

Towards a Multi-Representational Approach to Prediction, Understanding, and Discovery in Hydrology

Luis De la Fuente¹, Hoshin Vijai Gupta¹, and Laura Elizabeth Condon¹

¹University of Arizona

November 26, 2022

Abstract

A key step in model development is selection of an appropriate representational system, including both the representation of what is observed (the data), and the formal mathematical structure used to construct the input-state-output mapping. These choices are critical, because they completely determine the questions we can ask, the nature of the analyses and inferences we can perform, and the answers that we can obtain. Accordingly, a representation that is suitable for one kind of investigation might be limited in its ability to support some other kind.

Arguably, how different representational approaches affect what we can learn from data is poorly understood. This paper explores three complementary representational strategies as vehicles for understanding how catchment-scale hydrological processes vary across hydro-geo-climatologically diverse Chile. Specifically, we test a lumped water-balance model (GR4J), a data-based dynamical systems model (LSTM), and a data-based regression-tree model (Random Forest). Insights were obtained regarding system memory encoded in data, spatial transferability by use of surrogate attributes, and informational deficiencies of the dataset that limit our ability to learn an adequate input-output relationship. As expected, each approach exhibits specific strengths, with LSTM providing the best characterization of dynamics, GR4J being the most robust under informationally deficient conditions, and RF being most supportive of interpretation.

Overall, the complementary nature of the three approaches suggests the value of adopting a multi-representational framework in order to more fully extract information from the data. Our results show that a multi-representational approach better supports the goals of prediction, understanding, and scientific discovery in Hydrology.

Towards a Multi-Representational Approach to Prediction, Understanding, and Discovery in Hydrology

Luis A. De la Fuente¹, Hoshin V. Gupta¹, and Laura E. Condon¹

¹Department of Hydrology and Atmospheric Sciences
The University of Arizona, Tucson, AZ 85721, USA

Corresponding authors:

Luis A. De la Fuente (ldelafue@email.arizona.edu), ORCID: 0000-0001-6979-0547

Hoshin V. Gupta (hoshin@email.arizona.edu), ORCID: 0000-0001-9855-2839

Laura E. Condon (lecondon@email.arizona.edu), ORCID: 0000-0003-3639-8076

KEY POINTS

- The representation underlying a model pre-determines what can be learned, which argues for a flexible approach to scientific investigation.
- By employing multiple representational approaches, we improve our chances of properly understanding the underlying Data Generation Process.
- Such an approach helped in understanding how to model precipitation-streamflow response across hydro-geo-climatologically diverse Chile.

KEYWORDS

Representation, Machine Learning, LSTM, Random Forest, GR4J, conceptual model, lumped water balance model, Understanding, Discovery, Hydrological Processes, Catchments, Hydro-geo-climatology

ABSTRACT

A key step in model development is selection of an appropriate representational system, including both the representation of what is observed (the data), and the formal mathematical structure used to construct the input-state-output mapping. These choices are critical, because they completely determine the questions we can ask, the nature of the analyses and inferences we can perform, and the answers that we can obtain. Accordingly, a representation that is suitable for one kind of investigation might be limited in its ability to support some other kind.

Arguably, how different representational approaches affect what we can learn from data is poorly understood. This paper explores three complementary representational strategies as vehicles for understanding how catchment-scale hydrological processes vary across hydro-geo-climatologically diverse Chile. Specifically, we test a lumped water-balance model (GR4J), a data-based dynamical systems model (LSTM), and a data-based regression-tree model (Random Forest). Insights were obtained regarding system memory encoded in data, spatial transferability by use of surrogate attributes, and informational deficiencies of the dataset that limit our ability to learn an adequate input-output relationship. As expected, each approach exhibits specific strengths, with LSTM providing the best characterization of dynamics, GR4J being the most robust under informationally deficient conditions, and RF being most supportive of interpretation.

Overall, the complementary nature of the three approaches suggests the value of adopting a multi-representational framework in order to more fully extract information from the data. Our results show that a multi-representational approach better supports the goals of prediction, understanding, and scientific discovery in Hydrology.

PLAIN LANGUAGE SUMMARY

The representations we use when analyzing data and modeling systems completely determine the *questions* we can ask, the nature of the *analyses and inferences* we can perform, and the *answers* that we can obtain. So, any given modeling approach may be highly suitable for learning certain things about a system but be completely unsuitable for learning other things. To explore how different representational approaches can affect what we can learn from data, we explore how three complementary modeling approaches (one lumped water balance and two machine-learning methods) can support an improved understanding of how catchment-scale hydrological processes vary across the diverse hydro-geo-climatology of Chile. Each approach was found to exhibit specific strengths, and interesting insights were obtained regarding system memory, attributes that correlate with transferability across different regions, and informational deficiencies of the available dataset. Overall, this study suggests the value of adopting a general multi-representational framework to better support prediction, understanding and scientific discovery in the Earth and Environmental Sciences.

1 1 INTRODUCTION

2 1.1 The Problem of Selecting an Appropriate Representational System

3 [1] When developing any dynamical systems model, be it conceptual or data-based, a key
4 step is the selection of an appropriate representational system. This step includes two aspects:
5 (1) the choice of system inputs (driving variables) and boundary conditions relevant to
6 predicting the dynamical evolution of the system states and outputs, and (2) the formal
7 mathematical/algorithmic structure used to construct the input-output (or input-state-output)
8 mappings that are hypothesized to characterize the system (*Gupta et al, 2012; Gharari et al,*
9 *2021*).

10 [2] In hydrology, as in other environmental disciplines, the selection of system inputs and
11 boundary conditions determines the nature and quality of the information that can be brought
12 to bear on the prediction problem – without adequate and informationally relevant data, the
13 task of predicting the system outputs is doomed from the outset. Having done so, the
14 mathematical/algorithmic representational system selected for constructing the input-output
15 mapping is critical, because it completely determines the *questions* we can ask, the nature of
16 the *analyses and inferences* we can perform, and the *answers* that we can obtain.

17 [3] For example, a dynamical process-resolved (often called physically-based) catchment-
18 scale hydrological model is typically constructed to answer questions such as “*what kind and*
19 *magnitude of streamflow response can we expect to see when a specific catchment system is*
20 *perturbed by a certain sequence of rainfall (and temperature) inputs*”. If a mass/energy-
21 conserving spatially-lumped bucket-type state-space representation is implemented, then one
22 may be able to obtain insights into aggregate catchment-scale soil moisture storage variations
23 (and their vertical distribution in the soil zone), whereas a spatially-distributed finite-
24 element/difference representation may be used to infer the dynamic evolution of soil-
25 moisture (and other state-variables and fluxes) in three-dimensional space. Such models
26 focus on preserving and tracking “*mass and energy (sometimes also momentum) flows*”
27 through the system.

28 [4] On the other hand, if a data-based machine-learning type of representation is
29 implemented, then the focus is on preserving and tracking “*information flows*” through the
30 system. In this case, the model may not be as well suited (as a process-resolved
31 representation) to inferring state variables such as “*soil moisture*” or fluxes such as
32 “*percolation, recharge and interflow*” that are constrained (by physics) to obey conservation
33 principles, unless appropriate regularization constraints are also implemented.

34 [5] In summary, the representational structure (of both the data and the model) selected for
35 the analysis imposes strong constraints on the questions we can ask, the results we can get,
36 and the inferences we can reasonably hope to perform. Consequently, we can expect, a priori,
37 that different representational strategies may provide different perspectives on the factors and
38 processes governing the generation of system behaviors.

39 1.2 Models as Complementary Perspectives on Reality

40 [6] For any given application, it can be challenging to determine what the most appropriate
41 model structure might be. In hydrology, as in other fields, this situation has led to the
42 availability of a very large variety of models, each based on different assumptions (and even
43 philosophies), and often having been tested under different (sometimes very specific)

44 conditions. This diversity of modeling approaches brings to mind the classic story of the
 45 “*blind people and the elephant*” where each person’s interpretation of what constitutes an
 46 “*elephant*” is based on their experience being limited to some very specific aspect of the
 47 animal, while also being limited by their ability to map that experience onto their previous
 48 knowledge (i.e., they are limited by what they can “*recognize*”).

49 [7] Given that any model is a “*relevant simplified representation*” of the world, where the
 50 simplifications typically reflect personal biases and preferences, it is quite possible for there
 51 to be as many viable “*representations*” to choose from as there are people working on a given
 52 problem. In practice, however, only some of these viable representations will tend to perform
 53 well (when evaluated against data), thereby considerably decreasing the number of
 54 potentially suitable options. Nevertheless, it can remain difficult to identify a single “*best*”
 55 model, since many different representations can be found to exhibit similar levels of
 56 performance (*Clark et al, 2011*).

57 [8] Returning to the metaphor of the “*blind people and the elephant*”, rather than asking
 58 which of the representations is (somehow) “*the best*”, one might instead consider whether
 59 the multiple complementary perspectives offered by the different representations can provide
 60 information that can be used to develop a better overall understanding of the system under
 61 investigation. By taking a multi-representational perspective, within which each
 62 interpretation of the system is deemed to be valuable (in at least some partial sense that
 63 contributes to a more complete overall point of view), we can hope to make progress towards
 64 uncovering the “*real*” nature of the underlying *Data Generation Process* (DGP).

65 **1.3 The Potential Offered by Lumped Water Balance Modeling**

66 [9] Lumped water balance models, that are structurally and behaviorally isomorphic to the
 67 system, are the mainstay of how understanding is developed in science in a very simplified
 68 (conceptual) representation. Such representations enable theoretical prior knowledge (such
 69 as conservation and thermodynamic principles) to be imposed as physical constraints on the
 70 allowable input-state-output trajectories of a system.

71 [10] The development of such models is structured as a sequence of conditional hypotheses,
 72 beginning with the governing conservation laws, and proceeding through the specification of
 73 the system architecture, process parameterization equations, and property-to-parameter
 74 relationships; see extensive discussion in *Gharari et al (2021)*. Different choices at each stage
 75 of development can give rise to different “*compound hypotheses*”; for examples of modular
 76 modeling systems in hydrology, see (*Fenicia et al, 2011; Clark et al, 2015; Craig et al,*
 77 *2020*). This facilitates a multi-hypothesis approach to scientific investigation (*Clark et al,*
 78 *2011*), in which each model represents a point within a hypothesis space that is strongly
 79 constrained by physics, assumptions, and prior knowledge.

80 [11] The strength of this representational approach is the ability to constrain discovery to
 81 model structures that are consistent with physical principles. An important consequence is
 82 that “*meaning*” can be ascribed to the various components, fluxes, and state variables of the
 83 model, making it possible to transfer understanding between locations, and to generalize to
 84 classes of systems that share similar representational properties.

85 [12] However, this strength can become a weakness when, by imposing strong priors, we
 86 limit the ability of the model to learn explicitly and directly from data, and to discover things

87 that are inconsistent with the space of hypotheses explicitly covered by the priors (*Gharari*
88 *et al 2021*).

89 **1.4 The Potential Offered by Machine Learning (ML)**

90 [13] Conversely, due to its ability to extract complex relationships from large datasets,
91 Machine Learning (ML) has gained a reputation for being able to help address some of the
92 most challenging tasks in science, particularly where theoretical understanding is lacking or
93 is weak. Recently, applications to hydrology have demonstrated excellent results in different
94 areas. To mention just a few, *Long-Short Term Memory* networks were successfully applied
95 to large-scale streamflow prediction (*Kratzert et al, 2018*), to a long record for one catchment
96 (*Hu et al, 2018*), to the estimation of water table depth for five agricultural areas (*Zhang et*
97 *al, 2018*), and to 5-day-ahead prediction (*Sudriani et al, 2019*).

98 [14] Overall, the power of ML-based representations arises from their theoretical and
99 practical ability to approximate any input-state-output mapping to an arbitrarily high degree
100 of accuracy, given sufficient data. Consequently, ML has emerged as a powerful
101 complement/alternative to the hypotheses-driven process-based approach to hydrological
102 modeling. However, as with the lumped water balance approach, each ML
103 algorithm/approach is based on a different mathematical perspective about how to represent
104 the structures underlying a given data set, and/or on how to represent and extract information
105 contained in the data. Accordingly, when different ML algorithms are applied to a given data
106 set, each is (also) likely to provide a different and, in general, complementary perspective on
107 the underlying nature of the DGP. By understanding how different ML algorithms represent
108 and extract information from data, we can seek to understand the particular value offered by
109 each perspective and exploit it to obtain a more comprehensive picture of the underlying
110 system.

111 **1.5 Objectives and Scope of this Paper**

112 [15] The objective of this paper is to explore how a multi-representational approach can help
113 to extract relevant information from a dataset, with a view to improving prediction,
114 understanding, and discovery. Our specific goal is to use such an approach to develop an
115 understanding of catchment hydrology across the hydro-geo-climatologically diverse extent
116 of Chile. Rather than the traditional strategy of implementing a single pre-selected
117 computational model code to the entire country, or perhaps a different model code to
118 hydrologically different parts of the country, we implement three complementary
119 representational approaches (model structures) across the entirety of Chile. These include a
120 a lumped water-balance model (based in a physical understanding of watershed behavior)
121 and two machine learning models (based on information extracted from historical
122 observations). Our focus is on understanding the strengths and weaknesses associated with
123 each representational approach, and on exploring the potential richness of inferences that a
124 multi-representational approach can support.

125 [16] In the next section, we introduce the problem of catchment-scale hydrological
126 forecasting in the context of the particular hydro-geo-climatology of Chile. Section 3 will
127 discuss the study methodology. The study results are presented in Section 4. Finally, we
128 provide a discussion and some thoughts about the implications of this work in Section 5. To
129 be clear, this study should be considered to be exploratory, with a view to improving our

130 understanding of how a multi-representational approach can be exploited in the service of
131 enhanced scientific discovery.

132 **2 THE CHALLENGE OF STREAMFLOW PREDICTION ACROSS** 133 **HYDROLOGICALLY DIVERSE CHILE**

134 [17] Prediction of streamflow at national scales is challenging, due to the multitude of
135 relevant factors that can vary simultaneously across time and space. In particular, the ability
136 of hydrological models to generalize can be poor in regions where the spatial variability of
137 dynamical forcings and static attributes is large (*Malone et al., 2015*). This is especially
138 relevant to Chile, which is characterized by tremendous geo-hydro-climatic variability, both
139 along its 4,270 km (2,653 mi) North-South extent and also from East to West (*Figure 1*). At
140 one extreme, Northern Chile is home to the driest desert in the world, containing regions
141 where no precipitation has been recorded for more than 25 years. At the other extreme, more
142 than 5000 mm/year of precipitation has been recorded in parts of the south, where there are
143 also permanent icefields.

144 [18] Bordered by the Pacific Ocean to the West and Argentina to the East, the country
145 averages just 175 km (109 mi) in width, while the North-South running Andes Mountain
146 range rises to the highest elevation in South America (6,959 m or 22,831 ft). Moreover, a
147 second mountain range, with lower elevations, runs parallel to the coast along almost the
148 entire country. Owing to the high elevations of the mountain ranges, precipitation in the
149 headwater catchments occurs mainly as snow, due to which the corresponding streamflow
150 peak will appear many days or even weeks after the precipitation event. In contrast, where
151 liquid precipitation occurs in catchments with high slopes, the times of concentration can be
152 shorter than one day. Other factors, including the variability of forest fraction, degree of
153 human intervention, and valleys created between the two mountains range, are also strongly
154 related to the availability of water in the long term.

155 [19] This immense variability in geo-hydro-climatic conditions poses a considerable
156 challenge for any modeling system, and especially for lumped water balance representations
157 where the model structure must be selected in advance. As such, Chile presents a perfect
158 opportunity to explore the possibility of developing modeling techniques that can deal with
159 large geo-hydro-climatic variability, and even exploit it to achieve better model performance.

160 **3 STUDY METHODOLOGY**

161 [20] This section presents and discusses our study methodology, including the dataset used
162 (section 3.1), three representational methods used (section 3.2), and issues related to the
163 experimental design (section 3.3).

164 **3.1 Dataset**

165 [21] For the purposes of this study, we will use mainly the information provided by the
166 catchment-scale CAMELS-CL dataset (*Alvarez-Garreton et al, 2018*). This dataset includes
167 11 variables and 105 categorical and numerical attributes for 516 Chilean catchments. For
168 model development and evaluation, we selected 322 catchments selected to span the country
169 and to have a minimum streamflow record length of 7 years. The literature suggests that 2-3
170 water-years of daily data represents a minimum record length for calibration of conceptual
171 process-resolved models (*Gupta and Sorooshian 1985*) while around 8-10 years may be

172 required to ensure some degree of stability with respect to the estimated model (*Vrugt et al.*,
173 *2006*). On balance, therefore, 7 years represents a reasonable tradeoff between the availability
174 of the model development and spatial representation of catchments. Note also that the time-
175 periods of model development data selected for each catchment are not necessarily identical
176 or even overlapping, they simply represent whatever is available for those catchments. More
177 details on how the data were selected and partitioned appear in *De la Fuente (2021)*.

178 **3.2 Representations Examined**

179 [22] To develop an improved understanding of the nature of catchment-scale hydrology at
180 the national scale across Chile we will use a multi-representational approach. Clearly, this
181 approach must be consistent with the available data (Section 3.1). With this in mind, we
182 chose to investigate three complementary representational strategies – one being a lumped
183 water balance model, and the other two being ML-based modeling strategies. While
184 additional representational strategies could also have been included, these three arguably
185 represent sufficiently different approaches to extracting information from data to support the
186 objectives of this study.

187 [23] For the lumped water balance model, we chose the GR4J dynamical systems model
188 (*Perrin et al, 2003*), due to its relative parsimony and the reports of good performance in
189 other studies (*Kunnath-Poovakka & Eldho, 2019; Sezen & Partal, 2019; Pagano et al.,*
190 *2010*), and because the catchment-scale data required for its implementation is available (see
191 *Appendix Table A-2*). For the ML-based modeling strategies, we selected the LSTM network
192 (*Hochreiter & Schmidhuber, 1997*) and the RF regression-tree algorithm (*Breiman, 2001*).
193 Further details about these modeling strategies are provided below.

194 **3.2.1 The GR4J Lumped Water Balance Representation**

195 [24] The GR4J model is a parsimonious lumped water balance (process-based)
196 representation of daily time-step spatially-lumped catchment-scale hydrology, whose input-
197 state-output behavior can be controlled by adjusting four tunable parameters (*Figure 2a*).
198 GR4J is the outcome of several studies that evaluated a variety of model structures having
199 various levels of complexity, using data from 429 catchments with differing geo-hydro-
200 climatic conditions. While more complex representational structures are available, such as
201 GR5J (*Le Moine, 2008*) and GR6J (*Pushpalatha et al, 2011, CemaNeige; Valéry, 2010*),
202 GR4J provides a relatively simple representative of this class of models, and a search for the
203 “best” such model for Chile is not part of the scope of this study.

204 [25] Importantly, the GR4J representation seems well suited to application across Chile for
205 several reasons, including the fact that it contains a flow path without storage that gives it the
206 ability to simulate the very rapid precipitation-streamflow response that is characteristic of
207 the steep surface slopes that occur across portions of Chile. Further, it has the desirable
208 feature that it includes a parameter that can be tuned to permit the model to either “import”
209 or “export” water into its main routing storage tank to enable the model to better match the
210 input-output behavior of the system as expressed by the observed time-series data.

211 [26] As such the GR4J model serves as a kind of lumped water balance benchmark against
212 which the ML-based representations can be compared (no pre-existing benchmark model
213 calibrated across the entire country exists). Following the traditional approach, the
214 parameters of the GR4J model are tuned (calibrated) “locally” to be specific to each
215 catchment. This is in contrast to the ML-based approaches (see below) where the parameters

216 (network weights and/or biases) are tuned “*globally*” to represent all catchments across the
 217 study domain (in our case the entire country of Chile). Because GR4J is calibrated locally,
 218 it can be considered to represent a performance benchmark that one would want a globally-
 219 tuned data-based ML approach to be able to exceed, particularly when requiring the model
 220 to generalize well (i.e., to perform well on catchments that are withheld from the model
 221 development data set). For more details on the structure of the GR4J model, please see [Perrin](#)
 222 [et al \(2003\)](#).

223 3.2.2 *The LSTM ML-based Representation*

224 [27] The LSTM model ([Figure 2b](#)) is a fully connected Recurrent Neural Network (RNN)
 225 with the ability to learn from the past. The recurrence feature is akin to that of the tank-like
 226 components of physically-based catchment models, but where the nature of the relationships
 227 between inputs, state variables, and outputs is learned from the data. The structure of an RNN
 228 is based on the linear superposition of non-linear basis functions (known as activation
 229 functions), which enables the network to closely match the nature of the true relationships
 230 that underlie the data.

231 [28] The LSTM-based representation is somewhat more complicated than the traditional
 232 RNN because additional components (called gates, or gating functions) are used to
 233 contextually control how the flow of input, output, and previous state information affects the
 234 network response at each time step. In this sense, the LSTM cell-states can be thought of as
 235 non-linear extensions of the classic linear reservoir component commonly used for
 236 hydrological modeling (see [Table 1](#)). As such, the LSTM can be interpreted as a
 237 representation of Dynamic Information Storage, where the gating functions act as
 238 contextually variable resistances to the flows of different kinds of information. For more
 239 detail about the equations used in the LSTM model, please refer to [Kratzert et al. \(2018\)](#).

240 [29] By drawing an analogy with the linear reservoir ([Table 1](#)), it is possible to interpret each
 241 of the functions and components in the LSTM. The function $g(x,h)$ computes a linear
 242 combination of inputs and outputs, adds a bias term, and then non-linearly transforms the
 243 result onto the range $[-1, 1]$ using the hyperbolic tangent equation. The quantity $c(t)$ can be
 244 understood as an informational “*state variable*” of the system, where the information carried
 245 by $c(t)$ is contextually updated based on the values of weights $f(t)$ and $i(t)$ applied to the past
 246 storage information and drivers, respectively. Whereas the linear reservoir computes the
 247 output as a constant proportion of the storage $S(t)$, the LSTM defines this proportion
 248 dynamically through the gating function $o(t)$. In summary, the gating functions $f(t)$, $i(t)$, and
 249 $o(t)$ act as dynamic amplification factors that can take on values between $]0,1[$ controlled by
 250 a sigmoid-shaped activation function. Meanwhile, the variables $g(t)$ and $c(t)$ represent
 251 different informational representations of the input, output, and storage, normalized using the
 252 hyperbolic tangent function to take on values between $]-1,1[$.

253 [30] Due to the isomorphic relationship of the LSTM to the linear reservoir, it becomes
 254 possible for the structure of the GR4J model to be emulated using an LSTM. However,
 255 because the LSTM can exploit the information provided by inputs beyond precipitation and
 256 potential evapotranspiration, it becomes more difficult to apply physically-based
 257 interpretations to the behaviors of its state variables.

258 3.2.3 *The RF-based Representation*

259 [31] The RF is a regression-tree methodology (*Figure 2c*) that adopts the classic strategy of
260 “*divide and conquer*” to construct an approximation of the input-output mapping expressed
261 by the data. The RF algorithm searches for an “*optimal*” partitioning of the input space, that
262 maximizes the similarity within each output cluster that represents a leaf of the decision tree.
263 The similarity measure used is the weighted average dispersion of the outputs within a
264 cluster, where dispersion is measured as the sum of squared deviations from the mean of the
265 cluster members (L2 norm). It is typically assumed that the input-space partition resulting
266 in the smallest average dispersion, weighted by the number of elements, is best. This process
267 is repeated within each partition until a prespecified minimum number of elements remains
268 within a subset, and/or until a predefined number of splits have been conducted.

269 [32] The RF implements a piecewise-constant approximation of a complex continuous input-
270 output mapping, where for each new split, we look for the minimum difference between the
271 cluster means (predictions) and their corresponding target values (output data). To avoid this
272 deterministic process becoming highly biased by the specific data sample used to construct
273 the decision tree, the RF approach uses a random sample (selected with replacement) from
274 the data set to construct each decision tree and repeats the process multiple times to generate
275 a “forest” (ensemble) of decision trees. The use of randomized ensembles helps to reduce
276 overfitting, while randomness in the input selection helps to improve the accuracy of the
277 classifier and regressor algorithm (*Breiman, 2001*). The final prediction is generated as the
278 average prediction made by each of the trees in the forest. From the perspective of
279 interpretability, it is easy to examine each input-space split associated with the model
280 predictions, making it relatively easy to understand the steps connecting the input to the
281 outputs, without the need to track any intermediate variables or states.

282 [33] Another issue is related to the nature of the problem. Because streamflow is the result
283 of complex processes within a Dynamical Environmental System, knowledge of the system
284 state can be very important for characterizing how the system will respond to new inputs
285 given past conditions. In other words, the streamflow on a specific day is not just the result
286 of what happens on that day but also depends on what has happened in the past. Moreover,
287 the history of a catchment can be understood in terms of different time scales, such that the
288 current streamflow response is related to both what happened recently (short-term memory)
289 and what has happened in the past months or even years (long-term memory), due to
290 persistence in the behavior of the system. These two kinds of memory are not explicitly
291 represented by the structure of the RF. Accordingly, it is the responsibility of the modeler to
292 ensure that the input data contain variables that provide such memory-related information
293 that can be used by the RF in place of state variables to track both the short-term and long-
294 term dynamics of the system. This feature can be interpreted as both a “*pro*” and a “*con*” of
295 the RF approach; it is *pro* in the sense that it allows the modeler to exert better control over
296 the model by injecting physical understanding, while it is a *con* because it imposes higher
297 demands in the form of data preprocessing.

298 [34] The ML-based RF and lumped water balance GR4J representations are harder to
299 compare (than the LSTM and GR4J representations) because of RF’s lack of state variables
300 that mediate between the inputs and outputs, and because of the piecewise-constant nature of
301 the RF representation. However, this does not mean that memory and dynamics are not
302 considered by the RF, because state variables can be thought of as representing the aggregate

303 effects of an infinite number of past system inputs. Accordingly, provided that the RF is fed
304 with a sufficiently long history of past system inputs, the representation can learn to construct
305 input-space splits that emulate those that would have resulted from the tracking of state
306 variable information.

307 **3.3 Experimental Design**

308 [35] The main challenge to creating a unified model development methodology is that each
309 of the three representational strategies has different conceptual, mathematical, and coding
310 characteristics, and therefore different structures and processes of implementation, that must
311 be followed to obtain an operational model. It is, therefore, impossible to implement an
312 entirely uniform methodology for model development. Accordingly, we followed the
313 reasonable approach of implementing the recommended best model development practices
314 for each representational type and compare the results so obtained. However, the overall
315 methodology conforms to a common framework, so as to enable comparative analysis.
316 Accordingly, all comparisons are based on the use of the same data and performance metrics
317 for model development and evaluation.

318 [36] *Appendix A1.1* discusses how we partitioned the data into , that was done using three
319 periods for the purposes of model calibration, selection, and evaluation, consistent with the
320 ML literature. The *Appendices A1.2* and *A1.3* summarize the variables and attributes used in
321 model development. *Appendix A1.4* discusses how the warm-up period used for both GR4J
322 and LSTM was selected. *Appendix A1.5* discusses the process of parameter/hyperparameter
323 selection for each model. *Appendices A1.6* and *A1.7* describe the metrics and algorithms
324 used, and *Appendix A1.8* describes how the out-of-sample testing dataset was generated.

325 **4 EXPERIMENTAL RESULTS**

326 [37] In keeping with the objectives of this paper, our analysis pays special attention to how
327 the different representational approaches can be used to make inferences regarding various
328 characteristics of the hydrological processes that underlie the data. Accordingly, this section
329 consists of two parts. In section 4.1 we investigate issues of overall understanding and/or
330 discovery, such as system memory and feature importance, enabled by the multi-
331 representational approach. In section 4.2 we investigate how the different models
332 constructed using the different representational approaches performed in terms of the ability
333 to generalize in space and time.

334 **4.1 Understanding Enabled by the Multi-Representational Approach**

335 [38] Each representational approach responds differently to the fluxes of information through
336 the system, and that response can provide useful insights into the characteristics of the
337 system. This happens because each time that the representation assimilates a new piece of
338 information, it updates its internal structure/parameters, which can be understood as *learning*
339 about *changes in the internal state* of the system. Therefore, the final “learned” version of
340 the representation (the trained model) encapsulates a considerable amount of information that
341 can be subject to analysis. Here, we investigate how data skewness, system memory, and the
342 relative importance of surrogate variables provide insights into the underlying nature of the
343 DGP.

344 4.1.1 The Box-Cox Transformation Parameter

345 [39] The Box-Cox power transformation $Y = (y^\lambda - 1)/\lambda$ is commonly used in statistical
 346 analysis (*Box and Cox, 1964*) to account for skewness in the data (represented here by y). If
 347 $\lambda = 1$ the variable Y has essentially the same distributional properties as y (no transformation
 348 beyond a shift of origin), while by setting λ to be smaller or larger than 1.0, the skewness of
 349 Y can be reduced or increased, respectively, relative to y . In hydrology, strong skewness of
 350 the streamflow distribution (corresponding to values of $\lambda \rightarrow 0$) is indicative of precipitation
 351 being the main driver of system dynamics, while weak streamflow skewness may indicate
 352 the dominance of groundwater, snowmelt, or other processes that act as low-pass filters in
 353 the generation of streamflow dynamics.

354 [40] In our implementation of the GR4J model, the λ parameter was allowed to vary with
 355 location to allow the model development process to account for the hydro-geo-climatic
 356 variability of skewness in the streamflow data. *Figure 3* presents the distribution of spatially-
 357 varying ‘optimal’ λ values obtained for the 322 catchments in the model development dataset.
 358 While the optimized λ values vary across the full range tested (0.00 to 2.00), a high
 359 concentration (~35%) of values fall in the first bin ($\lambda \sim 0$), consistent with the traditional use
 360 of a logarithmic transformation when calibrating daily-time-step models in hydrology
 361 (*Hassan & Hassan, 2021*). This result is consistent with the large range of possible
 362 streamflow values that can occur across Chile, where streamflow can vary over several orders
 363 of magnitude in catchments that respond quickly to intense precipitation events. In contrast,
 364 ~10% of the catchments are associated with $\lambda > 1$, where very little variability in the range
 365 of streamflow magnitudes can be found (e.g., where precipitation-runoff is weak, or where
 366 baseflow tends to be the dominant streamflow generating mechanism).

367 [41] In the implementations of the LSTM and RF models, single values of lambda were
 368 applied across the entire country. *Figures 4a* and *b* present cumulative density functions
 369 (CDF) of KGEss metric performance (*Gupta et al 2009, Knoben et al 2019*; see definitions
 370 in *Appendix A1.7*) computed over the selection period (validation), for different choices of
 371 lambda. For each CDF, we report the area under the curve to serve as guidance for selecting
 372 the best value of lambda (treated as a hyperparameter), where smaller areas correspond to
 373 better overall performance. In contrast with GR4J, where the average value of λ is close to
 374 zero, *Figure 4a* indicates that better performance of the RF model is obtained with λ close
 375 to one, which corresponds to *not* applying a transformation to the streamflow data. This
 376 makes sense in retrospect, because when the decision tree splits the data into different clusters
 377 it is inherently able to account for skewness, and so the addition of a transformation does not
 378 necessarily provide any significant additional value to the model development process (note
 379 the relatively weak dependence of performance on λ). Based on this observation, when
 380 developing the RF model, we fixed the value to $\lambda = 1$ (corresponding to no transformation)
 381 and report only results obtained using this value.

382 [42] Note that when implementing the LSTM, we obtained better results (*Figure 4b*) by
 383 using a ‘global’ standardization of the data (subtracting the mean and dividing by the
 384 standard deviation of streamflow values computed from the entire Chile-wide dataset), rather
 385 than a ‘local’ standardization (where the mean and standard deviation were regressed against
 386 aridity index). Further, *Figure 4b* indicates that the performance of the LSTM model so
 387 obtained is not sensitive to the choice of λ , and so we again fixed the hyperparameter $\lambda = 1$

388 (corresponding to no transformation) for both standardization approaches, and report only
389 results obtained using this value.

390 [43] These results are interesting because the value of λ encodes information about the DGP
391 when using the GR4J representation but not when using the ML-based representations. One
392 might speculate that this result is a consequence of the fact that the ‘*local*’ GR4J modeling
393 approach permits each catchment to be represented by a different value for λ , while the
394 ‘*global*’ ML-based modeling approaches require the specification of a single country-wide
395 value. However, if this were true we might expect the ML-based ‘*globally optimal*’ λ values
396 to converge to something like the mean or median of the GR4J-based distribution of ‘*locally*
397 *optimal*’ λ values. Given that this is not the case for both the (quite different) LSTM and RF
398 representational approaches, it is more likely that the ML models are internally able to
399 address problems related to data skewness in some other manner. Nonetheless, this remains
400 an interesting issue for future investigation.

401 **4.1.2 System Memory**

402 [44] For the ML-based representations, an important property is the manner in which the
403 “*system memory*” is characterized, in terms of the number of previous time-steps of input
404 data (meteorological variables) that are determined to provide useful information about the
405 current value of streamflow. Note that this “*lag-time*” hyperparameter is not relevant to the
406 GR4J representation, which tracks system memory exclusively through its state variables.
407 **Figures 5a** and **b** show how the CDFs of the model performance vary with different values
408 of the lag-time hyperparameter for the RF and LSTM models respectively.

409 [45] Consider first the RF model. For catchments with KGEss better than ~ 0.45 , the CDFs
410 move progressively to the right (indicating improved performance) as the lag-time is
411 increased from 2 to 32 days, whereas for catchments with KGEss less than ~ 0.45 the results
412 are insensitive to the value of the lag-time hyperparameter. Further, the marginal
413 performance improvement declines, on average, as the lag-time is increased. A similar result
414 is found for the LSTM model, but now we see an additional region of improvement when
415 going from 128 to 270 days, occurring mainly in catchments with KGEss values below ~ 0.85 .

416 [46] Taken together, these results suggest that the ML-based models are detecting the
417 expression of two different kinds of processes giving rise to streamflow generation across
418 the country, one related to a system “*memory*” of around 32 days and the other related to one
419 of around 270 days. We will revisit this topic in the next section, where we see that this
420 difference in length of system memory is correlated with hydro-geo-climatic attributes. Note
421 that this kind of information about systemic differences between different catchments in the
422 study region is somewhat more difficult to infer from the relatively simple GR4J
423 representation used in this study.

424 **4.1.3 Feature Importance**

425 [47] Another interesting aspect of ML-based approaches is the manner and flexibility by
426 which the relative informativeness/importance of different (spatially-varying) hydro-geo-
427 climatic attributes can be assessed. As a consequence, it is (in general) easier to detect which
428 attributes are more/less important when explaining the predictive power of an ML-based
429 model. Whereas this is, in principle, also possible using a conceptual/lumped modeling
430 approach, such an inference would have to be done indirectly through an analysis of the

431 spatial patterns of calibrated values of the model parameters, which is arguably a less direct
432 and somewhat more complicated process.

433 [48] In particular, for ML-based models, tools such as the Scikit-learn Python library
434 (*Pedregosa et al., 2011*) facilitate a simple sensitivity analysis that permits a relatively
435 straightforward exploration of the importance of each input variable or system attribute in
436 contributing to the predictions. For conceptual/lumped representations, this exploration is
437 complicated by the fact that the importance/informativeness of a given variable or attribute
438 is mediated by the specific structural assumptions encoded by the system architecture and
439 process parameterization equations chosen for the model. In contrast, the ML-based
440 representational structure is not quite so strongly pre-determined and is therefore, arguably,
441 less likely to bias any inferences of relative feature importance. Of course, this is not *entirely*
442 true since different ML-based approaches also (unavoidably) encode different
443 representational assumptions about how to map system inputs to outputs. However, the
444 relative flexibility of ML-based representations (as well as their focus on the strengths of
445 “*informational*” relationships) should, in principle, enable interesting (and hopefully useful)
446 insights regarding relative feature importance to be inferred. That inference is neither as
447 simple or as direct as in a lumped water balance model, due to the fact that attributes and
448 variables are connected through hundreds (even thousands) of parameters. Nonetheless, tools
449 such as the one mentioned above are increasingly making such analysis possible.

450 [49] The consequence is that an ML-based assessment of the relative importance of hydro-
451 geo-climatic attributes can provide potentially valuable information regarding what system
452 attributes are likely to be important when constructing a (better) conceptual/lumped model,
453 and regarding how these attributes are likely to vary across space, and must therefore be
454 considered in order to achieve a lumped water balance representation that generalizes well at
455 the large scale.

456 [50] Here we analyze the information about feature importance that is an inherent property
457 of the RF-based model, where the variables selected for data-space thresholding earlier in the
458 tree (e.g. at the first split) can be interpreted as being more ‘*fundamentally*’ or ‘*globally*’
459 important to the construction of a decision tree. *Figure 6a* indicates that the most important
460 attribute is an “*aridity*” index (aridity_cr2met, computed using the CR2MET precipitation
461 product), which strongly suggests that the form of the relationship between the availability
462 of water and the generation of streamflow is different in different parts of the country (e.g.,
463 in humid versus arid regions). While this observation is not novel (*Neto et al, 2020; Booij et*
464 *al, 2019, Chen et al, 2019*), it is consistent with the fact that lumped catchment-scale water
465 balance representations, with their fixed architectures and process-parameterization
466 equations, are typically not able to generalize well across different hydro-geo-climatic
467 conditions.

468 [51] Of course, this does not imply that failure to account for “*aridity*” is, per se, a complete
469 and meaningful explanation for poor performance of any given model type. In general,
470 spatio-temporal changes in aridity index are likely to simply indicate relative changes in the
471 importance of various drivers of streamflow. From *Figure 6a* we see that the second, fourth,
472 fifth, and seventh most important attributes are daily precipitation values, which indicates
473 that the behavior of the RF model is mainly controlled by aspects related to precipitation,
474 once aridity has been accounted for. This is consistent with the interpretation for the Box-
475 Cox transformation parameter λ used with the GR4J model.

476 [52] In our case, the first split (*Figure 6b*) that occurs in most trees of the RF model occurs
477 at an aridity index threshold of 0.6 mm/day. While, for any given study area, the precise value
478 at which this split occurs will depend on the distribution of wet and arid catchments in the
479 dataset, this observation suggests that different streamflow generating representations may
480 be required for the model to perform well in regions that are “*energy-limited*” as opposed to
481 “*water-limited*”. When a similar analysis is performed for the precipitation threshold, we find
482 that the nature of the streamflow response is different for values above/below ~10 mm/day.
483 From this, we could hypothesize that more than ~10mm/day of precipitation is required (on
484 average) to generate surface runoff, but of course, much more analysis would need to be done
485 to test such a hypothesis.

486 [53] Finally, we note that of the top ten most important attributes, the only ones that are not
487 related to aridity and/or precipitation are the “*month of year*” (Month) and “*forested fraction*”
488 (nf_frac). The month of year attribute conveys information related to hydro-climatic cycling
489 (annual periodicity), whereas the forest fraction conveys (among other things) information
490 about infiltrability and soil water retention capacity of the soil.

491 [54] The main point of these two rather simple (even trivial) examples shown in *Figure 6* is
492 that the RF representation facilitates a kind of analysis that can provide interesting
493 information that is not easily obtained using either the GR4J or LSTM representational
494 approaches. In this sense, the RF approach provides a strong complement to other
495 representational approaches when our goal is to use modeling in support of scientific
496 discovery and understanding.

497 **4.2 Comparative Analysis of Similarities and Differences in Performance**

498 [55] The relative ability of any properly trained model to perform well on independent
499 “*evaluation period*” datasets can be considered indicative of well the corresponding
500 representational approach supports discovery about the underlying DGP. However, even if
501 all of the models tested on the evaluation period provide essentially identical values for some
502 aggregate performance metric (such as KGEss or NSE; see Appendix A.1), deeper analysis
503 may reveal systematic differences in model simulated behaviors that the aggregate metric is
504 not capable of distinguishing between (*Gupta et al 2008, 2009*).

505 [56] For our rainfall-runoff modeling case study, such behaviors may include things such as
506 the simulated-to-observed long-term water balance and variability ratios, and the timing and
507 shape (measured, for example by cross-correlation strength between the simulated and
508 observed time-series of model output response). By examining how well each
509 representational approach reproduces such behaviors (when trained, as is customary, on an
510 aggregate performance metric), we can hope to obtain insights into the strengths and
511 weaknesses of each, which is the objective of this paper. This section investigates overall and
512 spatial patterns of such differences in model behavior/performance, with a view to
513 understanding the manner and extent to which the models (developed using different
514 representational approaches) are able to generalize well in space and time.

515 **4.2.1 Overall Performance**

516 [57] First, we examine the distributions of overall model performance across the country.
517 *Figure 7a* shows the CDFs of evaluation period performance (as measured by KGEss) for
518 all locations where $KGEss > 0$ (where predictions are, on average, better the “no-model”
519 prediction that simply uses the observed mean; *Knoben et al 2019*). Similar results were

520 obtained using NSE (not shown; for details see *De la Fuente, 2021*). Two interesting points
521 can be noted:

522 1) The LSTM curve (blue line) is significantly further to the right (~85% of the
523 catchments) over most of the range, indicating statistically better overall performance.

524 2) The GR4J model fails to meet the $KGE_{ss} > 0$ threshold at only ~5% of the
525 catchments, as opposed to ~11% for LSTM and ~22% for RF.

526 [58] Regarding the first result, the superior performance of the LSTM model over most of
527 the range is (arguably) expected given that the LSTM can both a) explicitly learn about
528 system dynamics and memory through its representation of state variable recurrence, and b)
529 learn the functional form of the input-state-output mapping due to its structural flexibility.
530 Note that the former ability is not explicitly enabled by the RF architecture (green line), while
531 the latter ability is not possible for the fixed GR4J architecture (red line).

532 [59] Regarding the second point, given that all three representations are trained using
533 (almost) the same input-output information (GR4J model uses only precipitation and
534 evapotranspiration), this result suggests that there are hydro-geo-climatic conditions under
535 which the GR4J representation provides useful (lumped water balance) information that is
536 not directly inferable from the available data by the LSTM and RF representations. Of course,
537 whether this benefit comes from the specific mass-conserving and process-equation nature
538 of the GR4J architecture, or from its ability to compensate for mass-balance errors by
539 importing/exporting groundwater (or some other reason) is not immediately clear, and will
540 require more detailed investigation. In a recent study by *Hoedt et al. (2021)*, a “*mass-*
541 *conservative*” LSTM model was found to be able to learn a good state-variable representation
542 of the dynamics of snow storage, but such findings would need to be tested at larger scales
543 over a variety of hydro-geo-climatic conditions before more general conclusions can be
544 drawn.

545 [60] Next, we examine the distributions of the decomposition components of KGE_{ss} (see
546 definitions in Appendix A1.7). While aggregate metrics such as KGE_{ss} (and NSE etc.) can
547 provide a good overall idea of model performance, they can often be poor at revealing
548 important differences in characteristic model behaviors, particularly when overall
549 performance is poor (*Gupta et al., 2009*). *Figure 7b* provides further discriminatory
550 information by plotting the CDF of model Bias Ratio, where values larger (smaller) than
551 $10^0 (= 1)$ indicates a tendency to overestimation (underestimation).

552 [61] This plot reveals that the GR4J and LSTM models, that have the explicit ability to
553 simulate system dynamics, tend (on average) to be unbiased, whereas the RF model tends to
554 be positively biased. Interestingly, for situations where the models tend to overestimate the
555 mean (Bias Ratio > 1.0), the GR4J model tends to do better (have lower bias) than the two
556 ML-based models, with the RF model being the worst. However, for situations where the
557 models tend to underestimate the mean (Bias Ratio < 1.0) that situation is reversed and the
558 two ML-based models perform better than GR4J, with the RF model being the best. Similar
559 results were found for the Standard Deviation Ratio (results not shown).

560 [62] So, while the LSTM is statistically superior in terms of overall KGE_{ss} performance for
561 the majority of catchments, the situation is clearly more nuanced – with each representational
562 type providing different characteristic abilities to simulate various attributes of streamflow,
563 despite the fact that all the model types were trained using (almost) the same data. This

564 supports our contention that a multi-representational approach can aid in scientific
565 investigation and discovery, particularly when faced with significant hydro-geo-climatic
566 variability.

567 [63] Meanwhile, the use of multiple metrics, that target different (ideally complementary)
568 signature properties of the data (*Gupta et al 2008*), can assist in the extraction of different
569 kinds of useful information, enabling inferences about different aspects of the input-(state)-
570 output response of the system.

571 **4.2.2 Spatial Patterns of Performance**

572 [64] The previous statistical analysis is informative about the overall properties and
573 capabilities of the different representational types. However, it is of little value when needing
574 to make statements about actual performance at any catchment. In this section, we investigate
575 how the different representational types perform across the variety of different hydro-geo-
576 climatic conditions that characterize Chile.

577 [65] *Figures 8a* and *b* explore the relationship between model performance and two
578 interesting hydro-geo-climatic factors – *Latitude*, and *Aridity*. Given the long narrow shape
579 and North-South orientation of the country, these two factors serve as useful surrogates for
580 hydro-geo-climatological variability, with the Northern extent of the country being
581 characterized by very dry conditions and high elevations, the Southern extent being
582 characterized by extreme precipitation and permanent icefields, and the central region being
583 characterized by intermediate degrees of wetness and considerable variability in elevation.

584 [66] The curves in *Figure 8a* show smoothed trajectories (using a moving average of 15
585 catchments) of the variation in KGEss performance with *Latitude* from South to North (left
586 to right across the x-axis). First, we see that, while all three models exhibit relatively good
587 performance in the mid- and south-central (moderately wet) parts of the country [latitude
588 -45° to -35°], performance of GR4J decreases sharply relative to the ML-based RF and
589 LSTM as we move to the southernmost regions [latitude -55° to -45°]. This decline in
590 GR4J performance makes sense given that the south is characterized by the existence of
591 glaciers and lakes, which can introduce significant time-lags into the dynamics of the system
592 that cannot easily be reproduced by the existing GR4J architecture. In contrast, the flexibility
593 of the ML-based representations enables them to better account for such phenomena.

594 [67] Meanwhile, all three models exhibit relatively poor performance ($\text{KGEss} < 0.5$) across
595 the north-central parts of the country [latitude -35° to -25°]. This region is characterized
596 by strong slopes (rapid elevation changes and very short times of concentration) and
597 relatively greater aridity (see next sub-section) than the mid/south-central and southern
598 regions. Here, RF performs particularly poorly, which may be attributable to the fact that it
599 does not have access to data with greater than 16 days lag time and is, therefore, unable to
600 account for longer (seasonal or annual time-scale) system memory, unlike GR4J and LSTM.

601 [68] Finally, the northern part of the country [latitude -25° to -18°] contains the Atacama
602 Desert, which is the aridest region in the world and has moderate slopes. Here, RF and LSTM
603 both exhibit better performance than GR4J. The inability of the latter to simulate the
604 hydrologic behavior of such extreme conditions is likely due to the fact that GR4J was
605 developed to represent the very different hydro-geo-climatic conditions that characterize
606 France. Meanwhile, the relatively poor performance achieved by both ML-based models

607 suggests that the variables that make up the existing CAMELS-CL dataset are not sufficiently
608 informative about the particular input-state-output dynamics of the catchments in this region
609 to enable a robust and accurate model to be developed, and that other variables and attributes
610 should be added to improve model performance (more on this later).

611 [69] The curves in *Figure 8b* show smoothed trajectories for how KGEss performance varies
612 with the *Aridity Index* (computed as the mean of *aridity_cr2met* and *aridity_mswep*). Here
613 we see a clear dependence of performance on aridity, with all three models exhibiting better
614 performance ($KGEss > 0.5$) under wet (i.e., energy limited) conditions but with performance
615 becoming progressively worse as the hydro-climatic conditions become increasingly more
616 arid (water-limited). Interestingly, the performance of both GR4J and LSTM (that have the
617 ability to simulate system dynamics) declines more or less linearly with increasing log-
618 aridity, but RF performance declines somewhat more rapidly and is significantly worse than
619 for GR4J and LSTM when the *Aridity Index* is between about 1.5 to 8.0. Given that GR4J is
620 designed to represent systems that are primarily driven by precipitation, it is understandable
621 that performance can decline as the direct dependence of streamflow on precipitation
622 becomes less, while the mediating effects of evapotranspiration and long-term groundwater
623 storage become more predominant.

624 [70] However, while the ML-based models have considerably more flexibility to discover
625 appropriate functional relationships in the data and would therefore normally be expected to
626 serve as indicators of upper-bounds on achievable model performance (*Nearing et al., 2020*),
627 they also show the same declining trend in performance with increasing aridity. This suggests
628 that the information content of the CAMELS-CL data set is biased towards a better
629 representation of the hydrological properties of wet (energy-limited) catchments and is
630 therefore not sufficiently complete to enable model development for arid parts of the country.
631 For example, it is noteworthy that the CAMELS-CL data set does not include information
632 about soil characteristics such as depth to bedrock, hydraulic conductivity, or soil fraction,
633 all of which are present in the US version (*Addor et al., 2017*), and which can be very
634 important in the characterization of the baseflow and streamflow-precipitation elasticity
635 (*Addor et al., 2018*).

636 [71] Another interesting observation is that the system memory associated with streamflow
637 generation from precipitation is different for energy-limited (wet) and water-limited (arid)
638 catchments. Referring to *Figure 5b*, we see that whereas the majority of catchments show
639 improvement of LSTM performance when provided with ~ 32 past days of input data
640 (reflecting short-time-scale memory processes), there is a smaller set of catchments with
641 poorer model performance that shows improvement only when provided with 270+ past days
642 of input data (reflecting longer-time-scale memory processes). The indication is, therefore,
643 that when investigating and modeling the streamflow response of catchments, our
644 representation – whether ML-based or conceptual/lumped – must contain structures that
645 make it possible to track memory processes at more than one dominant time-scale, depending
646 on the hydro-geo-climatology of the region. For example, one might consider the need to
647 track at least the short-term (weekly/monthly/seasonal), medium-term (annual), and possibly
648 longer-term (climatological) time scales. Of course, to discover and build representations of
649 longer-term (climatological) rainfall-streamflow response one would typically need more
650 than 7 years of data.

651 [72] The important point, however, is that the representational type selected for model
652 development should (ideally) make it possible for information about multiple hydro-climatic
653 time scales to be exploited. GR4J and LSTM contain explicit representations (through
654 dynamic state variables and multiple flow pathways) that – to some degree – facilitate this,
655 with the LSTM having a much greater degree of flexibility to do so (which may explain its
656 generally better performance in *Figure 8b*). However, for reasons explained in Appendix
657 A1.4, our implementation of the data-based RF only included data lagged up to 16 days,
658 which may explain why performance is worse than for the data-based LSTM when the aridity
659 index is on the range 1.5 to 8.0. Note that this kind of model-enabled analysis and discovery
660 is not easily achieved if only a purely conceptual/lumped approach had been used in this
661 study; by adopting a multi-representational approach that incorporates both
662 conceptual/lumped and a variety of complementary ML-based modeling strategies, the
663 process of analysis and discovery can be greatly enhanced.

664 [73] Finally, *Figures 9a-c* show evaluation-period KGEss performance for each of the three
665 models at each catchment used for model development (green indicates good performance,
666 yellow-orange indicates poor performance, and red indicates really bad performance).
667 Overall, all three models exhibit a tendency to good performance (KGEss >> 0.5) south of
668 latitude $\sim 33^\circ\text{S}$, and at the very northern tip of the country (north of latitude $\sim 20^\circ\text{S}$).

669 [74] Focusing specifically on the region between latitudes 27°S and 33°S , we see that RF
670 (*Figure 9c*) performs very poorly throughout this part of the country (see also *Figure 8a*).
671 However, LSTM performs quite well along a narrow strip of this region that borders
672 Argentina. This strip is located at higher elevations where temperatures are low and where
673 snowmelt processes dominate the generation of streamflow. The ability of LSTM to discover
674 and track longer-term memory processes is likely contributing to its good performance here.
675 As we move westward towards the coast, LSTM performance decreases, indicating that the
676 model no longer has access to the information needed to properly simulate the streamflow
677 response (which, in this case, is probably information about connections between
678 groundwater and streamflow). Turning to GR4J, we see that its KGEss performance across
679 the region is just barely better than 0.0, indicating that the model is mainly only able to
680 reproduce the long-term mean value of streamflow. Given that GR4J has explicitly neither
681 the ability to represent the dynamics of snow accumulation and melt nor the long-term
682 dynamics of groundwater, this result makes sense.

683 [75] *Figure 9d* indicates, for these same catchments, which model provides the best
684 evaluation-period KGEss performance (red=GR4J, blue=LSTM, green=RF) across the
685 country. Here we simply report the model with the best evaluation-period KGEss, regardless
686 of whether these KGEss values are statistically distinguishable. No clear pattern emerges,
687 but in general, the blue (LSTM) and red (GR4J) colors dominate, with LSTM generally being
688 the best-performing model across the country. This is consistent with the statistical results
689 (CDF plots) shown in *Figure 7a*.

690 [76] Some more nuanced findings emerge from a statistical analysis of KGEss performance
691 by model type, reported in *Table 2*. LSTM is the best performing model at 53% (172 of 322)
692 of the catchments, with an excellent median KGEss performance of 0.70. However, this
693 statistic masks the fact that where LSTM fails, it does so very badly – the worst KGEss value
694 is very poor and, consequently, the dispersion of performance is highly skewed. In contrast,
695 the distribution for GR4J, which performs best at only 30% of the catchments and has a

696 median KGES performance of 0.56, has much lower skewness and dispersion, and achieves
697 positive KGES values at a greater number (94%) of the catchments.

698 [77] So, while data-based representations may have a greater potential to learn from the data,
699 and thereby achieve greater predictive performance, the conceptual/lumped representations
700 contain valuable regularizing information that may help to prevent model performance from
701 becoming catastrophically poor under conditions where the data is insufficiently informative
702 about the dynamics of streamflow generation. We can speculate, therefore, that GR4J could
703 help to moderate the dispersion associated with the lower percentiles if an ensemble of these
704 three model types were to be used for operational streamflow prediction across Chile. Of
705 course, to implement such a system for Chile, further work would need to be done to
706 generalize the method for estimating parameter values to enable application at ungauged
707 catchment locations. We do not pursue this possibility further in this paper and leave it for
708 future work.

709 **4.2.3 Spatial Generalization**

710 [78] The results presented so far indicate that the data-based LSTM has the potential to
711 provide the “*best*” overall performance, while GR4J tends to provide more “*robust*” results
712 in cases where data-based approaches may fail. Meanwhile, the data-based RF is particularly
713 useful for enabling discovery, by providing clues that can lead to hypotheses about what
714 kinds of hydro-geo-climatic processes (and hence data sets) should be incorporated into
715 ongoing model development efforts.

716 [79] However, the previous analysis was for a “*pseudo-independent*” data set, consisting of
717 evaluation-period data from the same catchments that were used for model development. As
718 such, the results may not provide a reliable assessment of the quality of model performance
719 that might be expected at (other/new) catchments that are not part of the model development
720 dataset. **Figures 10a-c** and **Table 2** report the results of our “*out-of-sample*” analysis, where
721 model performance was assessed on the 167 CAMELS-CL catchments for which less than 7
722 years but more than 1 year of data were available (these catchments were withheld from the
723 model development dataset). Since GR4J parameter estimates are not available for these
724 catchments (an extra parameter regionalization step would be required, that was not pursued
725 in this study), this assessment was done only for LSTM and RF.

726 [80] Overall, the out-of-sample results indicate that LSTM and RF do not show significantly
727 different (relative to each other) spatial distributions of performance. This tends to conflict
728 with the in-sample evaluation results (**Figure 9**), despite the fact that both the in-sample and
729 out-of-sample catchment locations are distributed similarly with respect to the *Aridity Index*.
730 When we compare the CDF’s of in-sample and out-of-sample performance (**Figure 11a**) for
731 these models, we see that both RF and LSTM exhibit remarkably similar statistical
732 distributions of out-of-sample performance, which suggests that both of these ML-based
733 approaches have a similar ability to generalize to locations that were not included in the
734 model development dataset. There is, however, a larger deterioration in the statistical
735 distribution of model performance from in-sample to out-of-sample for LSTM than for RF.

736 [81] Meanwhile, the CDF of streamflow prediction bias (**Figure 11b**) shows that RF retains
737 the same tendencies in- and out-of-sample tendencies to overestimate the long-term mean
738 streamflow (compare with **Figure 7b**). This is encouraging, as it suggests the possibility of
739 being able to learn and correct for any long-term predictive bias at a given location.

740 [82] Finally, *Table 3* reports a more detailed statistical analysis of KGEss performance by
 741 model type, showing that RF slightly outperforms LSTM on most of the statistical indicators.
 742 So, while LSTM clearly achieved (in general) better temporal (in sample) generalization, the
 743 results for out-of-sample generalization are less definitive. It is possible that the tradeoff
 744 between temporal- and spatial-generalization ability is somehow different for each
 745 representational type. Further, this may be partially related to the differences in model
 746 development strategies – while the LSTM was sequentially fed with the information from
 747 different catchments (model parameters are updated using the data from each catchment in
 748 turn), the RF model development focuses on finding the best split for all catchments
 749 simultaneously, which may make it less sensitive to the new conditions encountered in out-
 750 of-sample testing. While this is simply speculative at this point, it would be interesting to
 751 further examine this issue using large-sample catchment-scale data sets from other parts of
 752 the world.

753 5 DISCUSSION

754 [83] An understanding of how hydrological processes vary at large (e.g. national) scales is
 755 important to the development of strategies for mitigating the effects of floods and droughts
 756 (and other natural hazards). Such understanding can be difficult to establish, given the large
 757 number of variables, attributes, and relationships that need to be considered. Under such
 758 circumstances, the traditional approach of attempting to model the entire diversity of hydro-
 759 geo-climatic conditions across an entire country/region with a single representational
 760 approach may not result in a sufficiently accurate characterization of the underlying *Data*
 761 *Generation Process (DGP)*. Through a case study, we have explored the possibility of using
 762 a multi-representational approach to address the challenge of large-scale model development,
 763 where the different representations are selected to have complementary strengths and with
 764 the goal of maximizing learning and discovery.

765 5.1 Challenges and Opportunities of a Multi-Representational Approach

766 [84] While each representation can support different kinds of discovery through the model
 767 development and evaluation process, adoption of a multi-representational approach brings
 768 forth both opportunities and challenges to be addressed.

- 769 1) It becomes difficult to implement a completely “*uniform*” strategy for model
 770 development since each representational approach may exploit the information in
 771 data differently and can have different requirements for inference.
- 772 2) For conceptual/process-resolved representations, discovery/learning about the spatial
 773 variability of hydrological processes is mediated through an analysis of spatial
 774 patterns or parameters.

775 [85] However, multiple parameter sets can give rise to similar model performance, thereby
 776 complicating our ability to make meaningful inferences. In contrast, for the ML-based
 777 models, the need for transformations and/or standardizations of the data was found to be
 778 unnecessary, and even to bring about declines in overall performance telling us how different
 779 representations are dealing with the data.

780 3) ML-based approaches facilitate an exploration of varying memory time scales.

781 [86] Our analysis suggested that, for energy-limited catchments in Chile, the ability to access
782 input information over the past 32 days was critical to achieving an optimal representation,
783 whereas for arid catchments the memory time-scale required was much longer (~270 days).
784 In this regard, 32 days is likely associated with rapid time-scale precipitation-driven
785 processes such as surface runoff and lateral flow, and while 270 days is likely associated with
786 slower time-scale groundwater driven processes such as baseflow. While further
787 investigation is needed to test these findings, such findings illustrate the power of ML-based
788 approaches to support learning and discovery.

789 4) The RF architecture enables an exploration of feature importance, potentially
790 enabling a higher degree of interpretability and discovery than the GR4J and LSTM
791 representations.

792 [87] In our case, aridity was seen to provide the highest-level segregation of catchments,
793 which makes sense given that the nature of the hydrological processes underlying the
794 generation of streamflow depends, unavoidably, on the availability of water. Beyond this,
795 various characteristic features associated with precipitation in the period just prior to the
796 streamflow event of interest were seen to provide strong explanatory power. By exploring
797 the structure of the decision tree model, it is possible to gain insight into the main
798 relationships or drivers governing the behavior of the system under investigation.

799 [88] Overall, by synthesizing the results obtained using a multi-representational approach,
800 we can obtain a more comprehensive overall picture of the underlying DGP, which in turn
801 creates a better context for a more in-depth investigation of the capabilities and performance
802 of each specific modeling approach.

803 **5.2 On the Issue of Data Informativeness**

804 [89] In terms of overall performance during the independent evaluation period (temporal
805 generalization), the LSTM model provided better overall (statistical) performance than GR4J
806 and RF. On the other hand, the GR4J model tended to be more robust, providing the best
807 performance for locations where KGE_{ss} was lower than 0.15. It might make sense, therefore,
808 to implement a lumped water balance model as a “*lower benchmark*” in any multi-
809 representational ensemble of component models, and to generally require that any ML-based
810 approach under consideration for inclusion in such an ensemble should demonstrate some
811 benefits over the benchmark. Further, when an ML-based model fails to perform well when
812 compared with the lumped water balance benchmark, this should alert us to the possibility
813 that the data may not be sufficiently informative regarding the processes we seek to model.

814 [90] In this regard, note also that RF performed slightly better than LSTM when tested out-
815 of-sample (spatial generalization). The reasons for this are not yet clear, but it is possible
816 that the LSTM model development strategy employed tends to overfit the temporal/sequence
817 patterns in the data. Regardless, this result also lends support to the idea that a multi-
818 representational ensemble has the potential to be superior to one that is less representationally
819 diverse.

820 [91] One common finding for all three models was their poor performance in one particular
 821 region of Chile. Given that the most important shared commonality of the three models is
 822 their access to the same dataset, coupled with the fact that ML-based approaches are highly
 823 flexible, this result strongly indicates that the dataset is not sufficiently informative to enable
 824 a suitable characterization of the streamflow response of this region, and that the main driver
 825 of local streamflow is not precipitation. Unfortunately, the Chilean CAMELS dataset does
 826 not include attributes from which it could be possible to infer groundwater-driven baseflow,
 827 or other related processes, and so proper characterization of the streamflow response of this
 828 region will require further investigation and exploration of alternative sources of relevant
 829 information.

830 [92] Overall, these observations point to the issue of whether the available data is
 831 informative enough for a sufficiently robust characterization of the underlying DGP to be
 832 achieved. While a multi-representation approach cannot (by itself) solve that issue, it can
 833 certainly help us to recognize the existence of the problem so that we can seek additional
 834 relevant information that may help us in the process of learning and discovery.

835 **5.3 Relationship of the Multi-Representational Approach to Hypothesis Testing**

836 [93] Given the tendency for each of the three representation types to provide better
 837 performance under different hydro-geo-climatic conditions, and the fact that each one
 838 facilitates different (complementary) kinds of information extraction and degree of
 839 interpretability, it seems clear that the three models can collectively be treated as a valuable
 840 tool for gaining insights regarding the underlying *DGP*. From this point of view, a
 841 meaningful answer to the question “*Does a single “correct” catchment-scale hydrological*
 842 *model exist at all?*” expressed by *Clark et al (2011)* may be that:

843 *“It seems sensible to abandon any concept of a “best” model, and instead consider the*
 844 *value of learning to live with a plurality of representations while developing strategies*
 845 *for extracting important relevant information from the representational ensemble”.*

846 This is, of course, because any model is unavoidably a “*simplified*” (and hopefully
 847 informationally relevant) representation of reality.

848 [94] Another way of think of this is that it is the “*ensemble of representations*” (and not each
 849 of the individual components thereof) that is actually “*the model*” per se, since it helps to
 850 meet the goal of incorporating within the “*Model*” (writ large) a representation of “*what we*
 851 *know that we do not know*” (i.e., our known uncertainties). From this perspective, our task
 852 is to populate this ensemble with representations that best support our investigative goals.
 853 This is clearly consistent with the idea of a “*multiple-hypothesis approach*” (*Clark et al,*
 854 *2011*), but one where the hypotheses are selected to be as (potentially) *informationally*
 855 *complementary* as possible, so that learning/discovery can be maximized. In contrast, an
 856 approach where the ensemble consists of hypotheses that may only be marginally different
 857 from each other (e.g, that all share the same system architecture while differing only in the
 858 forms of the process parameterization equations) may not lend itself to efficient and effective
 859 learning (*Gharari et al., 2021*).

860 [95] Such a perspective unavoidably affects how we think about the model development
861 process, and its role in a scientific investigation. Our view is that conceptual/process/theory-
862 based and ML-data-based approaches to model development must co-exist within such an
863 environment, with neither being the dominant approach, and that a multi-representational
864 strategy is a key to promoting model-based scientific discovery. While this perspective is
865 likely to promote (as is currently happening) interest in hybrid approaches that integrate
866 theory-based and data-based strengths, it is not clear that such a push towards reductionism
867 through integration will necessarily obviate the need for a continued multi-representational
868 approach in order for models to be tools that enable scientific discovery.

869 6 CONCLUSIONS

870 [96] In conclusion, while the metaphor of the “*blind people and the elephant*” is highly
871 suggestive, it is not completely accurate. In the metaphor, each person constructs a different
872 representation based on potentially different prior knowledge and clearly different sensory
873 information (data). In our case, all of the representational approaches have access to the same
874 sources of information (dataset) but differ in their abilities to fully exploit that information
875 due to (prior) representational restrictions.

876 [97] So, while one might debate how to improve the metaphorical story to match the current
877 situation, more important is the fact that an optimal strategy for scientific discovery would
878 seem to be one that combines multiple complementary model structural representations
879 (modeling strategies) with multiple complementary mechanisms for extracting information
880 from data (inferential strategies). In this regard, it is perhaps worth noting that the strategy of
881 “*multi-headed attention*” that has recently become the topic of intense inquiry in fields such
882 as text prediction, translation, and speech recognition (*Vaswani et al, 2017; Devlin et al,*
883 *2018; Luo et al, 2021*), is explicitly based on the notion that multiple attentional perspectives
884 bring considerable value to such tasks.

885 [98] This paper seeks to explicitly promote adoption of a multi-representational approach to
886 learning, understanding, and discovery in the hydrological sciences. We believe that the
887 multi-representational approach is fundamental to understanding hydrology at a large scale,
888 where the complexity of the system we seek to understand and represent demands access to
889 large and informationally diverse data sets and an analytical strategy that is purposefully
890 diverse. As always, we are keenly interested in dialogue and collaboration on this and related
891 issues of how we use models to support prediction, understanding, and scientific discovery.

892 ACKNOWLEDGMENTS

893 [99] This publication is the product of research done by *De la Fuente (2021)* to satisfy the
894 requirements for obtaining a Master of Science degree in Hydrology, while being funded by
895 the Chilean Government scholarship “*Beca de Magister en el Extranjero, Becas Chile en*
896 *Áreas Prioritarias, Convocatoria 2018*”. Gupta acknowledges partial support from the
897 Australian Research Council (ARC) through the Centre of Excellence for Climate Extremes
898 grant CE170100023. Condon acknowledges partial support from NSF Early Career Award
899 grant 1945195. The authors declare no conflicts of interest.

900 **CODE AND DATA AVAILABILITY**

901 [100] The CAMELS-CL dataset is freely available from
902 <https://doi.pangaea.de/10.1594/PANGAEA.894885>. The analytical methods are presented
903 as a Jupiter notebook freely available at
904 <http://www.hydroshare.org/resource/fc08997100fa4cd6abdd8a4f5731de15>.

905 **APPENDIX**

906 **A1. Model Development Strategies**

907 ***A1.1 Partitioning the Data***

908 [101] A key step in model development is to partition the available data into ‘*model*
909 *development*’ and ‘*evaluation*’ subsets, where the former is used for model structure selection
910 and parameter tuning, while the latter is used to assess the generalization performance that
911 can be expected from the developed model. However, no clear guidance exists for how to
912 achieve such a partitioning for data that represent dynamical hydrological systems (*Wu et al,*
913 *2013, Daggupati et al, 2015, Zheng et al 2018, Guo et al 2020*). In general, the hydrological
914 literature has traditionally assumed that the entire available dataset comes from a stationary
915 underlying data generating process, and that any split that preserves the full range of
916 hydrologic variability (dry, medium, and wet) in both sets is satisfactory. Based on this
917 assumption, it is common to use a continuous-time period that makes up ~60-80% of the
918 available data for model development, while allocating the remaining ~20-40% for an
919 evaluation of the generalization ability of the model.

920 [102] In this study, we adopt the strategy of further partitioning the ‘*model development*’
921 subset into ‘*calibration*’ and ‘*selection*’ subsets, where the calibration subset is used for
922 model/network parameter tuning (commonly called ‘*training*’ in the ML literature), and the
923 selection subset is used for model/network structure selection and/or hyperparameter tuning
924 (commonly called ‘*validation*’ in the ML literature). Note that we adopt this naming
925 convention to try and overcome the inconsistency in terminology between the ML and
926 hydrological modeling literature. Accordingly, the available data are partitioned into three
927 subsets, where the first 60% of the data is used for model calibration, the next 24% is used
928 for model selection, and the final 16% is used for model evaluation (commonly called
929 ‘*testing*’ in the ML literature)

930 ***A1.2 Variable Selection***

931 [103] The variables selected from the CAMEL-CL dataset include two sources of
932 precipitation (CR2MET and MSWEP, both having long records), three values characterizing
933 temperature (Maximum, mean, and minimum), and potential evapotranspiration (PET)
934 estimated via the *Hargreaves and Samani (1985)* method. The PET value derived from
935 MODIS was not used because its time step is higher than daily (8 days). Further, the snow
936 water equivalent (SWE) data does not cover the entire country and was therefore not
937 considered suitable for the current study.

938 [104] Because the GR4J model has a pre-defined input representation, it is unable to use any
939 other sources of data and so we used the weighted average of the two sources of precipitation
940 as input to the GR4J model. In contrast, the ML-based models are able to use the information
941 provided by all of the available variables and attributes, but in different ways. While the RF
942 model used lagged input variables as surrogates for system memory (lag memory), the LSTM
943 model used internal state variables to characterize system memory (sequential memory).
944 More details regarding the variables and attributes used for the development of each model
945 type are presented in Tables A-1, A-2, and A-3.

946 ***A1.3 Representing System Memory***

947 [105] For the RF representation, which does not explicitly include dynamical state variables,
948 system memory was included by concatenating past inputs (precipitation, evapotranspiration,
949 and temperature) to the inputs for the current time step. This follows the idea of a Markov
950 Process, where a state variable can be thought of as a summary property of an infinite number
951 of past inputs to the system. For the RF, the number of past input lags was treated as a model
952 hyperparameter. While this strategy enables important information to be made available to
953 the model, it results in a very high cost (in terms of computational and storage resources)
954 because huge system memory is required to manage the dataset as the number of lags is
955 increased. We found that at 32 days of lagged memory, the computation became unstable.
956 This prevented us from readily exploring longer memory time scales, such as 270 or 365 days
957 (or longer), and the results presented only consider a memory time-scale of 16 days.

958 ***A1.4 Model Warm-Up***

959 [106] It is recommended, regardless of representational strategy, to use a warm-up period
960 (during which performance metrics are not computed) to minimize errors associated with the
961 initialization of dynamical model states. For lumped water balance modeling it is common
962 to use a full year (365 days) of data for this purpose; for example, *Perrin et al. (2003)* used
963 a full year to initialize the GR4J model, following the suggestion of *Chiew and McMahon*
964 *(1994)*. For the LSTM machine-learning approach, *Kratzert et al. (2019)* used 270 days, after
965 testing 90, 180, 270, and 365 days as different options.

966 [107] In this study, we adopted the following strategy for warm-up period selection. For the
967 GR4J and LSTM representations, because preliminary testing suggested that the LSTM
968 requires the longest warm-up period to ensure stable results, we followed the strategy of first
969 tuning the LSTM to determine a suitable warm-up period length (as a model hyperparameter)
970 and then using that same period to “warm-up” the GR4J model.

971 ***A1.5 Parameters and Hyperparameters to be Tuned***

972 [108] Each representational strategy involves different sets of parameters and
973 hyperparameters, depending on its structural form. Whereas the original GR4J model
974 contains 4 tunable parameters that must be calibrated for each catchment, our implementation
975 includes an additional 3 parameters, two of which are used to facilitate driving the model by
976 a weighted average of the two available precipitation products (CR2MET and MSWEP), and
977 the third being the Box-Cox transformation parameter.

978 [109] For the LSTM model, in addition to a large number of system-wide network weights
 979 and biases, 6 hyperparameters must be tuned, namely the sequence length (memory from the
 980 past hidden states), number of hidden nodes, batch size, number of epochs, standardization
 981 parameters, and the Box-Cox transformation parameter. To standardize the data (centering
 982 by subtracting the mean, and rescaling by dividing by the standard deviation), we investigated
 983 two options – global and local standardization. In global standardization, for each variable,
 984 we use the mean and standard deviation computed from the entire dataset, whereas in local
 985 standardization (which is applied only to the precipitation and streamflow variables) we
 986 assume that the local means and standard deviations vary as functions of the aridity index.

987 [110] Finally, in addition to determining the nodal split “parameters”, the RF model requires
 988 the tuning of 4 hyperparameters for the entire set of catchments taken together, the first
 989 representing system “memory” (expressed as the number of days previous to the current day
 990 for which inputs are simultaneously presented to the model), the second being the Box-Cox
 991 transformation parameter, the third being the number of trees, and the fourth being the
 992 minimum number of elements that must be retained in the last leaf. To attempt to circumvent
 993 the problem of being unable to input lagged daily inputs beyond 32 days to the RF model to
 994 account for “memory” in the system, we augment the input data to include surrogate variables
 995 intended to be informative about the state of the system. Specifically, we included the
 996 “month-of-year” as an attribute, to enable the model to learn a representation of long-term
 997 memory as the average behavior associated with different months of the year. Meanwhile, as
 998 mentioned above, the short-term memory was treated as a hyperparameter.

999 *A1.6 Model Calibration/Training*

1000 [111] To calibrate the parameters of the GR4J model to each catchment in the calibration
 1001 period, we tested both the *Root Mean Square Error* (RMSE) and the *Kling-Gupta Efficiency*
 1002 (KGE, [Gupta et al., 2009](#)), as defined below. Overall, we found that KGE provided slightly
 1003 more robust results ([De la Fuente, 2021](#)), and therefore we present here only the results
 1004 obtained using KGE in this paper.

$$1005 \quad RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

$$1006 \quad KGE = 1 - \sqrt{(r - 1)^2 + (\alpha - 1)^2 + (\beta - 1)^2}$$

1007 y_i : Measured streamflow

1008 \hat{y}_i : Simulated streamflow

1009 n : Total number of data

1010 r : Linear correlation coefficient between y_i and \hat{y}_i

1011 α : σ_s/σ_o : relative variability between simulated and observed data.

1012 β : μ_s/μ_o : ratio between simulated and observed data.

1013 [112] For parameter optimization, we used three algorithms from the *Spotpy Python* library
 1014 ([Houska et al., 2015](#)), namely *Maximum Likelihood Estimation* (MLE), *Differential*
 1015 *Evolution Adaptive Metropolis* (DE-MCz), and *Shuffled Complex Evolution* (SCE-UA). In

1016 total, 22 independent optimization runs were done for each catchment, and the parameter set
 1017 that provided the best performance (out of the 22 parameter sets so obtained) on the ‘*selection*
 1018 (*hyperparameter tuning*)’ data subset was chosen.

1019 [113] To develop the RF model, we use the *Scikit-learn Python library* (Pedregosa et al.,
 1020 2011). The *RandomForestRegressor* module (version 0.23.1) has two options for
 1021 performance metrics – *Mean Squared Error* (MSE) and *Mean Absolute Error* (MAE). While
 1022 MAE can be used to reduce the tendency to emphasize larger streamflow values, because we
 1023 are implementing the Box-Cox transformation on streamflow we chose MSE to be the metric
 1024 used for RF calibration.

$$1025 \quad MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

$$1026 \quad MAE = \frac{\sum_{i=1}^n abs(y_i - \hat{y}_i)}{n}$$

1027 [114] To train the LSTM model, we used the implementation provided by Kratzert et al
 1028 (2019) and modified it to conform to the data structures and variables of the CAMELS-CL
 1029 dataset. Whereas the original code enables the choice of either MSE or NSE as the calibration
 1030 metric, we used only NSE because its normalization of the error enables better comparison
 1031 across catchments having different amounts of temporal variability.

$$1032 \quad NSE = 1 - \frac{MSE}{\sigma_o^2}$$

1033 ***A1.7 Model Performance Evaluation***

1034 [115] For performance evaluation, we use the KGE skill score (KGEss) (Knoben et al., 2019)
 1035 computed on the evaluation-period data. KGEss is a rescaled version of the KGE metric such
 1036 that a value of zero corresponds to the prediction being no better than simply using the mean
 1037 observed streamflow, in a manner analogous to NSE. While other metrics, including NSE
 1038 and RMSE, were also used for model evaluation (De la Fuente, 2021), we do not report them
 1039 here as the conclusions are similar to those obtained using KGEss. Importantly, we account
 1040 for sampling variability by computing the estimated posterior distributions of KGEss by
 1041 bootstrapping 100 times (Efron and Tibshirani, 1994) and using the median value of KGEss
 1042 in all comparisons.

$$1043 \quad KGEss = \frac{KGE - KGE_{benchmark}}{1 - KGE_{benchmark}} = \frac{KGE + \sqrt{2} - 1}{\sqrt{2}} = 1 - \frac{1 - KGE}{\sqrt{2}}$$

1044 ***A1.8 Out-of-Sample Testing***

1045 [116] For an additional ‘*out-of-sample*’ model evaluation step, we retained all of the
 1046 CAMELS-CL catchments for which less than 7 years but more than 1 year of data is
 1047 available; being less than 7 years of record length, none of these catchments are included in
 1048 the model development data set. The resulting 167 catchments facilitate a meaningful out-

1049 of-sample operational comparison of the generalization abilities of the LSTM and RF ML-
1050 based representations. Note that the GR4J model was not tested using this out-of-sample set
1051 of catchments since regional generalization of lumped water balance model parameters to
1052 ‘*ungauged*’ catchments is not within the scope of this paper.

1053 [117] Note that, because the additional ‘*out-of-sample*’ model evaluation data set is
1054 completely independent of the data set used for model development, while being similarly
1055 representative of the geo-hydro-climatic variability across the country (*Figure A-1*), the
1056 model performed using those data can be considered to be similar to the idea of “*Proxy-basin*
1057 *differential split-sample testing*” (*Klemeš, 1986*).

REFERENCES

- Addor, N., Newman, A. J., Mizukami, N., & Clark, M. P. (2017). The CAMELS data set: catchment attributes and meteorology for large-sample studies. *Hydrology and Earth System Sciences*, *21*(10), 5293-5313. <https://doi.org/10.5194/hess-21-5293-2017>
- Addor, N., Nearing, G., Prieto, C., Newman, A. J., Le Vine, N., & Clark, M. P. (2018). A ranking of hydrological signatures based on their predictability in space. *Water Resources Research*, *54*(11), 8792-8812. <https://doi.org/10.1029/2018WR022606>
- Alvarez-Garreton, C., Mendoza, P. A., Boisier, J. P., Addor, N., Galleguillos, M., Zambrano-Bigiarini, M., ... & Ayala, A. (2018). The CAMELS-CL dataset: catchment attributes and meteorology for large sample studies-Chile dataset. *Hydrology and Earth System Sciences*, *22*(11), 5817-5846. <https://doi.org/10.5194/hess-22-5817-2018>
- Breiman, L. (2001). Random forests. *Machine learning*, *45*(1), 5-32. <https://doi.org/10.1023/A:1010933404324>
- Booij, M. J., Schipper, T. C., & Marhaento, H. (2019). Attributing changes in streamflow to land use and climate change for 472 catchments in Australia and the United States. *Water*, *11*(5), 1059. <https://doi.org/10.3390/w11051059>
- Box, G. E., & Cox, D. R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society: Series B (Methodological)*, *26*(2), 211-243. <https://doi.org/10.1111/j.2517-6161.1964.tb00553.x>
- Chen, S. A., Michaelides, K., Grieve, S. W., & Singer, M. B. (2019). Aridity is expressed in river topography globally. *Nature*, *573*(7775), 573-577. <https://doi.org/10.1038/s41586-019-1558-8>
- Chiew, F., & McMahon, T. (1994). Application of the daily rainfall-runoff model MODHYDROLOG to 28 Australian catchments. *Journal of Hydrology*, *153*(1-4), 383-416. [https://doi.org/10.1016/0022-1694\(94\)90200-3](https://doi.org/10.1016/0022-1694(94)90200-3)
- Clark, M. P., Kavetski, D., & Fenicia, F. (2011). Pursuing the method of multiple working hypotheses for hydrological modeling. *Water Resources Research*, *47*(9). <https://doi.org/10.1029/2010WR009827>
- Clark, M. P., Nijssen, B., Lundquist, J. D., Kavetski, D., Rupp, D. E., Woods, R. A., ... & Rasmussen, R. M. (2015). A unified approach for process-based hydrologic modeling: 1. Modeling concept. *Water Resources Research*, *51*(4), 2498-2514. <https://doi.org/10.1002/2015WR017198>
- Craig, J. R., Brown, G., Chlumsky, R., Jenkinson, R. W., Jost, G., Lee, K., ... & Tolson, B. A. (2020). Flexible watershed simulation with the Raven hydrological modelling framework. *Environmental Modelling & Software*, *129*, 104728. <https://doi.org/10.1016/j.envsoft.2020.104728>

Daggupati, P., Pai, N., Ale, S., Douglas-Mankin, K. R., Zeckoski, R. W., Jeong, J., ... & Youssef, M. A. (2015). A recommended calibration and validation strategy for hydrologic and water quality models. *Transactions of the ASABE*, 58(6), 1705-1719. <https://doi.org/10.13031/trans.58.10712>

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint <https://arxiv.org/abs/1810.04805>

De la Fuente, L. (2021). *Using Big-Data to Develop Catchment-Scale Hydrological Models for Chile* (Master dissertation, The University of Arizona). <https://repository.arizona.edu/handle/10150/656824>

Efron, B., & Tibshirani, R. J. (1994). *An introduction to the bootstrap*. CRC press.

Fenicia, F., Kavetski, D., & Savenije, H. H. (2011). Elements of a flexible approach for conceptual hydrological modeling: 1. Motivation and theoretical development. *Water Resources Research*, 47(11). <https://doi.org/10.1029/2010WR010174>

Gharari, S., Gupta, H. V., Clark, M. P., Hrachowitz, M., Fenicia, F., Matgen, P., & Savenije, H. H. (2021). Understanding the Information Content in the Hierarchy of Model Development Decisions: Learning from data. *Water Resources Research*, <https://doi.org/10.1029/2020WR027948>

Guo, D., Zheng, F., Gupta, H., & Maier, H. R. (2020). On the Robustness of Conceptual Rainfall-Runoff Models to Calibration and Evaluation Data Set Splits Selection: A Large Sample Investigation. *Water Resources Research*, 56(3), e2019WR026752. <https://doi.org/10.1029/2019WR026752>

Gupta, V. K., & Sorooshian, S. (1985). The relationship between data and the precision of parameter estimates of hydrologic models. *Journal of Hydrology*, 81(1-2), 57-77. [https://doi.org/10.1016/0022-1694\(85\)90167-2](https://doi.org/10.1016/0022-1694(85)90167-2)

Gupta, H. V., Wagener, T., & Liu, Y. (2008). Reconciling theory with observations: elements of a diagnostic approach to model evaluation. *Hydrological Processes: An International Journal*, 22(18), 3802-3813. <https://doi.org/10.1002/hyp.6989>.

Gupta HV, H Kling, KK Yilmaz & GF Martinez (2009). Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling. *Journal of hydrology*, 377(1-2), 80-91. <https://doi.org/10.1016/j.jhydrol.2009.08.003>

Gupta, H. V., Clark, M. P., Vrugt, J. A., Abramowitz, G., & Ye, M. (2012). Towards a comprehensive assessment of model structural adequacy. *Water Resources Research*, 48(8). <https://doi.org/10.1029/2011WR011044>

Hargreaves, G. H., & Samani, Z. A. (1985). Reference crop evapotranspiration from temperature. *Applied engineering in agriculture*, 1(2), 96-99. <https://doi.org/10.13031/2013.26773>

Hassan, M., & Hassan, I. (2021). Improving Artificial Neural Network Based Streamflow Forecasting Models through Data Preprocessing. *KSCE Journal of Civil Engineering*, 1-13. <https://doi.org/10.1007/s12205-021-1859-y>

Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>

Hoedt, P. J., Kratzert, F., Klotz, D., Halmich, C., Holzleitner, M., Nearing, G., ... & Klambauer, G. (2021). MC-LSTM: Mass-Conserving LSTM. *arXiv preprint arXiv:2101.05186*. <https://arxiv.org/abs/2101.05186v3>

Houska, T., Kraft, P., Chamorro-Chavez, A., & Breuer, L. (2015). SPOTting model parameters using a ready-made python package. *PloS one*, 10(12). <https://doi.org/10.1371/journal.pone.0145180>

Hu, C., Wu, Q., Li, H., Jian, S., Li, N., & Lou, Z. (2018). Deep learning with a long short-term memory networks approach for rainfall-runoff simulation. *Water*, 10(11), 1543. <https://doi.org/10.3390/w10111543>

Klemeš, V. (1986). Operational testing of hydrological simulation models. *Hydrological sciences journal*, 31(1), 13-24. <https://doi.org/10.1080/02626668609491024>

Knoben, W. J., Freer, J. E., & Woods, R. A. (2019). Inherent benchmark or not? Comparing Nash–Sutcliffe and Kling–Gupta efficiency scores. *Hydrology and Earth System Sciences*, 23(10), 4323-4331. <https://doi.org/10.5194/hess-23-4323-2019>

Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall–runoff modelling using long short-term memory (LSTM) networks. *Hydrology and Earth System Sciences*, 22(11), 6005-6022. <https://doi.org/10.5194/hess-22-6005-2018>

Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019). Benchmarking a catchment-aware Long Short-Term Memory Network (LSTM) for large-scale hydrological modeling. *arXiv preprint arXiv:1907.08456*. <https://doi.org/10.5194/hess-2019-368>

Kunnath-Poovakka, A., & Eldho, T. I. (2019). A comparative study of conceptual rainfall-runoff models GR4J, AWBM and Sacramento at catchments in the upper Godavari river basin, India. *Journal of Earth System Science*, 128(2), 33. <https://doi.org/10.1007/s12040-018-1055-8>

Le Moine, N. (2008). *Le bassin versant de surface vu par le souterrain: une voie d'amélioration des performances et du réalisme des modèles pluie-débit?* (Doctoral dissertation, Doctorat Géosciences et Ressources Naturelles, Université Pierre et Marie Curie Paris VI).

Luo, H., Zhang, S., Lei, M., & Xie, L. (2021, January). Simplified self-attention for transformer-based end-to-end speech recognition. In *2021 IEEE Spoken Language*

Technology Workshop (SLT) (pp. 75-81). IEEE.
<https://doi.org/10.1109/SLT48900.2021.9383581>

Malone, R. W., Yagow, G., Baffaut, C., Gitau, M. W., Qi, Z., Amatya, D. M., ... & Green, T. R. (2015). Parameterization guidelines and considerations for hydrologic models. *Transactions of the ASABE*, 58(6), 1681-1703. <https://doi.org/10.13031/trans.58.10709>

Nearing, G. S., Ruddell, B. L., Bennett, A. R., Prieto, C., & Gupta, H. V. (2020). Does information theory provide a new paradigm for earth science? Hypothesis testing. *Water Resources Research*, 56(2). <https://doi.org/10.1029/2019WR024918>

Pagano, T., Hapuarachchi, P., & Wang, Q. J. (2010). Continuous rainfall-runoff model comparison and short-term daily streamflow forecast skill evaluation. *CSIRO*; 2010. <https://doi.org/10.4225/08/58542c672dd2c>

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. *The Journal of machine Learning research*, 12, 2825-2830. <https://jmlr.csail.mit.edu/papers/v12/pedregosa11a.html>

Perrin, C., Michel, C., & Andréassian, V. (2003). Improvement of a parsimonious model for streamflow simulation. *Journal of hydrology*, 279(1-4), 275-289. [https://doi.org/10.1016/S0022-1694\(03\)00225-7](https://doi.org/10.1016/S0022-1694(03)00225-7)

Pushpalatha, R., Perrin, C., Le Moine, N., Mathevet, T., & Andréassian, V. (2011). A downward structural sensitivity analysis of hydrological models to improve low-flow simulation. *Journal of hydrology*, 411(1-2), 66-76. <https://doi.org/10.1016/j.jhydrol.2011.09.034>

Sezen, C., & Partal, T. (2019). The utilization of a GR4J model and wavelet-based artificial neural network for rainfall-runoff modelling. *Water Supply*, 19(5), 1295-1304. <https://doi.org/10.2166/ws.2018.189>

Sudriani, Y., Ridwansyah, I., & Rustini, H. A. (2019, July). Long short term memory (LSTM) recurrent neural network (RNN) for discharge level prediction and forecast in Cimandiri river, Indonesia. In *IOP Conference Series: Earth and Environmental Science* (Vol. 299, No. 1, p. 012037). IOP Publishing. <https://doi.org/10.1088/1755-315/299/1/012037>

Valéry, A. (2010). *Modélisation précipitations débit sous influence nivale: Elaboration d'un module neige et évaluation sur 380 bassins versants* (Doctoral dissertation, Doctorat Hydrobiologie, Institut des Sciences et Industries du Vivant et de l'Environnement AgroParisTech).

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008). <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>

Vrugt, J. A., Gupta, H. V., Dekker, S. C., Sorooshian, S., Wagener, T., & Bouten, W. (2006). Application of stochastic parameter optimization to the Sacramento Soil Moisture Accounting model. *Journal of Hydrology*, 325(1-4), 288-307. <https://doi.org/10.1016/j.jhydrol.2005.10.041>

Wu, W., May, R. J., Maier, H. R., & Dandy, G. C. (2013). A benchmarking approach for comparing data splitting methods for modeling water resources parameters using artificial neural networks. *Water Resources Research*, 49(11), 7598-7614. <https://doi.org/10.1002/2012WR012713>

Zhang, J., Zhu, Y., Zhang, X., Ye, M., & Yang, J. (2018). Developing a Long Short-Term Memory (LSTM) based model for predicting water table depth in agricultural areas. *Journal of hydrology*, 561, 918-929. <https://doi.org/10.1016/j.jhydrol.2018.04.065>

Zheng, F., Maier, H. R., Wu, W., Dandy, G. C., Gupta, H. V., & Zhang, T. (2018). On lack of robustness in hydrological model development due to absence of guidelines for selecting calibration and evaluation data: Demonstration for data-driven models. *Water Resources Research*, 54(2), 1013-1030. <https://doi.org/10.1002/2017WR021470>

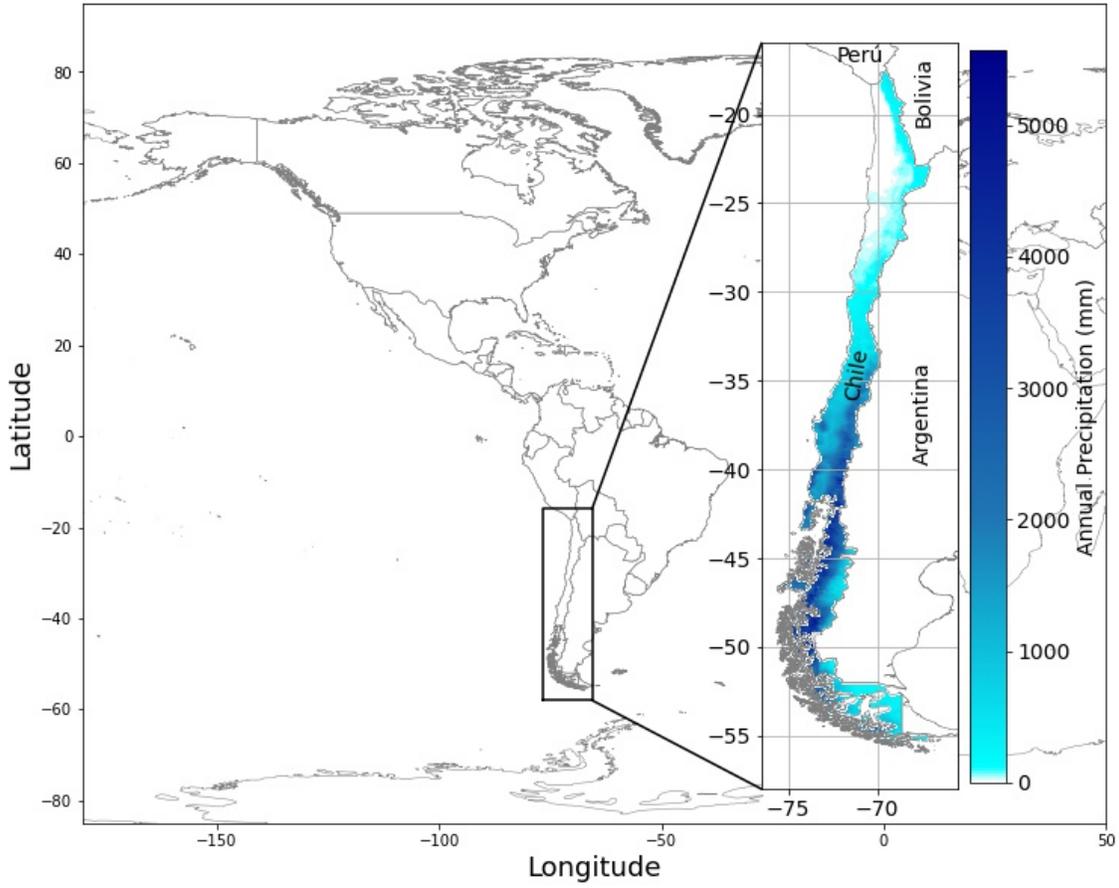


Figure 1. Map showing the geographic location of Chile, and its spatial distribution of annual precipitation.

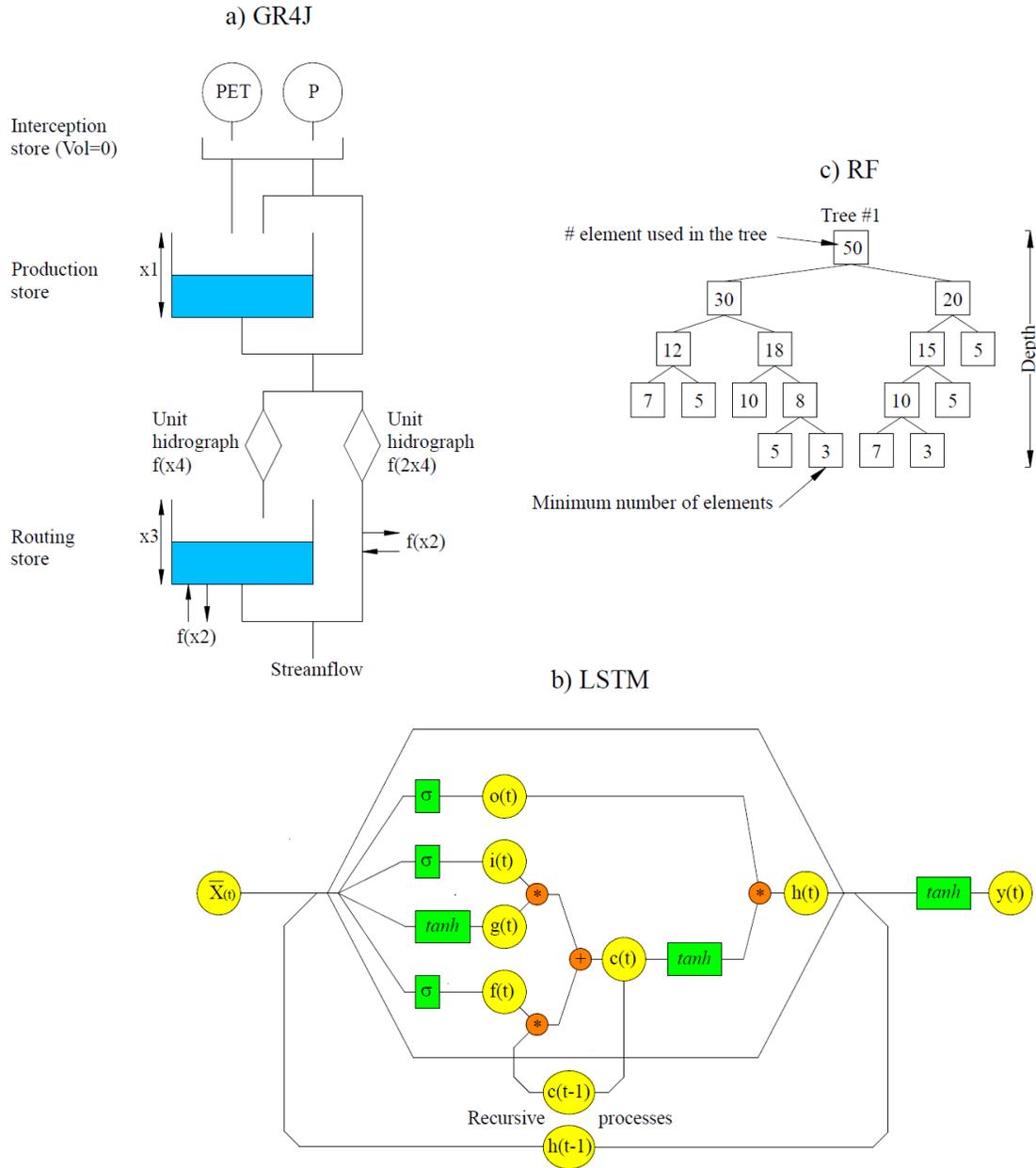


Figure 2. Representational structures of the three different models used; a) The GR4J lumped-water balance model; b) The *Long-Short Term Memory* (LSTM) machine learning model; and c) The *Random Forest* (RF) machine learning model.

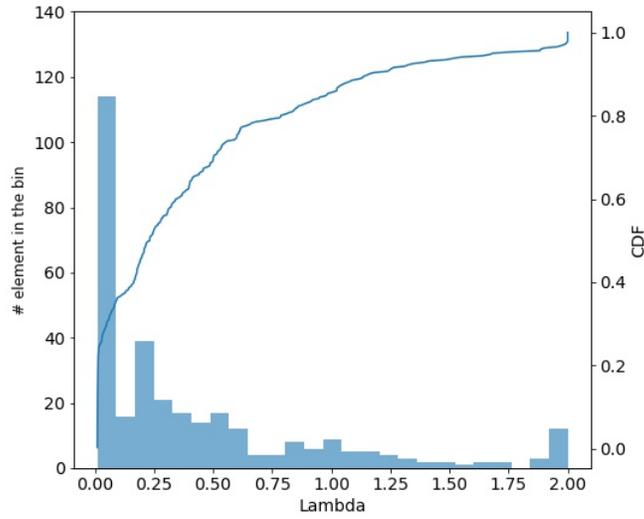


Figure 3. Frequency distribution of the λ hyperparameter for the 322 catchments when using the GR4J model.

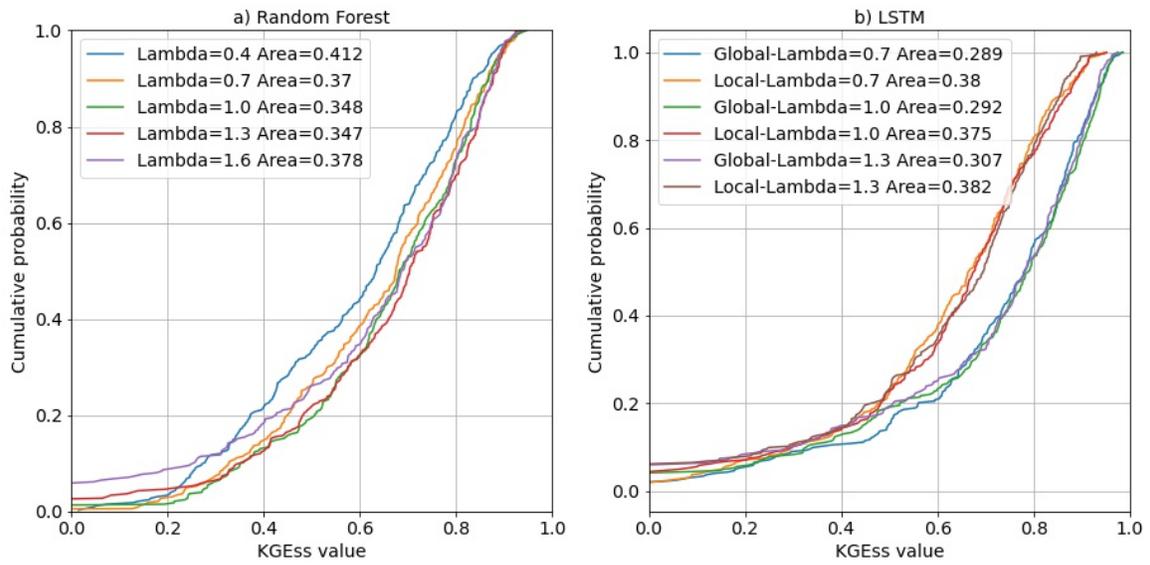


Figure 4. Selection period performance CDF's obtained using different values of the λ hyperparameter; the Left subplot is for RF and the right subplot is for LSTM.

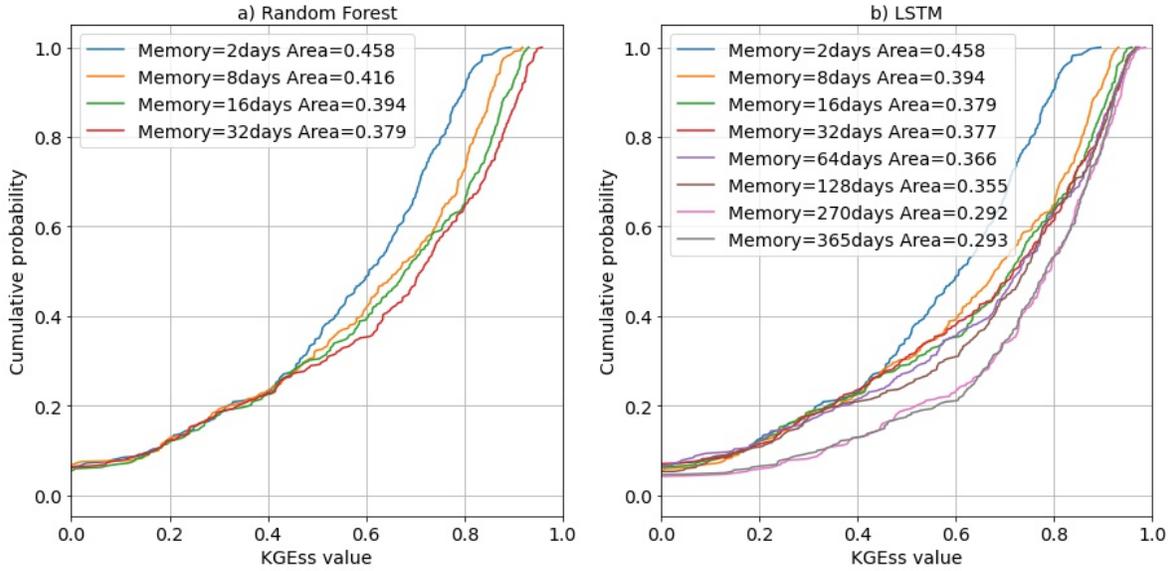


Figure 5. Selection period performance CDF's for the RF and LSTM models, showing dependence on memory lag.

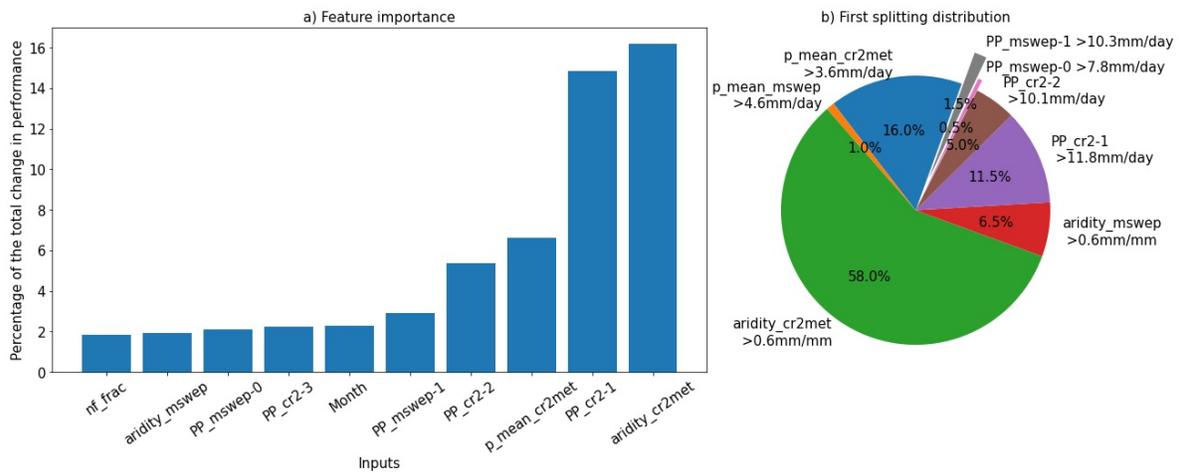


Figure 6. Feature importance and distribution of the first split of the RF model (For a description of the name of each attribute or variable, see Table A-3).

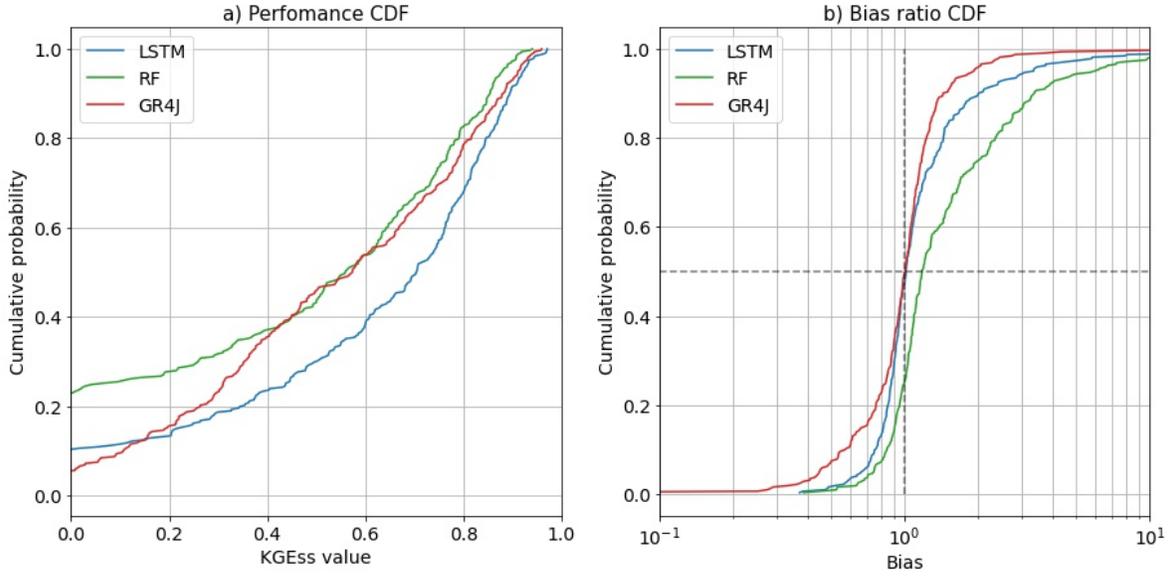


Figure 7. Evaluation period performance CDFs for the three models. The left subplot shows KGEss and the right subplot shows Bias Ratio.

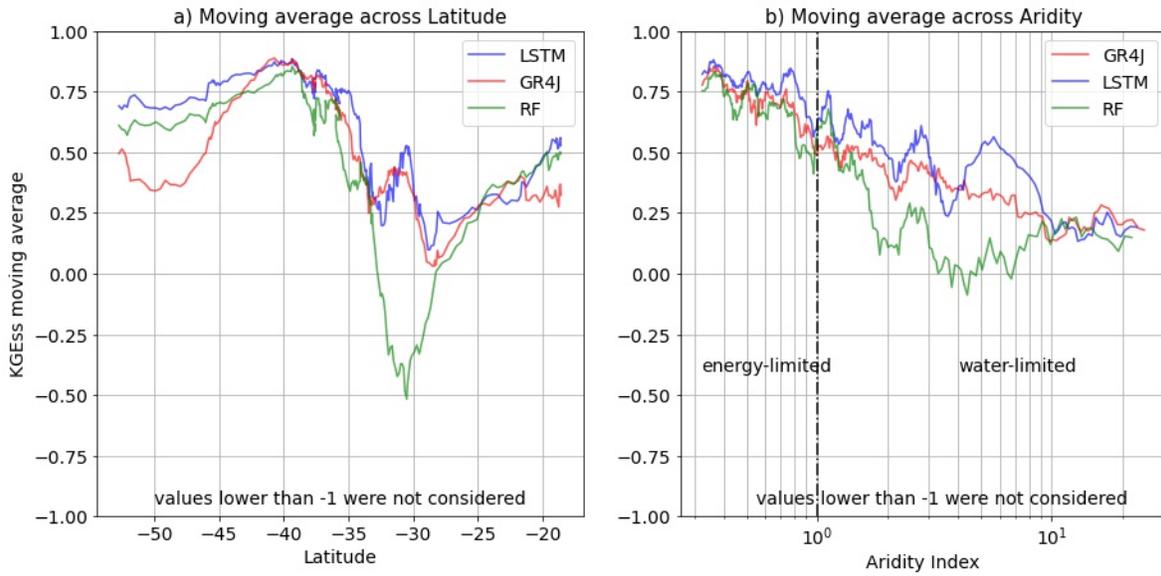


Figure 8. Variation of evaluation period KGEss performance with aridity and latitude.

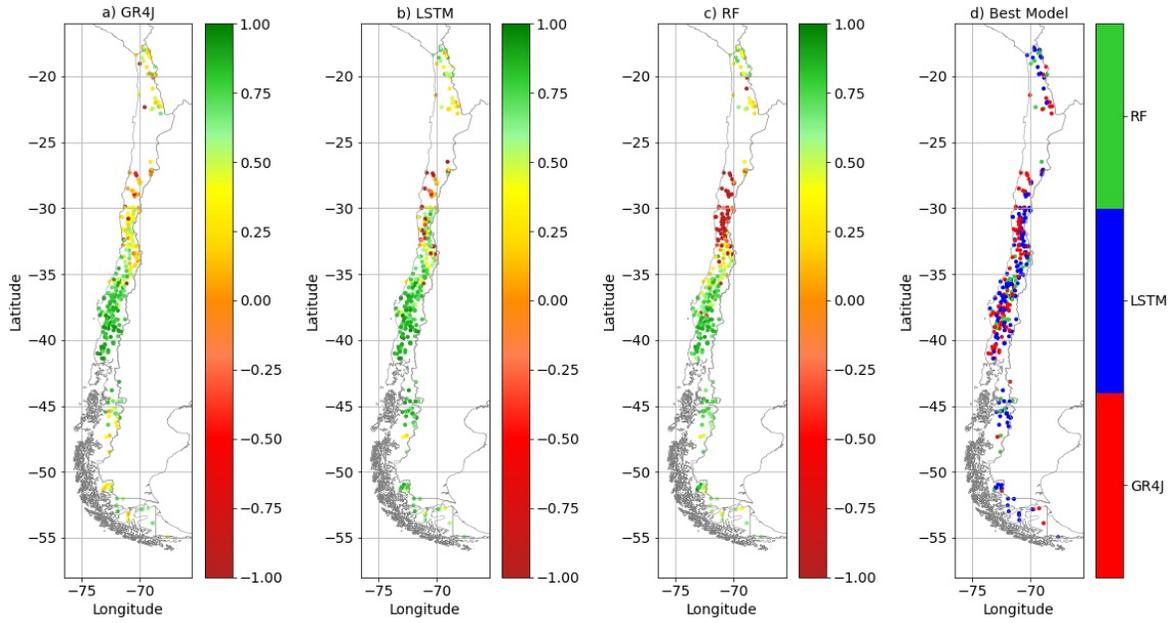


Figure 9. Spatial distributions of evaluation period model performance for a) GR4J, b) LSTM, c) RF, and d) the “best” performing model.

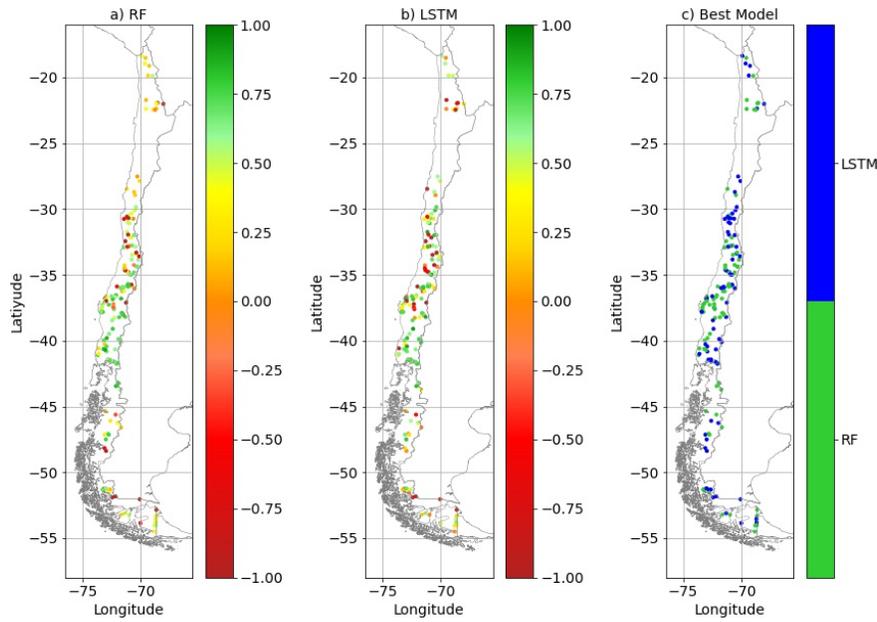


Figure 10. Spatial distributions “out-of-sample” performance for a) RF and b) LSTM. The right subplot (c) indicates the “best” performing model.

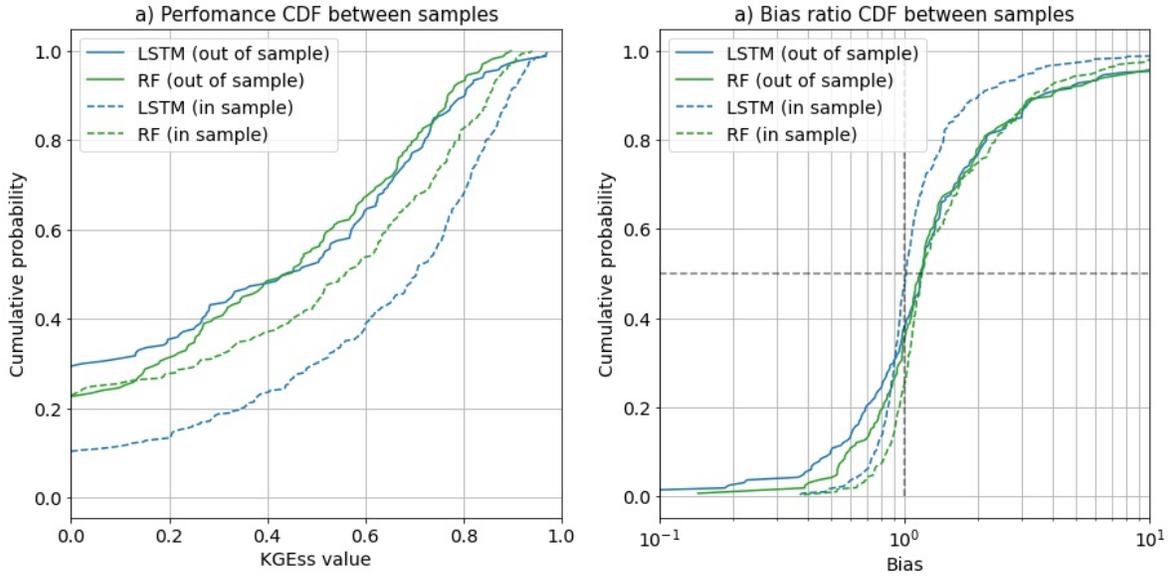


Figure 11. Performance CDFs for temporal (in-sample) and spatial (out-of-sample) generalization.

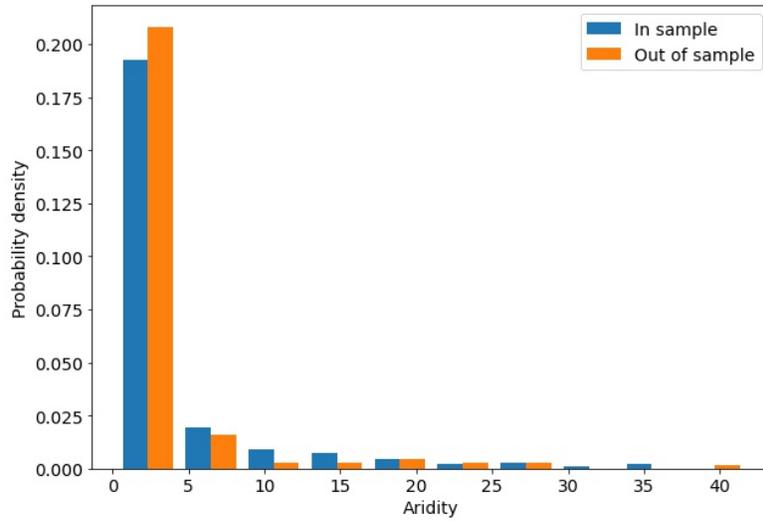


Figure A-1: Comparison of the histogram for both samples used in the performance analysis.

Table 1. Comparison between Linear Reservoir and LSTM

Linear Reservoir	LSTM
$\frac{dS}{dt} = I - O$	$\frac{dc}{dt} = g(x, h)$
<i>S</i> : water storage	<i>c</i> : information storage
<i>I</i> : input	<i>x</i> : input
<i>O</i> : output	<i>h</i> : output
$S(t) = \alpha \cdot S(t-1) + \beta \cdot (I - O)$	$c(t) = \alpha \cdot c(t-1) + \beta \cdot g(x, h)$
$\alpha = 1$	$\alpha = f(t)$
$\beta = 1$	$\beta = i(t)$
$O = k \cdot S(t)$	$h = k \cdot \bar{c}(t)$
$k =]0,1[$	$k = o(t)$
	$\bar{c}(t) = c(t)$ normalized between $] -1,1[$

Table 2. Evaluation period performance statistics.

Model	Mean	Std	Min	25%	50%	75%	Max	# Positive	# Best
GR4J	0.417	0.960	-11.621	0.311	0.561	0.789	0.960	303	98
RF	-2.411	39.310	-703.128	0.075	0.563	0.762	0.940	249	52
LSTM	-3.282	65.269	-1170.490	0.442	0.704	0.826	0.971	289	172

Table 3. Out-of-sample performance statistics.

Model	Mean	Std	Min	25%	50%	75%	Max	# Positive	# Better
RF	-0.356	2.605	-17.501	0.118	0.45	0.666	0.897	130	89
LSTM	-2.474	17.341	-203.124	-0.103	0.429	0.678	0.968	116	78

Table A-4 Parameters and search range used in the GR4J optimization.

Parameter	Description	Searching range
Alpha1	Amplification factor for CR2MET precipitation product	0-2.5
Alpha2	Amplification factor for MSWEP precipitation product	0-2.0
x1	Storage production capacity	0-5000
x2	Amplification of water exports	-10 to 10
x3	Storage routing capacity	0-1500
x4	Time-delay between the initial and maximum values of the hydrograph	0.5-4.5
Lambda	Exponent of Box-Cox transformation	0-2.0

Table A-5 Variables used in the GR4J model.

Variable	Description
PP_cr2-0	Precipitation in the same day (“0”) of the mean streamflow from CR2MET product
PP_mswep-0	Precipitation in the same day (“0”) of the mean streamflow from MSWEP product
ETP-0	Potential Evapotranspiration in the same day (“0”) of the mean streamflow
Q	Mean streamflow

Table A-6 Variables used in the Random Forest model.

n°	Attribute or variable	n°	Attribute or variable	n°	Attribute or variable	n°	Attribute or variable	n°	Attribute or variable
1	area	31	fp_frac	61	PP_cr2-3	91	PP_mswep-16	121	Tmean-6
2	aridity_cr2met	32	frac_snow_cr2met	62	PP_cr2-4	92	Q	122	Tmean-7
3	aridity_mswep	33	frac_snow_mswep	63	PP_cr2-5	93	shrub_frac	123	Tmean-8
4	big_dam	34	gauge_lat	64	PP_cr2-6	94	slope_mean	124	Tmean-9
5	carb_rocks_frac	35	gauge_lon	65	PP_cr2-7	95	snow_frac	125	Tmean-10
6	crop_frac	36	grass_frac	66	PP_cr2-8	96	sur_rights_flow	126	Tmean-11
7	Day	37	gw_rights_flow	67	PP_cr2-9	97	sur_rights_n	127	Tmean-12
8	elev_gauge	38	gw_rights_n	68	PP_cr2-10	98	Tmax-0	128	Tmean-13
9	elev_max	39	high_prec_dur_cr2met	69	PP_cr2-11	99	Tmax-1	129	Tmean-14
10	elev_mean	40	high_prec_dur_mswep	70	PP_cr2-12	100	Tmax-2	130	Tmean-15
11	elev_med	41	high_prec_freq_cr2met	71	PP_cr2-13	101	Tmax-3	131	Tmean-16
12	elev_min	42	high_prec_freq_mswep	72	PP_cr2-14	102	Tmax-4	132	Tmin-0
13	ETP-0	43	imp_frac	73	PP_cr2-15	103	Tmax-5	133	Tmin-1
14	ETP-1	44	lc_barren	74	PP_cr2-16	104	Tmax-6	134	Tmin-2
15	ETP-2	45	lc_glacier	75	PP_mswep-0	105	Tmax-7	135	Tmin-3
16	ETP-3	46	low_prec_dur_cr2met	76	PP_mswep-1	106	Tmax-8	136	Tmin-4
17	ETP-4	47	low_prec_dur_mswep	77	PP_mswep-2	107	Tmax-9	137	Tmin-5
18	ETP-5	48	low_prec_freq_cr2met	78	PP_mswep-3	108	Tmax-10	138	Tmin-6
19	ETP-6	49	low_prec_freq_mswep	79	PP_mswep-4	109	Tmax-11	139	Tmin-7
20	ETP-7	50	Month	80	PP_mswep-5	110	Tmax-12	140	Tmin-8
21	ETP-8	51	nf_frac	81	PP_mswep-6	111	Tmax-13	141	Tmin-9
22	ETP-9	52	p_mean_cr2met	82	PP_mswep-7	112	Tmax-14	142	Tmin-10
23	ETP-10	53	p_mean_mswep	83	PP_mswep-8	113	Tmax-15	143	Tmin-11
24	ETP-11	54	p_mean_spread	84	PP_mswep-9	114	Tmax-16	144	Tmin-12
25	ETP-12	55	p_seasonality_cr2met	85	PP_mswep-10	115	Tmean-0	145	Tmin-13
26	ETP-13	56	p_seasonality_mswep	86	PP_mswep-11	116	Tmean-1	146	Tmin-14
27	ETP-14	57	pet_mean	87	PP_mswep-12	117	Tmean-2	147	Tmin-15
28	ETP-15	58	PP_cr2-0	88	PP_mswep-13	118	Tmean-3	148	Tmin-16
29	ETP-16	59	PP_cr2-1	89	PP_mswep-14	119	Tmean-4	149	wet_frac
30	forest_frac	60	PP_cr2-2	90	PP_mswep-15	120	Tmean-5		

Table A-4 Variables used in the LSTM model.

n°	Attribute or variable	n°	Attribute or variable
1	area	31	p_mean_cr2met
2	aridity_cr2met	32	p_mean_mswep
3	aridity_mswep	33	p_mean_spread
4	big_dam	34	p_seasonality_cr2met
5	carb_rocks_frac	35	p_seasonality_mswep
6	crop_frac	36	pet_mean
7	elev_gauge	37	shrub_frac
8	elev_max	38	slope_mean
9	elev_mean	39	snow_frac
10	elev_med	40	sur_rights_flow
11	elev_min	41	sur_rights_n
12	forest_frac	42	wet_frac
13	fp_frac	43	PP_cr2-0
14	frac_snow_cr2met	44	PP_mswep-0
15	frac_snow_mswep	45	Tmin-0
16	grass_frac	46	Tmean-0
17	gw_rights_flow	47	Tmax-0
18	gw_rights_n	48	ETP-0
19	high_prec_dur_cr2met		
20	high_prec_dur_mswep		
21	high_prec_freq_cr2met		
22	high_prec_freq_mswep		
23	imp_frac		
24	lc_barren		
25	lc_glacier		
26	low_prec_dur_cr2met		
27	low_prec_dur_mswep		
28	low_prec_freq_cr2met		
29	low_prec_freq_mswep		
30	nf_frac		