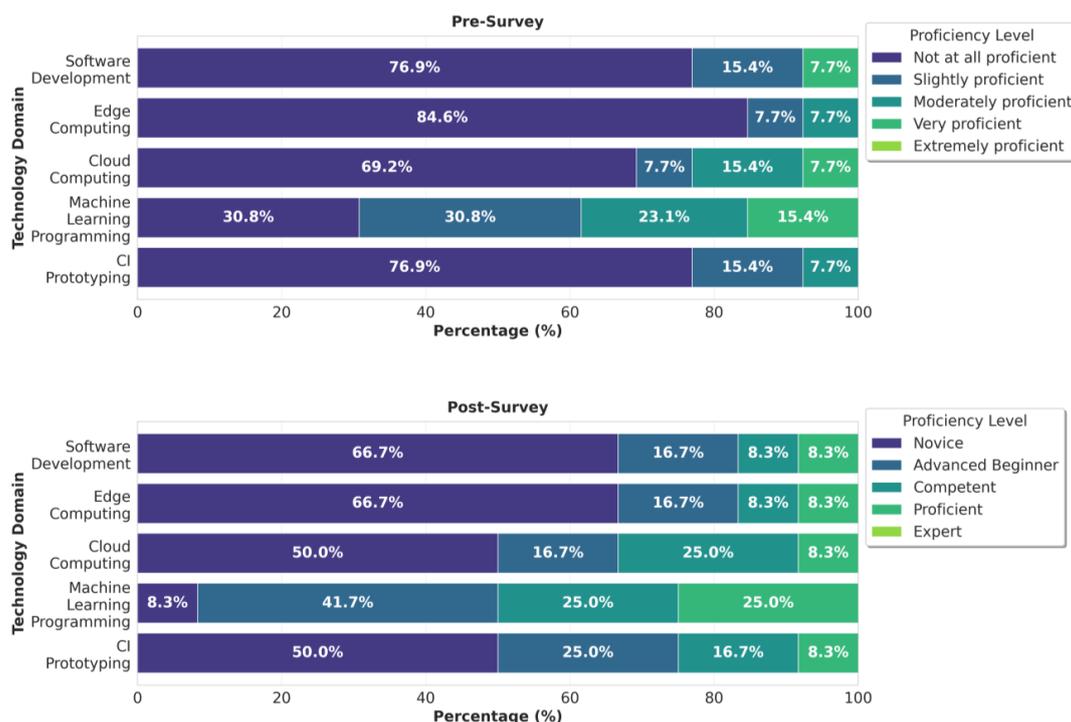
527 Table 2. Comparison of self-reported novice counts from the 2024 pre- and post WaterSoftHack-
 528 cybertraining surveys. "Expected Novices" is the number that would be expected in the post-survey if the
 529 proportion of novices had remained unchanged from the pre-survey.

Skill Category	Observed Novices (after cybertraining)	Expected Novices considering no improvement (%)
Software Development	6	9.23
Edge Computing	1	3.69
Cloud Computing	6	8.31
Machine Learning Programming	8	10.15
CI Prototyping	8	9.23

530



531



532

533

534

535

Figure 8. The results of pre- and post-cybertraining survey results conducted among the Fellows for WaterSoftHack 2024. 13 out of 15 Fellows responded to the pre-survey and 12 out of 15 Fellows responded to the post-survey.

536

537

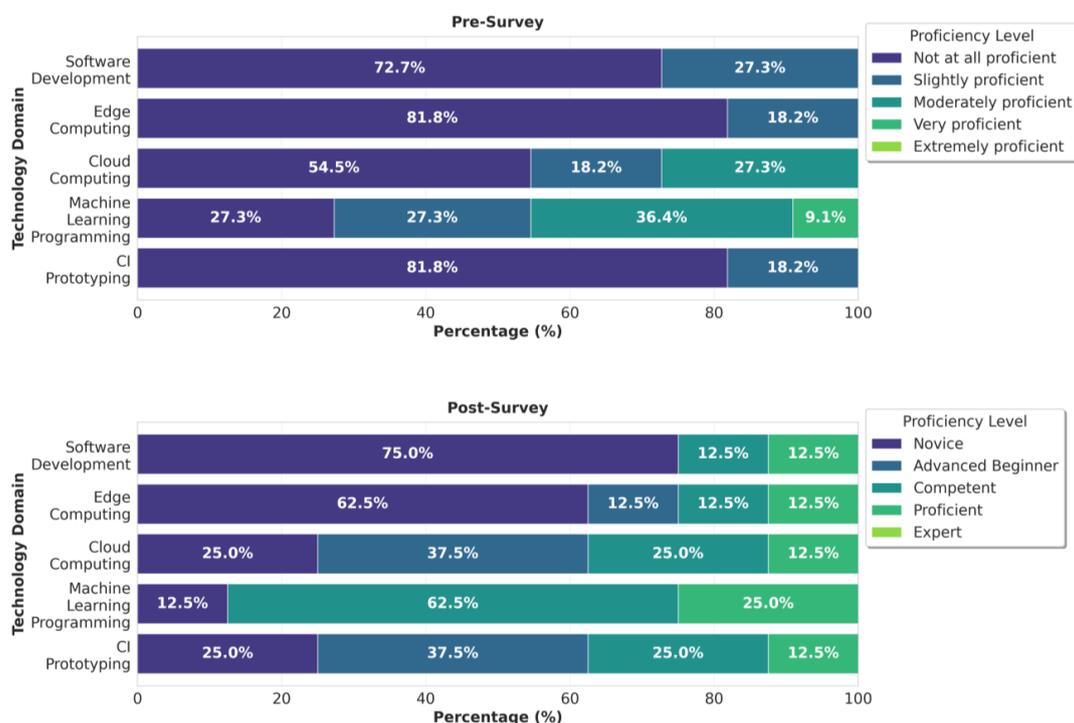
538

Table 3. Comparison of self-reported novice counts from the 2025 pre- and post- cybertraining surveys. "Expected Novices" is the number that would be expected in the post-survey if the proportion of novices had remained unchanged from the pre-survey.

Skill Category	Observed Novices (after cybertraining)	Expected Novices considering no improvement (%)
Software Development	6	5.81
Edge Computing	5	6.54
Cloud Computing	2	4.36
Machine Learning Programming	1	2.18
CI Prototyping	2	6.54

539

540



541
 542 Figure 9. The results of pre- and post-WaterSoftHack cybertraining survey results conducted among the
 543 Fellows for WaterSoftHack 2025. 11 out of 18 Fellows responded to the pre-survey and 8 out of 18 Fellows
 544 responded to the post-survey.

546 4.5. Participant Experience

547 Post-event surveys from the 2024 and 2025 editions of WaterSoftHack cybertraining consistently revealed
 548 strong positive feedback regarding networking and collaboration among participants. Throughout both
 549 years, attendees praised the event for providing valuable opportunities to connect with peers in related fields,
 550 which in turn fostered new research partnerships spanning various institutions and disciplines. Participants
 551 highlighted that these connections often extended beyond the event itself, with many engaging in ongoing
 552 communication or team projects started during the program.

553 A particularly notable strength of WaterSoftHack, as identified in participant feedback, was the practical
 554 focus of its hands-on components. The majority of respondents reported that these sessions helped them
 555 apply and learn computational techniques, such as machine learning and advanced cyberinfrastructure
 556 workflows, which they planned to integrate into future research and teaching. These practical sessions were
 557 frequently mentioned as the most useful and engaging, giving participants the confidence and skills needed
 558 to apply new methods in their work.

559 Survey responses further indicated that the program’s collaborative environment contributed to tangible
 560 outcomes. Several participants from both years stated specific intentions to develop their hackathon or
 561 capstone projects into publishable studies. Some were already collaborating on joint publications or follow-
 562 on research initiatives spurred by connections made during the event. The structure of the program,



563 emphasizing interdisciplinary teamwork and accessible mentorship, was credited with expanding
564 participants' research capacity and confidence in using new computational tools. Statistical outcomes
565 support these themes in 2024 and 2025 post-survey results.

566 WaterSoftHack post-event survey indicated that the majority of respondents recognized the program as
567 highly effective in connecting traditional water science with modern computational methodologies, with
568 60% rating it as "Very Effective" or "Extremely Effective". Participants consistently reported expanded
569 knowledge in applying machine learning to water science, including gaining foundational theoretical
570 understanding, enhancing practical coding abilities, and successfully integrating new methods into
571 authentic research problems such as streamflow prediction and data-driven hydrologic modeling.

572 Survey results revealed that open-source machine learning models and collaborative tools led 59% of
573 participants to significantly change their approach to computing in water science, notably improving their
574 predictive capacities and supporting real-time, cloud-enabled collaborative modeling efforts. Over half of
575 the respondents reported increased competence in building and using advanced cyberinfrastructure
576 workflows utilizing platforms like Google Colab, NAIRR, and GitHub, which were praised for their
577 accessibility and direct relevance to research projects.

578 Participants highlighted several key strengths of WaterSoftHack cybertraining program. Major strength
579 among them was the hands-on, applied learning environment, which combined with expert mentorship,
580 created a powerful setting for skill development and confidence-building. Most respondents indicated that
581 the sessions and tools learned were immediately transferable to their work, either supporting ongoing
582 research or facilitating future teaching and outreach. Collaboration emerged as a major benefit, with many
583 initiating new partnerships, sharing resources, and launching joint projects during and after the program.

584 While satisfaction and session ratings were overwhelmingly positive, especially for interactive, project-
585 based learning that reinforced concepts and skills, participants also identified opportunities for future
586 improvement. These included extending the program duration beyond two weeks, providing more
587 discipline-specific case studies, improving team matching, and clarifying educational materials to increase
588 inclusiveness for diverse interdisciplinary backgrounds. Furthermore, several participants commented on
589 scheduling, noting that the 8:00 a.m. Pacific Time start was challenging for those on the US West Coast,
590 and recommended more flexible timing to accommodate a wider geographic audience.

591 Feedback also reflected a strong appreciation for the program's collaborative and supportive culture,
592 impactful technical content, and opportunities for networking with peers and mentors. Overall, participants
593 described substantial professional growth and valued the chance to develop new skills, apply them in
594 research, and form lasting collaborations. The program's applied and collaborative approach was seen as a
595 model for advanced training at the intersection of water science and computational methods. Comparing
596 results from both 2024 and 2025, the WaterSoftHack has consistently satisfied key needs of the water
597 science community. While some 2024 participants noted that late delivery of training materials limited their
598 preparation, this issue was addressed in 2025 by providing resources six months in advance which markedly
599 improved pre-training readiness. Many also recommended earlier team formation and project initiation
600 during week one to boost collaboration efficiency, a suggestion that is being considered for implementation
601 in the 2026 WaterSoftHack iteration.



602 Findings from both years' surveys consistently underscore WaterSoftHack's significant impact on
603 professional development, collaborative research, and the application of computational methods in water
604 science. This shows the program's ongoing value as a model for interdisciplinary, cybertraining-driven
605 scientific capacity building for water science and engineering.

606

607 **4.6. Qualitative Insights from WaterSoftHack Fellows**

608 CUAHSI conducted post-event reflections with WaterSoftHack Fellows to illustrate the program's
609 influence on research practices and professional development (CUAHSI, 2025; 2026). One doctoral
610 participant reported that training in HydroLang significantly enhanced their ability to process and visualize
611 large remote-sensing datasets. These newly acquired skills are now being applied in their dissertation
612 research on water quality monitoring. The participant described plans to develop a web-based API for
613 automated monitoring using techniques learned during the workshop. Another fellow emphasized the
614 networking opportunities facilitated by the program, noting that connections formed during WaterSoftHack
615 evolved into ongoing professional relationships reinforced at subsequent scientific meetings.

616

617 A postdoctoral scholar highlighted the value of centralized and curated training resources that reduced the
618 fragmentation commonly encountered in self-directed learning. Collaborative work during the hackathon
619 resulted in the creation of shared resources on HydroShare and a continental-scale assessment of stream
620 temperature variability, which was later presented at a major scientific conference and is currently being
621 prepared for journal submission (CUAHSI, 2025). Reflections from the 2025 fellows further reinforced
622 these themes. Participants reported that hands-on machine learning modules covering model development,
623 hyperparameter tuning, and evaluation directly improved their ability to build predictive models for
624 environmental applications such as sediment forecasting and water quality monitoring (CUAHSI, 2026).
625 Fellows also highlighted the importance of interdisciplinary teamwork during capstone projects, which
626 enabled participants from different institutions and disciplinary backgrounds to collaborate on real-world
627 water science challenges and expand their professional networks.

628

629 Several fellows indicated that the skills gained during WaterSoftHack were immediately transferable to
630 their graduate research activities. Examples included developing satellite-based models for suspended
631 sediment estimation using machine learning, applying LSTM models to downscale satellite hydrologic
632 datasets, and experimenting with ensemble machine learning approaches for environmental forecasting
633 (CUAHSI 2025, 2026). These reflections collectively illustrate how the program not only builds technical
634 competencies but also fosters collaborative research environments that improve the adoption of data-driven
635 methods in water science.

636

637 **4.7. Candidate Follow-up and Career Outcomes**

638 Across the two candidate cohorts (2024–2025), a total of 27 WaterSoftHack Fellows were recruited into
639 the program, with each cohort comprising 15–18 candidates. In addition to the Fellows, the program
640 provided training to over 350 participants, of whom approximately 40% identified as female and 60% as
641 male. Due to the constraints of NSF funding, participation in the WaterSoftHack Fellowships was limited
642 to U.S.-based candidates (100%). However, the general training (week 1) attracted a diverse international
643 audience, including participants from Europe, Asia, the Middle East, and Africa. Most candidates entered
644 the program after completing a postgraduate master's degree, while a smaller subset had completed a 4-
645 year undergraduate degree or a bachelor's degree. Consistent with the program's interdisciplinary mission,



673 teamwork, and empowering researchers to apply these approaches in their ongoing and future research.
674 Many participants credited the training with advancing their technical skills, research productivity, and
675 openness to innovation.

676 At the same time, the process of developing and implementing WaterSoftHack revealed important
677 challenges such as the need for longer-duration training, improved preparation and team formation, and
678 more adaptive materials for a diverse audience spanning earth and water sciences. One of the most
679 significant challenges was fostering deep collaboration in a short-term virtual setting. Participants
680 appreciated the program's networking opportunities, collaborative spirit, and supportive environment,
681 which enabled valuable peer connections across institutions and time zones. However, many noted that
682 earlier and more structured team formation could further enhance these connections. The fast pace and
683 intensity of the one-week hackathon offered limited time for teams to develop cohesion, balance different
684 technical backgrounds, and converge on shared research objectives. These observations suggest that
685 allocating additional time for team building and project ideation potentially extending collaborative
686 activities over longer periods, such as an academic semester can further enhance the effectiveness of hands-
687 on cybertraining.

688 Building on these lessons, future iterations of WaterSoftHack incorporate pre-cybertraining virtual mixers
689 and dedicate deliberate time in the first week for project ideation, team building, and familiarization with
690 collaborative tools. Structured activities to refine research goals and finalize teams before the official start
691 of the hackathon would give participants a critical head start, allowing them to enter the intensive phase
692 with a clear direction and stronger rapport.

693 Another overarching challenge was evaluating program impact in a meaningful way. As highlighted in the
694 results section, small cohort sizes and the need for anonymous, unmatched survey design limited the ability
695 to derive statistically valid conclusions about improvements in skill and knowledge. This limitation
696 underscores a larger issue facing intensive, short-term training initiatives. The experience emphasized two
697 key lessons: qualitative feedback that was gathered through interviews and open-ended survey responses
698 proves invaluable and should be given prominent consideration. In addition, robust evaluation required both
699 matched pre- and post-surveys as well as longer-term tracking. Future evaluations should therefore
700 incorporate follow-up surveys conducted six to twelve months after the program to better capture long-term
701 outcomes. This approach can reveal the broader impacts, such as successful application of learned skills,
702 the development of new open-source tools, and research publications emerging from program participation.

703 Integrating computational skills into water science marks more than just a technical enhancement. It signals
704 a transformative paradigm shift in how scientific knowledge within the field of hydrology is generated,
705 interpreted, and translated into decision. By weaving computational competencies into the fabric of water
706 science education, the community stands poised to fully harness today's data-rich environment to unlock
707 the capacity for more sophisticated, predictive insights into critical issues in hydrology.

708 The WaterSoftHack program demonstrates that intensive, hands-on cybertraining can meaningfully
709 improve the transition toward data-driven hydrologic research. Participants not only gained foundational
710 and advanced cyberinfrastructure expertise but also demonstrated how such skills can immediately impact
711 real-world research, team projects, and lasting collaborations across institutional boundaries. Through
712 mentorship, applied learning, and continual refinement, WaterSoftHack reinforces the importance of
713 moving beyond traditional hydrological training programs to more modular and collaborative learning



714 frameworks as the field evolves. With this vision, WaterSoftHack provides a scalable model for preparing
715 the next generation of water scientists to effectively operate in an increasingly digital and data-rich research
716 environment.

717
718 **6. Acknowledgements:** This work is supported by the U.S. National Science Foundation (NSF) Office of
719 Advanced Cyberinfrastructure (AOR2320979). All opinions, findings, and conclusions or
720 recommendations expressed in this material are those of the authors and do not necessarily reflect the views
721 of the NSF. NSF NAIRR is acknowledged for its generous allotment of computing time for WaterSoftHack
722 cybertraining. Alan Alda Center at Stony Brook University (USA) is acknowledged for hosting the
723 scientific communication session. We acknowledge the participation of students, researchers and early
724 career faculty members across CUAHSI's >130 university partner networks.

725 **7. Software Availability**

726 **Name of package:** WaterSoft Python Package

727 **Year first available:** 2024

728 **Developers:** WaterSoftHack Team, Clemson University, Tulane University, University of Iowa, &
729 CUAHSI

730 **Package Availability:** <https://github.com/watersofthack/WaterSoft>

731 **License:** MIT License

732 **Software requirements:** Python ≥ 3.10 , numpy, pandas, scipy, matplotlib, seaborn, plotly, tensorflow,
733 torch, scikit-learn, statsmodels, geopandas, folium, requests, optuna, jupyter, notebook, dash, websockets,
734 python-socketio
735

736

737 **8. Data availability**

738 The hydrological data used in this study are publicly available from the U.S. Geological Survey National
739 Water Information System (NWIS) at <https://waterdata.usgs.gov/nwis/rt>. All human-subject-related data
740 referenced in this paper are fully aggregated and are provided within the tables of the manuscript.

741

742 **9. Author contributions**

743 KP and VS prepared and edited the draft manuscript. All authors contributed to the final version of the
744 manuscript.

745

746 **10. Ethics Statement**

747 The studies involving human participants were reviewed and approved by Clemson University Institutional
748 Review Board. The participants provided their written informed consent to participate in this study.

749

750 **11. Competing interests**

751 The authors declare that they have no conflict of interest.

752



753 **12. References**

- 754 Abbott, M. B., Bathurst, J. C., Cunge, J. A., O’Connell, P. E., & Rasmussen, J. (1986). An introduction to
755 the European Hydrological System—Systeme Hydrologique Europeen, “SHE”, 1: History and philosophy
756 of a physically-based, distributed modelling system. *Journal of Hydrology*, 87(1–2), 45–59.
- 757 Abbott, M. B., Bathurst, J., Cunge, J., O’connell, P., & Rasmussen, J. (1986). An introduction to the
758 European Hydrological System—Systeme Hydrologique Europeen, “SHE”, 2: Structure of a physically-
759 based, distributed modelling system. *Journal of Hydrology*, 87(1–2), 61–77.
- 760 Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A Next-generation
761 Hyperparameter Optimization Framework. *Proceedings of the 25th ACM SIGKDD International
762 Conference on Knowledge Discovery & Data Mining*, 2623–2631.
763 <https://doi.org/10.1145/3292500.3330701>
- 764 Anslow, C., Brosz, J., Maurer, F., & Boyes, M. (2016). Datathons: An experience report of data
765 hackathons for data science education. *Proceedings of the 47th ACM Technical Symposium on
766 Computing Science Education*, 615–620.
- 767 Bergstra, J., Bardenet, R., Bengio, Y. and Kégl, B., 2011. Algorithms for hyper-parameter
768 optimization. *Advances in neural information processing systems*, 24.
- 769 Beven, K. J., & Kirkby, M. J. (1979). A physically based, variable contributing area model of basin
770 hydrology/Un modèle à base physique de zone d’appel variable de l’hydrologie du bassin versant.
771 *Hydrological Sciences Journal*, 24(1), 43–69.
- 772 Bjerknes, V. (1904). Das Problem der Wettervorhersage, betrachtet vom Standpunkte der Mechanik und
773 der Physik. *Meteor. Z.*, 21, 1–7.
- 774 Blumenfeld, P. C., Soloway, E., Marx, R. W., Krajcik, J. S., Guzdial, M., & Palincsar, A. (1991).
775 Motivating project-based learning: Sustaining the doing, supporting the learning. *Educational
776 Psychologist*, 26(3–4), 369–398.
- 777 Briscoe, G. (2014). *Digital innovation: The hackathon phenomenon*.
- 778 Brooks-Harris, J. E., & Stock-Ward, S. R. (1999). *Workshops: Designing and facilitating experiential
779 learning*. Sage Publications.



- 780 Brown, J. D., Wu, L., He, M., Regonda, S., Lee, H., & Seo, D.-J. (2014). Verification of temperature,
781 precipitation, and streamflow forecasts from the NOAA/NWS Hydrologic Ensemble Forecast
782 Service (HEFS): 1. Experimental design and forcing verification. *Journal of Hydrology*, 519,
783 2869–2889.
- 784 Charney, J. G., Fjørtoft, R., & Neumann, J. von. (1950). Numerical integration of the barotropic vorticity
785 equation. *Tellus*, 2(4), 237–254.
- 786 Cikmaz, B. A., Yildirim, E., & Demir, I. (2025). Flood susceptibility mapping using fuzzy analytical
787 hierarchy process for Cedar Rapids, Iowa. *International Journal of River Basin Management*,
788 23(1), 1-13.
- 789 Creswell, J. W., & Clark, V. L. P. (2017). *Designing and conducting mixed methods research*. Sage
790 publications.
- 791 CUAHSI. 2025. Community guest spotlight with 2024 WaterSoftHack Fellows. Retrieved January 26,
792 2026 from [https://www.cuahsi.org/community/news/april-e-newsletter-community-guest-](https://www.cuahsi.org/community/news/april-e-newsletter-community-guest-spotlight-with-2024-watersofthack-fellows)
793 [spotlight-with-2024-watersofthack-fellows](https://www.cuahsi.org/community/news/april-e-newsletter-community-guest-spotlight-with-2024-watersofthack-fellows)
- 794 CUAHSI. 2026. Community guest spotlight with 2025 WaterSoftHack Fellows. Retrieved March 13,
795 2026, from [https://www.cuahsi.org/community/news/march-e-newsletter-community-guest-](https://www.cuahsi.org/community/news/march-e-newsletter-community-guest-spotlight-with-2025-watersofthack-fellows)
796 [spotlight-with-2025-watersofthack-fellows](https://www.cuahsi.org/community/news/march-e-newsletter-community-guest-spotlight-with-2025-watersofthack-fellows)
- 797 Demargne, J., Wu, L., Regonda, S. K., Brown, J. D., Lee, H., He, M., Seo, D.-J., Hartman, R., Herr, H.
798 D., & Fresch, M. (2014). The science of NOAA’s operational hydrologic ensemble forecast
799 service. *Bulletin of the American Meteorological Society*, 95(1), 79–98.
- 800 Demir, I., Xiang, Z., Demiray, B., & Sit, M. (2022). WaterBench-Iowa: a large-scale benchmark dataset
801 for data-driven streamflow forecasting, *Earth Syst. Sci. Data*, 14, 5605–5616.
- 802 Demir, I., Conover, H., Krajewski, W. F., Seo, B.-C., Goska, R., He, Y., McEniry, M. F., Graves, S. J., &
803 Petersen, W. (2015). Data-enabled field experiment planning, management, and research using
804 cyberinfrastructure. *Journal of Hydrometeorology*, 16(3), 1155–1170.



- 805 Design-Based Research Collective. (2003). Design-based research: An emerging paradigm for
806 educational inquiry. *Educational Researcher*, 32(1), 5–8.
- 807 Entekhabi, D., Njoku, E. G., O’neill, P. E., Kellogg, K. H., Crow, W. T., Edelstein, W. N., Entin, J. K.,
808 Goodman, S. D., Jackson, T. J., & Johnson, J. (2010). The soil moisture active passive (SMAP)
809 mission. *Proceedings of the IEEE*, 98(5), 704–716.
- 810 Gochis, D., Yu, W., & Yates, D. (2015). The WRF-Hydro model technical description and user’s guide,
811 version 3.0. NCAR Technical Document, 120.
- 812 Haw, W., & Crawford, A. (2025). Using Hackathons to Enhance University-Level Water Curriculum for
813 Students in Minority Communities. *Journal of Learning Development in Higher Education*.
- 814 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey,
815 C., Radu, R., & Schepers, D. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal*
816 *Meteorological Society*, 146(730), 1999–2049.
- 817 Hou, A. Y., Kakar, R. K., Neeck, S., Azarbarzin, A. A., Kummerow, C. D., Kojima, M., Oki, R.,
818 Nakamura, K., & Iguchi, T. (2014). The global precipitation measurement mission. *Bulletin of*
819 *the American Meteorological Society*, 95(5), 701–722.
- 820 Huppenkothen, D., Arendt, A., Hogg, D. W., Ram, K., VanderPlas, J., & Rokem, A. (2017). Hack Weeks
821 as a model for data science education and collaboration. *arXiv*. <https://arxiv.org/abs/1711.00028>
822
- 823 Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M., Saha, S., White,
824 G., & Woollen, J. (2018). The NCEP/NCAR 40-year reanalysis project. In *Renewable energy* (p.
825 Vol1_146-Vol1_194). Routledge.
- 826 Kirkpatrick, D., & Kirkpatrick, J. (2006). *Evaluating training programs: The four levels*. Berrett-Koehler
827 Publishers.
- 828 Krajewski, W. F., Ghimire, G. R., Demir, I., & Mantilla, R. (2021). Real-time streamflow forecasting: AI
829 vs. Hydrologic insights. *Journal of Hydrology X*, 13, 100110.



- 830 Kummerow, C., Barnes, W., Kozu, T., Shiue, J., & Simpson, J. (1998). The tropical rainfall measuring
831 mission (TRMM) sensor package. *Journal of Atmospheric and Oceanic Technology*, 15(3), 809–
832 817.
- 833 Lara, M., & Lockwood, K. (2016). Hackathons as community-based learning: A case study. *TechTrends*,
834 60(5), 486–495.
- 835 Liang, X., Lettenmaier, D. P., Wood, E. F., & Burges, S. J. (1994). A simple hydrologically based model
836 of land surface water and energy fluxes for general circulation models. *Journal of Geophysical*
837 *Research: Atmospheres*, 99(D7), 14415–14428.
- 838 Lighthill, M. J., & Whitham, G. B. (1955). On kinematic waves I. Flood movement in long rivers.
839 *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*,
840 229(1178), 281–316.
- 841 Likert, R., 1932. A technique for the measurement of attitudes. *Archives of psychology*.
842
- 843 Nash, J. (1957). The form of the instantaneous unit hydrograph. *Comptes Rendus et Rapports Assemblée*
844 *Generale de Toronto*, 3, 114–121.
- 845 Seibert, J. and Bergström, S., 2022. A retrospective on hydrological catchment modelling based on half a
846 century with the HBV model. *Hydrology and Earth System Sciences*, 26(5), pp.1371-1388.
847
- 848 Ramirez, C. E., Sermet, Y., & Demir, I. (2024). HydroCompute: An open-source web-based
849 computational library for hydrology and environmental sciences. *Environmental Modelling &*
850 *Software*, 175, 106005.
- 851 Ramirez, C. E., Sermet, Y., Molkenhain, F., & Demir, I. (2022). HydroLang: An open-source web-based
852 programming framework for hydrological sciences. *Environmental Modelling & Software*, 157,
853 105525.
- 854 Richardson, L. F. (1922). *Weather prediction by numerical process*. Franklin Classics.



- 855 Sherman, L. K. (1932). The relation of hydrographs of runoff to size and character of drainage-basins.
856 Eos, Transactions American Geophysical Union, 13(1), 332–339.
- 857 Seo, B. C., Keem, M., Hammond, R., Demir, I., & Krajewski, W. F. (2019). A pilot infrastructure for
858 searching rainfall metadata and generating rainfall product using the big data of NEXRAD.
859 Environmental modelling & software, 117, 69-75.
- 860 Tapley, B. D., Bettadpur, S., Ries, J. C., Thompson, P. F., & Watkins, M. M. (2004). GRACE
861 measurements of mass variability in the Earth system. Science, 305(5683), 503–505.
- 862 Wagener, T., Savic, D., Butler, D., Ahmadian, R., Arnot, T., Dawes, J., Djordjevic, S., Falconer, R.,
863 Farmani, R., Ford, D., Hofman, J., Kapelan, Z., Pan, S., and Woods, R.: Hydroinformatics education – the
864 Water Informatics in Science and Engineering (WISE) Centre for Doctoral Training, Hydrol. Earth Syst.
865 Sci., 25, 2721–2738, <https://doi.org/10.5194/hess-25-2721-2021>, 2021.
- 866