

1 **NOTE:**

2

3 Reviewers' original comments are in Cambria type with yellow highlighting.

4

5 *Text from the revised manuscript is in bold italic font.*

6

7 Page/Line/Figure numbers given by reviewers were left unchanged in their text and
8 refer to the originally submitted version of the manuscript

9

10 [Line XXX] indicates line number in the annotated manuscript

Reply to Reviewer #1

I. Summary of the most important scientific findings

(1) First, I would like to express that I very much like the subject of the paper and the research questions that have been investigated. The regionalization of hydrological model parameters based on catchments attributes has been a long history in the PUB initiative, and, over the last decade with the developments of deep learning tools gained some new momentum.

(2) This paper addresses this area of research, by using symbolic regression methods and combination with genetic programming, to not only map the (transfer) relationship of model parameters with physio-geographic properties of the catchments, but also make these relations interpretable in terms of concrete mathematical expressions that are derived.

(3) In this way, the paper provides a valuable approach towards the development of “explainable AI” in hydrology and water management.

II. Novelty and contributions of this paper

(4) In detail, the paper investigates how the single parameter N of the Smooth-Minima approach to perform base flow separation can be estimated and regionalized from tracer-based measurements and physio-geographical data for 855 catchments in the US.

(5) Optimal algebraic equations and parameters were estimated using a genetic programming approach and are compared to standard literature formulations demonstrating superior performance.

Thank you for this positive and encouraging summary of our work. We appreciate the reviewer’s recognition of the relevance of the research question, the methodological contribution of symbolic regression, and the potential of the proposed approach for developing interpretable and explainable AI tools in hydrology.

III. Critical comments and recommendations to the editor

While I very much like the general scope of the paper as stated before, here are a number of critical questions, concerns or suggestions I would like to make and that are listed in the following:

(6) Why did you use an genetic programming approach? It is well known that this approach is limited to a few number of variables and levels of branches. You are generating a discrete optimization problem that is in my opinion very difficult to handle. In recent work by e.g. Feigl et al. (2020, cited by you) this is approached/solved by generation a grammar and “compressing” it with a VAE (and some constrains) into a latent space which is continuous and therefore usable for “gradient-based” search methods. Why didn’t you use such an approach?

Thank you for this insightful and constructive comment. We used the relatively simple genetic programming approach instead of more complex grammar-based encoding

49 combined with variational autoencoders because the search space of our problem is
50 restricted in both dimensionality and structural complexity. First, prior to symbolic
51 regression, we applied a mutual information-based feature selection procedure to reduce
52 the candidate predictors to nine relevant variables. Second, the set of mathematical
53 operators was constrained to addition, subtraction, multiplication, division, power-law, and
54 logarithm. Nested operators were not allowed except for addition and subtraction. The
55 internal complexity of expressions within the power-law and logarithmic operators was
56 limited to a maximum of 3 and the total expression complexity was capped at 20. These
57 constraints substantially reduce the search space and mitigate the limitations of GP in
58 highly complex symbolic structures. In addition, the GP-based symbolic regression is
59 computationally manageable and have been successfully applied in other studies (Dorgo et
60 al., 2021; Pumo and Noto, 2023; Razaq et al., 2016).

61 To assess whether the derived SR expressions depend on the choice of optimization
62 methods, we implemented the reinforcement learning-guided search. The resulting
63 expressions are: $A^{0.22} + 3.70$, $(A + 901)^{0.24} + K_{sat}^{0.28}$, $0.31 * (A + 549)^{0.31} + K_{sat}^{0.34} +$
64 $3.15f_{SWE} + 1.84$, with R^2 values of 0.47, 0.52, and 0.55, respectively. These expressions
65 exhibit high similarity and comparable predictive performance to those obtained using the
66 GP-based approach. This consistency suggests that the identified SR expressions are robust
67 and independent from the symbolic regression optimization methods.

68

69 (7) You used $N=5$ (FD) and $N=1.6*A^{0.2}$ (FPL) as benchmark models taken as typical
70 parameterizations from the literature. I think that is ok to see the limitations and the lower
71 end of performance. But why didn't you use that equation structure ($a*A^b$) and looked ho
72 well it worked in comparison?

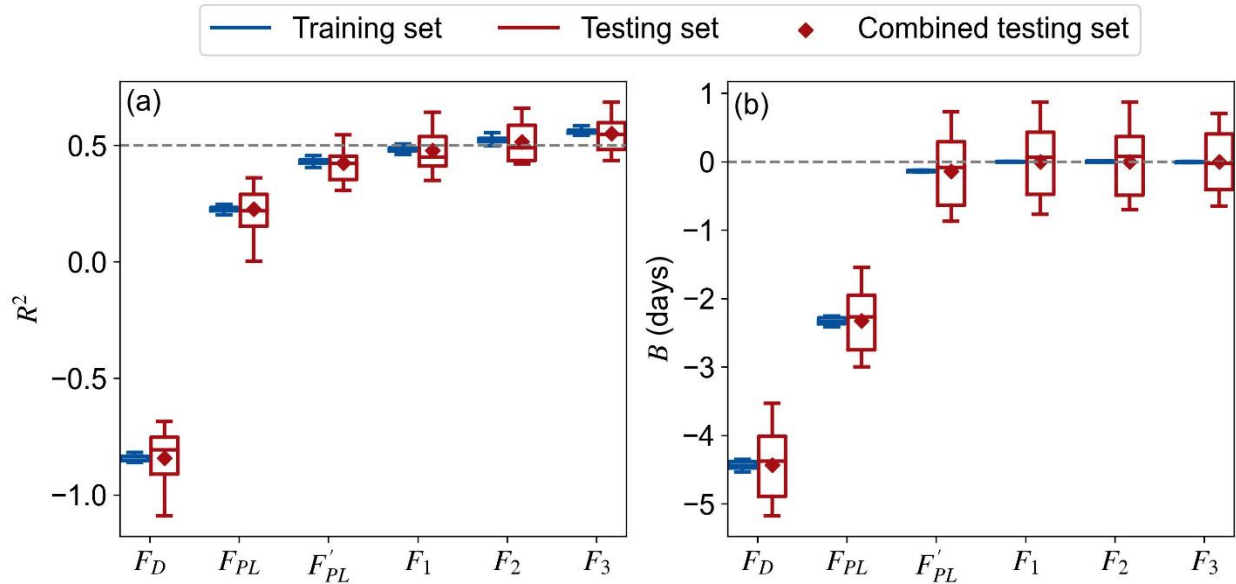
73 We appreciate this insightful comment. We tested a calibrated version of $N = a \cdot A^b$ as
74 suggested. The calibrated power-law model achieves an R^2 of 0.42, lower than those of the
75 SR-derived ones (F_1 , F_2 , and F_3 ranging from 0.48-0.55). This implies the limitation in the
76 $N = a \cdot A^b$ function, highlighting the advantage of using SR framework for the prediction.
77 We added the calibrated power-law results in Figure 4 and some descriptions in section
78 3.1, 4.2, and 4.3.

79 [Section 3.1, lines 168-171]

80 *To further examine the explanatory capacity of the areal-based power-law relationship, we considered*
81 *a calibrated form of $N = a \cdot A^b$ (denoted as F'_{PL}), where the coefficients a and b are estimated from*
82 *the data by minimizing the squared error between the reference and predicted N_s .*

83 [Section 4.2, lines 329-332]

84 *The calibration with respect to SEC data improves the performance for F'_{PL} to median R^2 values of*
85 *0.43 and 0.42 for the training and testing sets, respectively, and R^2 of 0.42 for the combined testing set.*



86

87 **Figure 4. Performance of N predictions using the constant (F_D), power-law (F_{PL}), calibrated power-**
 88 **law (F'_{PL}), and the three SR formulas (F_1 , F_2 , and F_3) for the ten training and testing sets of the 10-**
 89 **fold cross-validation: coefficient of determination R^2 (a) and mean bias B (b).**

90 [Section 4.3, lines 402-404]

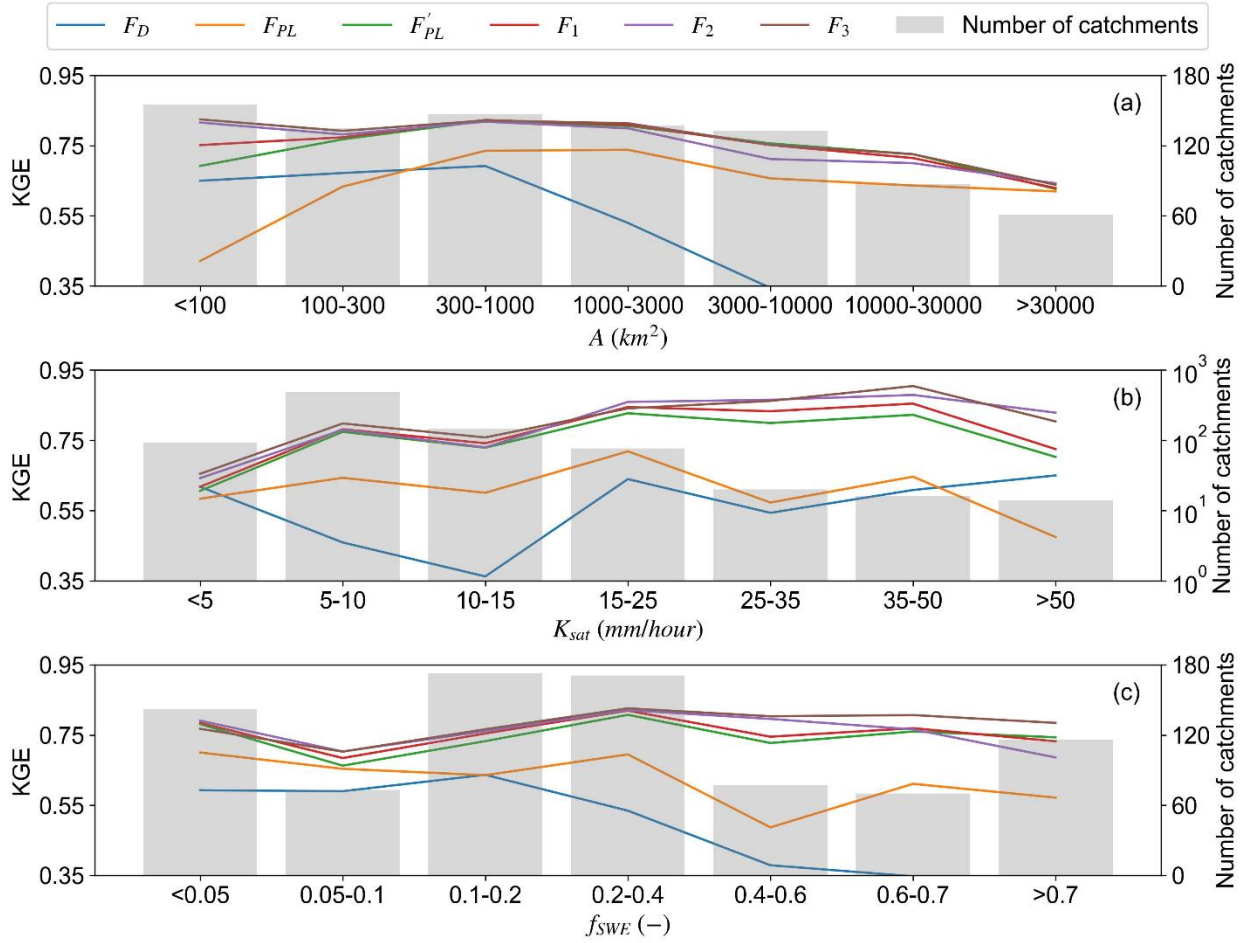
91 **The calibrated power-function (F'_{PL}) reveals improved performance to a similar level with the SR**
 92 **formulas except for the $A < 100 \text{ km}^2$ bin, indicating the necessity regional calibration.**

93 [Section 4.3, lines 412-416]

94 **In the intermediate ranges (5–25 mm/hour), the two regionally fitted power-law formulas (F_1 and F'_{PL})**
 95 **outperform the conventional power-law formula (F_{PL}). Although F_1 and F'_{PL} do not explicitly include**
 96 **K_{sat} , their performance remains comparable to F_2 and F_3 , which incorporate K_{sat} . This indicates that**
 97 **regional fitting may compensate the effects of K_{sat} for small and medium values.**

98 [Section 4.3, lines 423-425]

99 **In contrast, for $f_{SWE} < 0.4$, F'_{PL} , F_1 , and F_2 perform similarly to F_3 , demonstrating that regionally**
 100 **calibrated coefficients can partially offset the lack of explicit consideration of snow-related variables.**



101
 102 **Figure 6. Performance of baseflow separation for different ranges of catchment area (a), catchment-**
 103 **averaged saturated hydraulic conductivity (b), and snow day fraction (c). The right y-axis for panel b is**
 104 **in logarithmic scale.**
 105

106 (8) For such benchmark experiments, we often use a ML approach (e.g. Xgboost) in order to
 107 analyze the “maximum” possible performance given the calibration data and catchment
 108 properties, to see how well the CF approach does in comparison.

109 Thank you for this valuable suggestion. We evaluated a random forest (RF) model, which
 110 achieved a relatively high predictive performance ($R^2 = 0.80$). This result is reported in the
 111 Discussion to provide a reference. However, since predictive optimization is not the
 112 primary focus, we did not further expand on additional ML models or extensive
 113 benchmarking experiments. The primary objective is to identify generalizable and
 114 physically interpretable relationships that can inform parameter regionalization and
 115 hydrological understanding. While RF achieves higher predictive skill, its tree-based
 116 ensemble structure does not provide an explicit functional form linking catchment
 117 attributes to N . In contrast, SR yields closed-form equations that enables derivative
 118 analysis. For instance, differentiating F_3 gives the marginal effects of each attribute on N :
 119 $\frac{\partial N}{\partial A} = a_3 d_3 (A + c_3)^{a_3 - 1}$, $\frac{\partial N}{\partial K_{sat}} = e_3 K_{sat}^{e_3 - 1}$, and $\frac{\partial N}{\partial f_{SWE}} = f_3$. These derivatives can quantify

120 the contributions from the different geomorphological, edaphic, and climatic factors on
121 baseflow separation parameter. This mathematical functions are particularly valuable for
122 understanding the physical relationship between catchment attributes and baseflow
123 processes. Therefore, RF and SR should be viewed as complementary rather than
124 competing approaches: RF provides a benchmark for achievable predictive performance,
125 whereas SR offers structural transparency that facilitates theoretical interpretation and
126 model integration.

127 We clarified this scope and the complementary roles of RF and SR in the introduction and
128 discussion sections.

129 [Section 1, lines 107-111]

130 *This study should not be viewed as an effort to assert a superior utility of SR over other machine
131 learning models in the regionalization of baseflow parameters. Instead, the SR formulas serve as post-
132 hoc interpretability tools to complement other black box models, enhancing the transparency of the
133 underlying relationship between hydrological signatures and catchment attributes (Rudin, 2019).*

134 [Section 5.1, lines 471-494]

135 *In this study, we used SR to derive mathematical expressions for the predictions of N using 9
136 catchment attributes. Across ten cross-validation iterations, the identified expressions
137 exhibited consistent structures, predictors, and nearly identical regression coefficients, indicating that
138 SR can yield stable functional relationships between catchment attributes and N. Compared to the RF-
139 based predictions reported by Lin et al. (2026), the SR-based approach showed lower predictive skill
140 ($R^2 = 0.54$ vs. 0.80), reflecting the trade-off between predictive accuracy and interpretability. While RF
141 achieves superior predictive performance, it functions as a 'black-box' ensemble, offering no explicit
142 functional form to clarify whether environmental controls operate additively, multiplicatively, or
143 through nonlinear transformations. In contrast, SR provides structural transparency by yielding a
144 closed-form equation, facilitating direct analytical insights (Häfner et al., 2023; Karpatne et al., 2024).
145 This explicit representation enables rigorous sensitivity analysis via differentiation; for instance, the
146 marginal effects derived from equation F_3 quantify how geomorphic and climatic factors jointly govern
147 N. By trading a degree of predictive skill for parsimony, SR transforms the problem from simple
148 estimation into a hypothesis-generating exercise, providing compact transfer functions that are easily
149 integrated into regionalization frameworks (Feigl et al., 2022; Samaniego et al., 2010). Therefore, RF
150 and SR should be viewed as complementary rather than competing approaches: RF provides a
151 benchmark for predictive performance, while SR offers structural transparency that facilitates
152 theoretical interpretation and model integration.*

153

154 (9) While Appendix A is given some text book information on the principle of streamflow
155 separation, I would be much more interested in how the separation was performed and
156 especially what kind of uncertainties are involved.

157 Thank you for this important comment. We added descriptions on the SMM procedures for
158 baseflow separation in Section 3.1.

159 [Section 3.1, lines 157-162]

160 *The SMM procedure involves partitioning daily streamflow into non-overlapping N-day intervals and
161 identifying the minimum value within each segment. These minimum points ($Q_1, Q_2, \dots, Q_b, \dots$) are*

162 *then screened using a filtering coefficient (M): a point is discarded if $M \cdot Q_i$ exceeds the value of either*
163 *adjacent minimum. Finally, the baseflow series is constructed by linearly interpolating the remaining*
164 *minima. The method involves two key parameters: the segment length parameter (N) and the filtering*
165 *coefficient parameter (M).*

166 Studies found that the block length parameter of SMM constitute the primary source of
167 uncertainty (Lin et al., 2026; Stoelzle et al., 2020). That is why we proposed to regionalize
168 the SEC-optimized N parameters using SR instead of adopting the default values. Using
169 SEC-optimized parameter in baseflow separation ensures the different streamflow
170 components are consistency with the tracer dynamics.

171

172 **Some technical questions of concern are the following:**

173 **(10) The equations in Table 2 look to me strongly depend on what units you are choosing**
174 **for Ks and A. Also, in F2 and F3 you add A and Ks which makes no sense in terms of units!**
175 **What are the impacts of different choices on the overall estimation process. Is there**
176 **guidance needed for any other user application?**

177 Thank you for this insightful comment. We acknowledge that the additive structure $A +$
178 K_{sat} lacks hydrological justification and its appearance may be the coincidence of unit
179 inconsistency of K_{sat} . In the previous version, we didn't take out the $\times 100$ scaling factor
180 from K_{sat} , which may be the cause of this dimensionally inconsistency. To address this
181 issue, we revised the SR training configurations in the following ways: A) the unit of
182 saturated hydraulic conductivity (Ksat) was converted to mm/hour; B) multiplication,
183 division, power-law, and logarithmic operators were not allowed to be nested within
184 operators of the same type; C) the internal complexity of expressions inside power-law and
185 logarithmic operators was restricted to a maximum value 3; and D) recursive formulations
186 were not allowed. These new configurations result in more clear SR expressions: $N =$
187 $A^{a_1} + b_1$, $N = d_2(A + c_2)^{a_2} + (K_{sat})^{e_2} + b_2$, and $N = d_3(A + c_3)^{a_3} + (K_{sat})^{e_3} + f_3 \cdot f_{SWE} +$
188 b_3 . Compared to the previous expressions, the new ones reveal no direct additions between
189 variables with different units and every variable is in a single term with clear form.

190 We updated Table 2 and paragraphs in section 3.2 to reflect the changes.

191 [Section 3.2, Lines 186-188]

192 *The nine representative catchment attributes in original units (Table 1) were used as*
193 *predictors without normalization to maintain interpretability for the SR expressions.*

194 [Section 3.2, Lines 191-206]

195 *To control the search space and ensure physically interpretable expressions, several structural*
196 *constraints were imposed in SR model training. Multiplication, division, power-law, and*
197 *logarithmic operators were not allowed to be nested within operators of the same type. The*
198 *internal complexity of expressions inside power-law and logarithmic operators was restricted*
199 *to a maximum value of 3. The maximum allowable total complexity was set to 20. Expression*
200 *complexity is defined as the sum of the complexity index assigned to each component in the*
201 *equation. Take $N = 1.6 \times A^{0.2}$ as an example, if multiplication and power-law operators are*

202 *each assigned a complexity of 2 and constants and input variables are assigned a complexity*
203 *of 1, the total complexity of the expression is calculated as $2 + 2 + 1 + 1 + 1 = 7$. In this*
204 *study, all operators were assigned a uniform complexity index of 1 to avoid bias toward*
205 *specific functional forms. Recursive formulations (i.e., expressions where the output variable*
206 *appears as an input to itself) were not permitted to ensure model interpretability and avoid*
207 *trivial or ill-posed solutions.*

208 [Table 2]

209 **Table 2. The calibrated F'_{pL} and the three stable forms of SR expressions across all ten iterations of the cross-validation. The**
 210 **numbers within the brackets of the first row are the complexity indices of the expressions.**

Iteration	F'_{LR} (5)	F_1 (5)	F_2 (13)	F_3 (17)
1	$3.00 \times A^{0.15}$	$A^{0.23} + 3.66$	$0.40(A + 529)^{0.31} + K_{sat}^{0.30} + 2.39$	$0.27(A + 640)^{0.33} + K_{sat}^{0.35} + 3.71f_{SWE} + 1.79$
2	$3.10 \times A^{0.15}$	$A^{0.23} + 3.76$	$0.62(A + 642)^{0.28} + K_{sat}^{0.32} + 1.37$	$0.40(A + 652)^{0.30} + K_{sat}^{0.36} + 3.38f_{SWE} + 1.26$
3	$3.10 \times A^{0.15}$	$A^{0.23} + 3.74$	$0.47(A + 526)^{0.30} + K_{sat}^{0.28} + 2.19$	$0.32(A + 575)^{0.32} + K_{sat}^{0.33} + 3.55f_{SWE} + 1.80$
4	$3.25 \times A^{0.14}$	$A^{0.23} + 3.86$	$0.39(A + 334)^{0.30} + K_{sat}^{0.31} + 2.78$	$0.26(A + 361)^{0.33} + K_{sat}^{0.35} + 3.50f_{SWE} + 2.32$
5	$3.17 \times A^{0.15}$	$A^{0.23} + 3.79$	$0.39(A + 371)^{0.31} + K_{sat}^{0.31} + 2.64$	$0.35(A + 547)^{0.31} + K_{sat}^{0.35} + 3.60f_{SWE} + 1.61$
6	$3.03 \times A^{0.15}$	$A^{0.23} + 3.76$	$0.56(A + 602)^{0.28} + K_{sat}^{0.29} + 1.77$	$0.41(A + 754)^{0.30} + K_{sat}^{0.34} + 3.63f_{SWE} + 1.21$
7	$3.14 \times A^{0.15}$	$A^{0.23} + 3.83$	$0.44(A + 401)^{0.30} + K_{sat}^{0.29} + 2.58$	$0.30(A + 443)^{0.32} + K_{sat}^{0.33} + 3.28f_{SWE} + 2.18$
8	$3.16 \times A^{0.15}$	$A^{0.23} + 3.85$	$0.44(A + 421)^{0.30} + K_{sat}^{0.31} + 2.49$	$0.31(A + 491)^{0.32} + K_{sat}^{0.35} + 3.54f_{SWE} + 1.95$
9	$3.13 \times A^{0.15}$	$A^{0.23} + 3.81$	$0.46(A + 481)^{0.30} + K_{sat}^{0.28} + 2.42$	$0.30(A + 525)^{0.32} + K_{sat}^{0.33} + 3.38f_{SWE} + 2.06$
10	$3.14 \times A^{0.15}$	$A^{0.23} + 3.79$	$0.43(A + 545)^{0.30} + K_{sat}^{0.31} + 2.29$	$0.29(A + 627)^{0.32} + K_{sat}^{0.35} + 3.57f_{SWE} + 1.86$

211

212 (11) Looking at Fig. 2, I find some information on the derived equations, but I do not see
213 what kind of function space is actually generated to be used in the genetic programming
214 approach – does the GE approach in your case for recursion? Please clarify.

215 We thank the reviewer for pointing this out. We added the function space used in the
216 symbolic regression framework in section 3.2. The candidate expressions were constructed
217 from input variables (catchment attributes), free constants, and a predefined operator set
218 (addition, subtraction, multiplication, division, power-law, and logarithm). Recursive
219 formulations were not allowed. In addition, structural constraints were imposed to prevent
220 excessive nesting of nonlinear operators and to maintain physical interpretability.

221 [Section 3.2, lines 188-206]

222 *The function space is consisted of the catchment attributes, free constants, and a set of*
223 *mathematical operators: addition (+), subtraction (-), multiplication (×), division (÷), power-*
224 *law (power), and logarithm (log).*

225 *To control the search space and ensure physically interpretable expressions, several structural*
226 *constraints were imposed in SR model training. Multiplication, division, power-law, and*
227 *logarithmic operators were not allowed to be nested within operators of the same type. The*
228 *internal complexity of expressions inside power-law and logarithmic operators was restricted*
229 *to a maximum value of 3. The maximum allowable total complexity was set to 20. Expression*
230 *complexity is defined as the sum of the complexity index assigned to each component in the*
231 *equation. Take $N = 1.6 \times A^{0.2}$ as an example, if multiplication and power-law operators are*
232 *each assigned a complexity of 2 and constants and input variables are assigned a complexity*
233 *of 1, the total complexity of the expression is calculated as $2 + 2 + 1 + 1 + 1 = 7$. In this*
234 *study, all operators were assigned a uniform complexity index of 1 to avoid bias toward*
235 *specific functional forms. Recursive formulations (i.e., expressions where the output variable*
236 *appears as an input to itself) were not permitted to ensure model interpretability and avoid*
237 *trivial or ill-posed solutions.*

238

239 Overall, I would recommend to the editor to accept the manuscript after some (major)
240 revisions.

241

242 IV. Minor and specific comments

243 L95: Does the data set included nested catchments/information?

244 Thank you for the comment. Yes, the dataset may include some nested catchments. In this
245 study, all catchments were treated as independent units, and no explicit consideration of
246 upstream–downstream relationships was incorporated in the modeling, as the focus is on
247 deriving generalizable relationships at the catchment scale.

248

249 L117: Units of the variables in the equation?

250 We thank the reviewer for pointing this out. The units for all variables in the equation were
251 added.

252 [Section 3.1, lines 165-167]

253 *In the literature, N is often default to 5 days or predicted using a power-law relationship with*
254 *catchment area, namely $N = 1.6 \times A^{0.2}$, where A is the catchment area in km^2 and N is*
255 *expressed in days.*

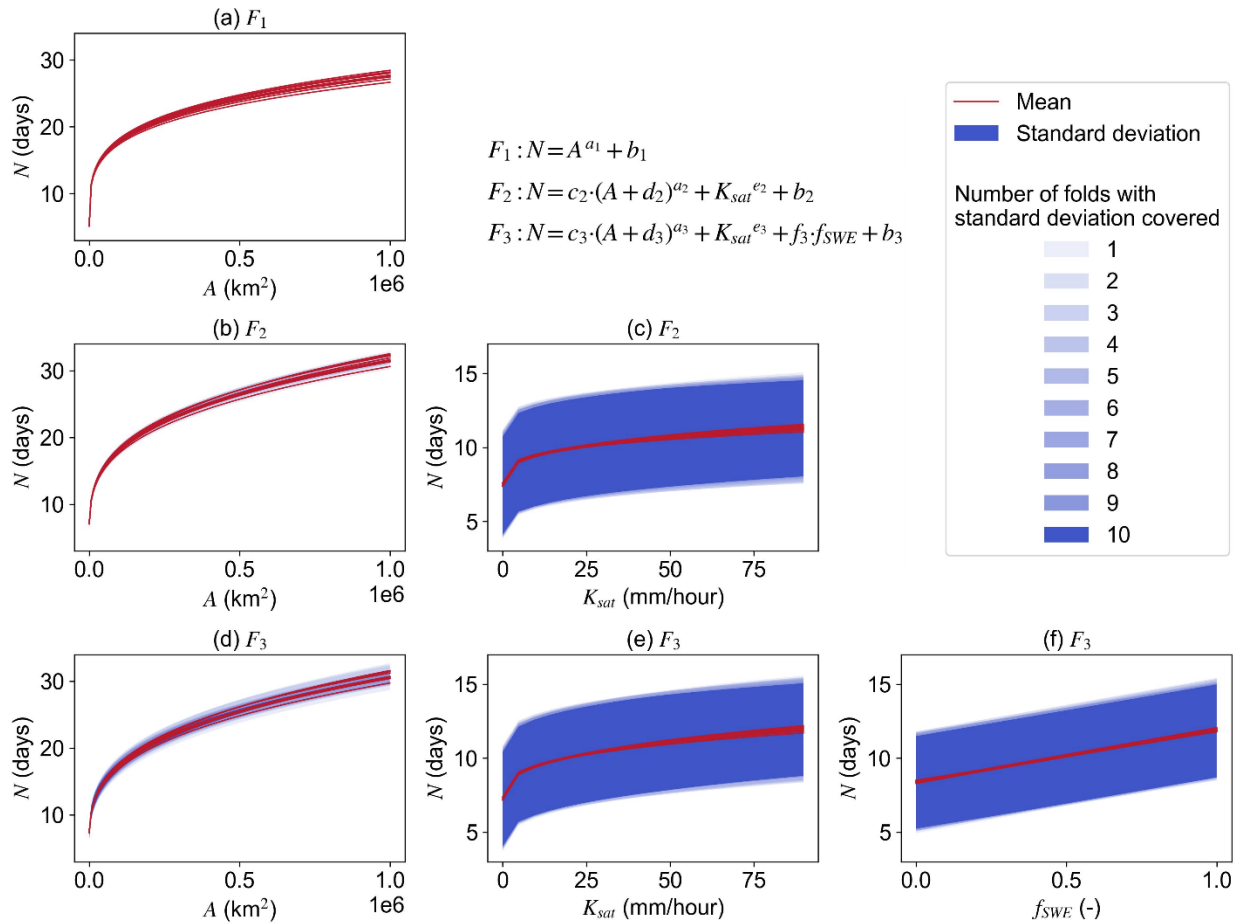
256

257 L199: In Fig. 3 the individual (10) lines are hardly visible or distinguishable

258 Thank you for this comment. The individual lines correspond to the consistent SR formulas
259 of the ten cross-validation folds. They appear closely clustered because the coefficients of
260 the formulas are highly similar. This behavior reflects the stability of the SR-derived
261 relationships rather than a visualization limitation. We added some descriptions to
262 clarified this point.

263 [Section 4.1, lines 287-299]

264 *Figure 3 further illustrates the behavior of these formulas. For F_1 , the nearly identical exponents*
265 *(~ 0.23) and intercepts (3.66–3.86 days) result in almost overlapping curves (Figure 3a), indicating a*
266 *stable power-law relationship between N and A , with a diminishing rate of increase. For F_2 , the*
267 *similarly constrained exponents (0.28–0.32) and intercepts (1.37–2.78 days) produce tightly clustered*
268 *response curves (Figure 3b-c), showing that both A and K_{sat} contribute positively to N . F_3 extends F_2*
269 *by introducing f_{SWE} as a linear term. The replicated formulas still exhibit closely grouped slopes*
270 *(3.28–3.71 days) and intercepts (1.21–2.32 days), which explains the clustering of curves in Figure 3d–*
271 *f. The marginal relationships of N with A and K_{sat} in F_3 remain consistent with those in F_2 , whereas*
272 *increasing f_{SWE} leads to an approximately linear increase in N , at a rate of about 0.3–0.4 days per 0.1*
273 *increment in snow fraction. Overall, SR identifies A as the most influential factor in predicting N , as*
274 *evidenced by its presence in all SR-derived formulas. The narrower ranges of predicted N s in Figure*
275 *3b and d also suggest that A exerts greater influence than K_{sat} and f_{SWE} .*



276

277 **Figure 3. Marginal relationship of N on different predictors (A , K_{sat} , and f_{SWE}) that consist of the SR**
 278 **expressions (F_1 , F_2 , and F_3). Each line represents one of the ten instances of F_1 , F_2 , and F_3 . Panels a,**
 279 **b, and d are for A ; panels c and e are for K_{sat} ; panel f is for f_{SWE} .**

280

281 **L268: I am very surprised by the relatively small spread of K_s values that is observed in the**
 282 **catchments?**

283 Thank you for this comment. The K_{sat} values were derived by taking spatial average over
 284 the catchment, which smoothes the extremes, resulting in a narrow value range. In
 285 addition, these K_{sat} values (ranging from 0 to 80 mm/hour) are in consistent with the
 286 ranges reported by others (Choulga et al., 2024; Gupta et al., 2021).

287

288 **L336-340: Can you provide some reference for the statements.**

289 Thanks for this suggestion. We added two references to support the statements.

290 [Section 5.3, lines 540-544]

291 **Moreover, the effects of K_{sat} and f_{SWE} are more evident when A is smaller than 100 km^2 (Figure 6a).**
 292 **As catchment area increases ($A > 100 \text{ km}^2$), their influence is outweighed by A , highlighting the**

293 *dominant role of drainage area in shaping the streamflow response time (McGlynn et al., 2004;*
294 *Sólyom and Tucker, 2007).*

295

296 References

297 Choulga, M. et al., 2024. Technical note: Surface fields for global environmental modelling.
298 Hydrol. Earth Syst. Sci., 28(13): 2991-3036. DOI:10.5194/hess-28-2991-2024

299 Dorgo, G., Kulcsar, T., Abonyi, J., 2021. Genetic programming-based symbolic regression for
300 goal-oriented dimension reduction. Chemical Engineering Science, 244: 116769.

301 DOI:<https://doi.org/10.1016/j.ces.2021.116769>

302 Feigl, M. et al., 2022. Automatic Regionalization of Model Parameters for Hydrological
303 Models. Water Resources Research, 58(12): e2022WR031966.

304 DOI:<https://doi.org/10.1029/2022WR031966>

305 Gupta, S., Lehmann, P., Bonetti, S., Papritz, A., Or, D., 2021. Global Prediction of Soil
306 Saturated Hydraulic Conductivity Using Random Forest in a Covariate-Based GeoTransfer
307 Function (CoGTF) Framework. Journal of Advances in Modeling Earth Systems, 13(4):

308 e2020MS002242. DOI:<https://doi.org/10.1029/2020MS002242>

309 Häfner, D., Gemmrich, J., Jochum, M., 2023. Machine-guided discovery of a real-world rogue
310 wave model. Proceedings of the National Academy of Sciences, 120(48): e2306275120.

311 DOI:doi:10.1073/pnas.2306275120

312 Karpatne, A., Jia, X., Kumar, V., 2024. Knowledge-guided machine learning: Current trends
313 and future prospects. arXiv preprint arXiv:2403.15989.

314 Lin, Y. et al., 2026. Regionalization of Optimal Baseflow Separation using Catchment-scale
315 Characteristics. Water Resources Research, 62(3).

316 McGlynn, B.L., McDonnell, J.J., Seibert, J., Kendall, C., 2004. Scale effects on headwater
317 catchment runoff timing, flow sources, and groundwater-streamflow relations. Water
318 Resources Research, 40(7). DOI:<https://doi.org/10.1029/2003WR002494>

319 Pumo, D., Noto, L.V., 2023. Exploring the use of multi-gene genetic programming in regional
320 models for the simulation of monthly river runoff series. Stochastic Environmental
321 Research and Risk Assessment, 37(5): 1917-1941. DOI:10.1007/s00477-022-02373-1

322 Razaq, S.A. et al., 2016. Prediction of Flow Duration Curve in Ungauged Catchments Using
323 Genetic Expression Programming. Procedia Engineering, 154: 1431-1438.

324 DOI:<https://doi.org/10.1016/j.proeng.2016.07.516>

325 Rudin, C., 2019. Stop Explaining Black Box Machine Learning Models for High Stakes
326 Decisions and Use Interpretable Models Instead. Nat Mach Intell, 1(5): 206-215.

327 DOI:10.1038/s42256-019-0048-x

328 Samaniego, L., Kumar, R., Attinger, S., 2010. Multiscale parameter regionalization of a grid-
329 based hydrologic model at the mesoscale. *Water Resources Research*, 46(5).
330 DOI:<https://doi.org/10.1029/2008WR007327>

331 Sólyom, P.B., Tucker, G.E., 2007. The importance of the catchment area–length relationship
332 in governing non-steady state hydrology, optimal junction angles and drainage network
333 pattern. *Geomorphology*, 88(1): 84-108.
334 DOI:<https://doi.org/10.1016/j.geomorph.2006.10.014>

335 Stoelzle, M., Schuetz, T., Weiler, M., Stahl, K., Tallaksen, L.M., 2020. Beyond binary baseflow
336 separation: a delayed-flow index for multiple streamflow contributions. *Hydrology and
337 Earth System Sciences*, 24(2): 849-867. DOI:10.5194/hess-24-849-2020

338

339

Reply to Reviewer #2

340 The manuscript titled addresses the challenge of accurately estimating the segment length
341 parameter (N) for the Smooth Minima Method (SMM) in baseflow separation. By employing
342 Symbolic Regression (SR) on a dataset of 855 catchments across the Contiguous United
343 States, the authors aim to overcome the limitations of fixed default values or simple power-
344 law functions traditionally used in the literature. The study successfully identifies three
345 interpretable mathematical expressions that link parameter N to catchment area, saturated
346 hydraulic conductivity, and snow day fraction. The validation results demonstrate that the
347 SR-derived formulas significantly improve predictive performance (improving R2 from
348 0.23 to 0.54) and achieve higher Kling-Gupta Efficiency in baseflow separation compared to
349 existing methods. Additionally, the use of Specific Electrical Conductance (SEC) mass
350 balance serves as a robust independent validation of the proposed regionalization
351 approach. Overall, this work provides valuable insights into the physical controls of
352 baseflow processes and demonstrates the utility of interpretable machine learning in
353 hydrology.

354 However, I have several concerns regarding the physical consistency of the derived
355 formulas, the justification for choosing SR over other methods, and the robustness of the
356 performance evaluation thresholds. These issues, along with comments on terminology
357 and methodology, are detailed below.

358 Thank you for this comprehensive summary and for the positive evaluation of our work.
359 We appreciate the reviewer's recognition of the motivation of this study, the use of
360 symbolic regression to derive interpretable mathematical relationships between the SMM
361 parameter N and catchment characteristics, and the improvement in predictive
362 performance achieved by the proposed approach. We also thank the reviewer for the
363 constructive comments. These important points have been carefully considered, and
364 detailed responses and corresponding revisions have been made in the manuscript, as
365 described below.

366

367 1. The abstract states that "N increases exponentially with catchment area...". However, the
368 derived formulas are power-law functions, not exponential functions.

369 Thank you for pointing this out. We corrected the statement in the abstract and throughout
370 the manuscript.

371 [Abstract, lines 30-33]

372 *The three expressions reveal that N increases following a power-law relationship with catchment area*
373 *(A) and catchment-averaged soil saturated hydraulic conductivity (K_{sat}) with decreasing rates, while it*
374 *increases linearly with snow day fraction (f_{SWE}).*

375 [Section 4.1, lines 287-289]

376 *For F_1 , the nearly identical exponents (~ 0.23) and intercepts (3.66–3.86 days) result in almost*
377 *overlapping curves (Figure 3a), indicating a stable power-law relationship between N and A, with a*
378 *diminishing rate of increase.*

379

380 2. The introduction of SR in the current manuscript is relatively brief (lines 61-71). While
381 the authors highlight the interpretability of SR compared to Random Forest (RF), there is a
382 lack of review regarding the application of SR in the broader field of hydrology. It would
383 strengthen the rationale of the study if the authors could briefly mention successful
384 applications of SR in other hydrological contexts (Chadalawada et al., 2020) before
385 narrowing down to baseflow separation. This would demonstrate that SR is a proven tool
386 in the domain.

387 We thank the reviewer for this constructive suggestion. We added some descriptions in the
388 Introduction to address the successful applications of SR in hydrological model parameter
389 regionalization.

390 [Section 1, lines 92-102]

391 *In recent years, SR has been increasingly applied in hydrology to uncover governing*
392 *relationships in complex environmental systems, owing to its ability to balance predictive*
393 *performance and interpretability (Chadalawada et al., 2020; Sheta et al., 2023). One example*
394 *is the use of SR to extract explicit functional relationships between catchment attributes and*
395 *hydrological model parameters for ungagged catchments (Feigl et al., 2022; Li et al., 2024).*
396 *Unlike a “black-box” model such as RF, SR derives explicit and concise equations that*
397 *identify underlying data patterns, while mitigating overfitting through complexity control*
398 *(Kronberger et al., 2022; Wilstrup and Kasak, 2021). This structural transparency enables*
399 *direct interpretation of how catchment attributes govern baseflow parameter values, which in*
400 *terms influence the partitioning of streamflow (Feigl et al., 2020; Klotz et al., 2017; Sheta et*
401 *al., 2023).*

402

403 3. The Introduction highlights the role of environmental tracer data in parameter
404 optimization (Line 48) but lacks specificity. I suggest briefly listing common tracers,
405 specifically including SEC, in this paragraph. Since the reference dataset used in this study
406 is SEC-optimized (as detailed in Section 2), explicitly mentioning SEC early on would
407 provide the reader with necessary context regarding the physical basis of the ground truth
408 parameters.

409 Thank you for this helpful suggestion. We listed several commonly used tracers (e.g.,
410 specific electrical conductivity, turbidity, and stable isotopes) in the Introduction and
411 explicitly mentioned SEC.

412 [Section 1, lines 52-54]

413 *Commonly used tracers include specific electrical conductivity (SEC), turbidity, and stable*
414 *isotopes, among which SEC is the most widely applied due to its routine availability in many*
415 *monitoring programs (Mei et al., 2024).*

416

417 4. The statement in Line 84 that optimal N shows “no clear spatial patterns” is ambiguous,
418 particularly since the subsequent analysis links N to spatially structured attributes. I
419 suggest clarifying this phrasing to avoid apparent contradiction.

420 Thank you for this comment. We refined our descriptions to indicate the multifactor
421 controls on N .

422 [Section 2, lines 123-126]

423 *The optimal N parameters for these catchments are depicted in Figure 1d with most values*
424 *smaller than 17 days. Larger N values tend to occur in mountainous regions characterized by*
425 *drier climates and greater snow persistence with implicit spatial patterns, which may be*
426 *attributed to the multifactor controls of N (Lin et al. 2026).*

427

428 5. Section 3.1 describes the parameters N and M well but briefly lacks a description of the
429 algorithmic procedure. Adding one or two sentences explaining how the method iterates
430 (e.g., dividing the hydrograph into blocks of N days, identifying minima, and comparing
431 adjacent minima using M) would make the method section more self-contained.

432 Thank you for this helpful suggestion. We added the main steps of SMM as requested.

433 [Section 3.1, lines 157-161]

434 *The SMM procedure involves partitioning daily streamflow into non-overlapping N -day*
435 *intervals and identifying the minimum value within each segment. These minimum points ($Q_1,$*
436 *Q_2, \dots, Q_i, \dots) are then screened using a filtering coefficient (M): a point is discarded if $M \cdot Q_i$*
437 *exceeds the value of either adjacent minimum. Finally, the baseflow series is constructed by*
438 *linearly interpolating the remaining minima.*

439

440 6. In Section 3.2, the manuscript does not specify whether the input variables (catchment
441 characteristics) were normalized or standardized prior to SR training. Based on the
442 discussion in Section 4.1 (where specific thresholds like $A < 100 \text{ km}^2$ are mentioned), it
443 implies that raw data with physical units were used.

444 If raw data were used, the term $\left(A+K_{\text{sat}}\right)$ in Formulas F_2 and F_3 is
445 dimensionally inhomogeneous, adding Area (L^2) to Hydraulic Conductivity (L/T). This
446 makes the formula valid only for the specific combination of units chosen in this study
447 (km^2 and mm/h) and undermines the claim of physical meaningfulness.

448 Even if the variables were normalized (rendering them dimensionless), the authors should
449 interpret the physical meaning of the additive structure $\left(A+K_{\text{sat}}\right)$. This
450 structure implies a direct substitutability between catchment size and soil permeability in
451 generating baseflow delay. What is the hydrological justification for this specific interaction
452 form, rather than, for example, a multiplicative interaction $\left(A+K_{\text{sat}}\right)$ which is
453 more common in hydraulic process laws?

454 Thank you for this insightful comment. We acknowledge that the $A + K_{sat}$ structure lacks
455 hydrological justification and its appearance may be the coincidence of unit inconsistency
456 of K_{sat} . In the previous version, we didn't take out the $\times 100$ scaling factor from K_{sat} , which
457 distorted its magnitude and contributed to the emergence of this dimensionally
458 inconsistent form during the SR search.

459 To address this issue, we retrained the SR models using K_{sat} (mm/hr) in its actual
460 magnitude. In addition, we imposed stricter structural constraints to guide the search
461 toward physically interpretable forms. Specifically, multiplication, division, power-law, and
462 logarithmic operators were not allowed to be nested within operators of the same type.
463 The internal complexity of expressions within power-law and logarithmic operators was
464 limited to a maximum value of 3, and the overall expression complexity was capped at 20.
465 Recursive formulations were prohibited to avoid trivial or ill-posed solutions.

466 With these modifications, the updated SR expressions no longer contain direct additions
467 between variables with different units. The revised formulas are: $N = A^{a_1} + b_1$, $N =$
468 $d_2(A + c_2)^{a_2} + (K_{sat})^{e_2} + b_2$, and $N = d_3(A + c_3)^{a_3} + (K_{sat})^{e_3} + f_3 \cdot f_{SWE} + b_3$. In these
469 expressions, A and K_{sat} contribute through separate terms, ensuring dimensional
470 consistency and improving physical interpretability. From a hydrological perspective,
471 catchment area reflects the characteristic spatial scale of storage and drainage
472 organization, whereas K_{sat} represents soil permeability and subsurface flow
473 transmissivity. Their additive contributions represent independent influences on the
474 characteristic delay parameter N .

475 For the comment on if SR should be trained with normalized variables, we prefer to
476 retained the units and the magnitudes to preserves their physical meanings for better
477 interpretation of the resulting expressions.

478 [Section 3.2, lines 186-188]

479 *The nine representative catchment attributes in original units (Table 1) were used as*
480 *predictors without normalization to maintain interpretability for the SR expressions.*

481 [Section 3.2, lines 191-206]

482 *To control the search space and ensure physically interpretable expressions, several structural*
483 *constraints were imposed in SR model training. Multiplication, division, power-law, and*
484 *logarithmic operators were not allowed to be nested within operators of the same type. The*
485 *internal complexity of expressions inside power-law and logarithmic operators was restricted*
486 *to a maximum value of 3. The maximum allowable total complexity was set to 20. Expression*
487 *complexity is defined as the sum of the complexity index assigned to each component in the*
488 *equation. Take $N = 1.6 \times A^{0.2}$ as an example, if multiplication and power-law operators are*
489 *each assigned a complexity of 2 and constants and input variables are assigned a complexity*
490 *of 1, the total complexity of the expression is calculated as $2 + 2 + 1 + 1 + 1 = 7$. In this*
491 *study, all operators were assigned a uniform complexity index of 1 to avoid bias toward*
492 *specific functional forms. Recursive formulations (i.e., expressions where the output variable*
493 *appears as an input to itself) were not permitted to ensure model interpretability and avoid*
494 *trivial or ill-posed solutions.*

495 [Table 2]

496 **Table 2. The calibrated F'_{pL} and the three stable forms of SR expressions across all ten iterations of the cross-validation. The**
 497 **numbers within the brackets of the first row are the complexity indices of the expressions.**

Iteration	F'_{LR} (5)	F_1 (5)	F_2 (13)	F_3 (17)
1	$3.00 \times A^{0.15}$	$A^{0.23} + 3.66$	$0.40(A + 529)^{0.31} + K_{sat}^{0.30} + 2.39$	$0.27(A + 640)^{0.33} + K_{sat}^{0.35} + 3.71f_{SWE} + 1.79$
2	$3.10 \times A^{0.15}$	$A^{0.23} + 3.76$	$0.62(A + 642)^{0.28} + K_{sat}^{0.32} + 1.37$	$0.40(A + 652)^{0.30} + K_{sat}^{0.36} + 3.38f_{SWE} + 1.26$
3	$3.10 \times A^{0.15}$	$A^{0.23} + 3.74$	$0.47(A + 526)^{0.30} + K_{sat}^{0.28} + 2.19$	$0.32(A + 575)^{0.32} + K_{sat}^{0.33} + 3.55f_{SWE} + 1.80$
4	$3.25 \times A^{0.14}$	$A^{0.23} + 3.86$	$0.39(A + 334)^{0.30} + K_{sat}^{0.31} + 2.78$	$0.26(A + 361)^{0.33} + K_{sat}^{0.35} + 3.50f_{SWE} + 2.32$
5	$3.17 \times A^{0.15}$	$A^{0.23} + 3.79$	$0.39(A + 371)^{0.31} + K_{sat}^{0.31} + 2.64$	$0.35(A + 547)^{0.31} + K_{sat}^{0.35} + 3.60f_{SWE} + 1.61$
6	$3.03 \times A^{0.15}$	$A^{0.23} + 3.76$	$0.56(A + 602)^{0.28} + K_{sat}^{0.29} + 1.77$	$0.41(A + 754)^{0.30} + K_{sat}^{0.34} + 3.63f_{SWE} + 1.21$
7	$3.14 \times A^{0.15}$	$A^{0.23} + 3.83$	$0.44(A + 401)^{0.30} + K_{sat}^{0.29} + 2.58$	$0.30(A + 443)^{0.32} + K_{sat}^{0.33} + 3.28f_{SWE} + 2.18$
8	$3.16 \times A^{0.15}$	$A^{0.23} + 3.85$	$0.44(A + 421)^{0.30} + K_{sat}^{0.31} + 2.49$	$0.31(A + 491)^{0.32} + K_{sat}^{0.35} + 3.54f_{SWE} + 1.95$
9	$3.13 \times A^{0.15}$	$A^{0.23} + 3.81$	$0.46(A + 481)^{0.30} + K_{sat}^{0.28} + 2.42$	$0.30(A + 525)^{0.32} + K_{sat}^{0.33} + 3.38f_{SWE} + 2.06$
10	$3.14 \times A^{0.15}$	$A^{0.23} + 3.79$	$0.43(A + 545)^{0.30} + K_{sat}^{0.31} + 2.29$	$0.29(A + 627)^{0.32} + K_{sat}^{0.35} + 3.57f_{SWE} + 1.86$

498

499 7. Section 3.2 outlines the SR setup but omits key hyperparameters required for
500 reproducibility, such as population size, number of generations, and mutation/crossover
501 rates used in the PySR configuration. I suggest including these details, perhaps in the
502 Appendix.

503 We appreciate the reviewer's suggestion. We added detailed descriptions regarding the
504 genetic programming hyperparameters used for modeling.

505 [Section 3.2, lines 188-212]

506 *The function space is consisted of the catchment attributes, free constants, and a set of*
507 *mathematical operators: addition (+), subtraction (-), multiplication (\times), division (\div), power-*
508 *law (power), and logarithm (log).*

509 *To control the search space and ensure physically interpretable expressions, several structural*
510 *constraints were imposed in SR model training. Multiplication, division, power-law, and*
511 *logarithmic operators were not allowed to be nested within operators of the same type. The*
512 *internal complexity of expressions inside power-law and logarithmic operators was restricted*
513 *to a maximum value of 3. The maximum allowable total complexity was set to 20. Expression*
514 *complexity is defined as the sum of the complexity index assigned to each component in the*
515 *equation. Take $N = 1.6 \times A^{0.2}$ as an example, if multiplication and power-law operators are*
516 *each assigned a complexity of 2 and constants and input variables are assigned a complexity*
517 *of 1, the total complexity of the expression is calculated as $2 + 2 + 1 + 1 + 1 = 7$. In this*
518 *study, all operators were assigned a uniform complexity index of 1 to avoid bias toward*
519 *specific functional forms. Recursive formulations (i.e., expressions where the output variable*
520 *appears as an input to itself) were not permitted to ensure model interpretability and avoid*
521 *trivial or ill-posed solutions.*

522 *The SR search process was configured with the following hyperparameters: a population size*
523 *of 33, populations of 15, the crossover rate of 0.066, and evolved over 40 generations. The*
524 *goodness of fit between the reference and the predicted N_s is evaluated using the mean*
525 *squared error (MSE)*

526

527 8. The legend in Figure 7b includes categories like F_1&F_2, implying a tie in performance.
528 However, the manuscript lacks an explicit definition of the threshold used to categorize
529 these ties (e.g., $\Delta R^2 < 0.02$?). Crucially, the choice of this threshold is intrinsic to the
530 discussion on the complexity-accuracy trade-off (Section 5.3). A strict threshold (e.g.,
531 $\Delta R^2 = 0$) favors complex models F3 even for negligible gains, while a practical
532 threshold (e.g., $\Delta R^2 < 0.01$) might reveal that the simpler F1 is good enough for a
533 much larger portion of the CONUS. Therefore, beyond simply clarifying the current
534 threshold, I suggest the authors perform a brief sensitivity analysis on this threshold. For
535 instance, how does the spatial pattern of the best formula change if the tolerance for a tie is
536 set to 0.01 vs. 0.03? This would provide a more robust and transparent visualization of
537 where the added complexity of F3 yields substantial benefits versus marginal statistical
538 gains.

539 We thank the reviewer for this insightful comment. Instead of applying an arbitrary
 540 threshold to quantify performance differences, we applied the Diebold–Mariano test
 541 (Diebold and Mariano, 1995) to define ties. We ranked the R^2 s of F_1 , F_2 , and F_3 , and
 542 performed Diebold–Mariano tests for the pairwise R^2 differences at the 0.01 significance
 543 level. If the largest R^2 is significantly different than the second-largest R^2 , the SR formula
 544 associated with the largest R^2 is the single best formula of the catchment. If the difference
 545 between the largest and the second-largest R^2 s is insignificant but that between the second
 546 and the third ones is significant, the first two SR formulas are tied. Otherwise, the three SR
 547 are tied. We believe that this statistical approach offers a more objective quantification
 548 than using arbitrary R^2 as thresholds.

549 The descriptions regarding DM and how we applied it to our task were added to Section 3.3
 550 and Test S2. In addition, Figure 7 and the corresponding descriptions were updated
 551 accordingly.

552 [Section 3.3, lines 253-263]

553 *To assess whether the differences in SR-based SEC prediction performance are statistical*
 554 *significant, we apply the Diebold–Mariano (DM) test (Diebold and Mariano, 1995). Details on*
 555 *procedure and test statistics are provided in Text S2. For each catchment, pairwise DM tests*
 556 *are performed among the three SR formulas to determine whether their SEC R^2 differences*
 557 *are statistical significance at the 0.01 significance level. If the largest R^2 is significantly*
 558 *different from the second-largest R^2 , the SR formula associated with the largest R^2 is the*
 559 *single best formula of the catchment. If the difference between the largest and the second-*
 560 *largest R^2 s is insignificant but that between the second and the third ones is significant, the*
 561 *first two SR formulas are tied. Otherwise, the three SR are tied. Based on these catchment-*
 562 *scale rankings, the best-performing formula for each region is determined as the one most*
 563 *frequently ranked as the best (including ties) among all catchments within the region.*

564 [Supporting Information, Text S2]

565 *To evaluate whether the predictive performance of different formulas is statistically*
 566 *distinguishable, we applied the Diebold–Mariano (DM) test (Diebold and Mariano, 1995). The*
 567 *DM test compares predictive accuracy based on the loss differential between two competing*
 568 *models. First, the loss differential at time t is defined as*

$$d_t = L(e_{A,t}) - L(e_{B,t}) , \quad (4)$$

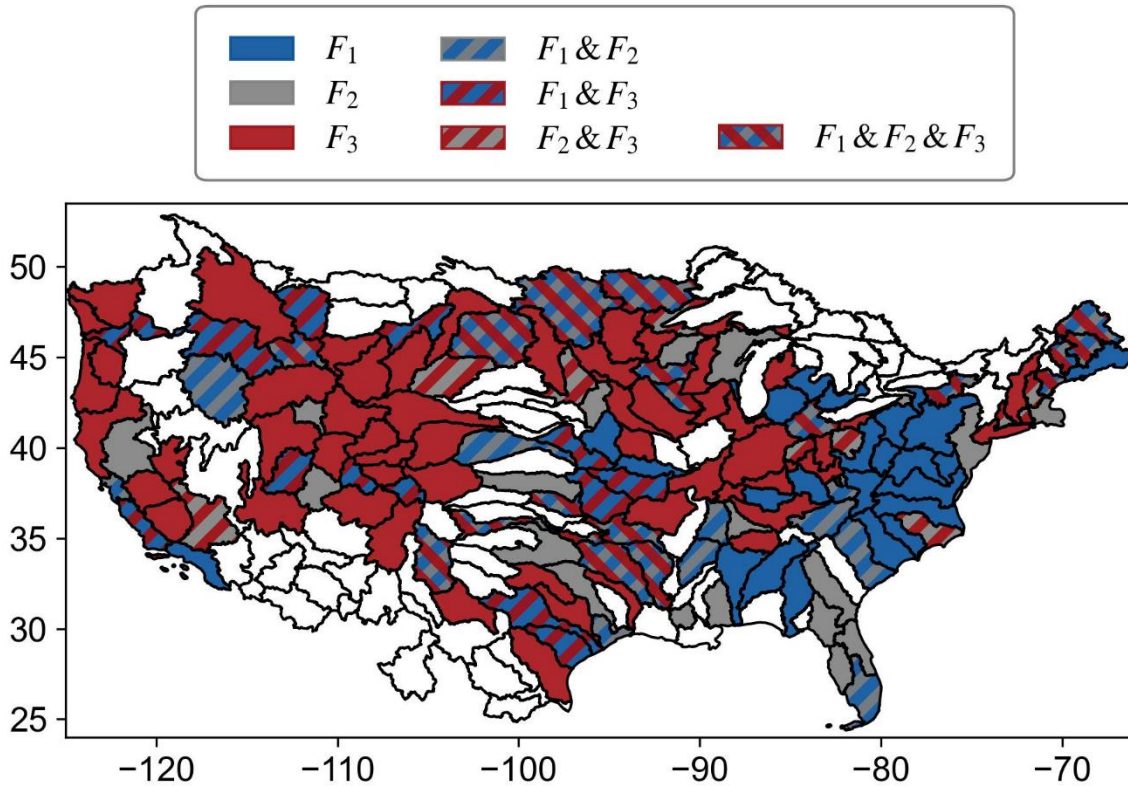
569 *where $e_{A,t}$ and $e_{B,t}$ denote the prediction errors of formulas A and B at time t , respectively,*
 570 *and $L(\cdot)$ is the loss function. In this study, the squared error was used to test for the significant*
 571 *differences between different R^2 values:*

$$L(e_t) = e_t^2 . \quad (5)$$

572 *The DM test statistic is then computed as*

$$DM = \frac{\bar{d}}{\sqrt{\text{var}(\bar{d})}} , \quad (6)$$

573 where \bar{d} is the mean loss differential and $\widehat{\text{Var}}(\bar{d})$ is the long-run variance of \bar{d} . Under the
 574 null hypothesis of equal predictive accuracy, the DM statistic asymptotically follows a
 575 standard normal distribution. A two-sided test was used to determine whether the predictive
 576 performance of two SR formulas differed significantly at the 0.01 significance level.



577
 578 *Figure 7. Spatial distribution of the best-performing formulas for SEC estimations across different*
 579 *HUC4 regions. At the HUC4 scale, the best-performing formula is the one(s) appears most frequently as*
 580 *the best among all catchments within the region.*

581
 582 9. In Section 5.1, the manuscript motivates the use of SR by contrasting it with the black
 583 box nature of RF. I suggest the authors provide a more precise definition of what
 584 constitutes interpretability in this context.

585 In hydrological literatures (Chen et al., 2022; Guillon et al., 2020), tree-based models are
 586 often considered interpretable via tools like SHAP values, which can quantify driver-
 587 response relationships much like the derivatives in SR. Given the significant drop in
 588 predictive skill (RF $R^2=0.80$ vs. SR $R^2=0.54$), it would be helpful to clarify why the explicit
 589 mathematical form of SR is preferred over the higher accuracy of RF. Is the interpretability
 590 here strictly defined as having a closed-form equation? A brief discussion about the
 591 definition of interpretability would provide a more balanced perspective.

592 **References**

593 Chadalawada, J., Herath, H.M.V.V., Babovic, V., 2020. Hydrologically Informed Machine
 594 Learning for Rainfall-Runoff Modeling: A Genetic Programming-Based Toolkit for

595 Automatic Model Induction. Water Resour. Res. 56, e2019WR026933.
596 <https://doi.org/10.1029/2019WR026933>

597 Chen, Y., Li, D., Zhao, Q., Cai, X., 2022. Developing a generic data-driven reservoir operation
598 model. Adv. Water Resour. 167, 104274.
599 <https://doi.org/10.1016/j.advwatres.2022.104274>

600 Guillon, H., Byrne, C.F., Lane, B.A., Sandoval Solis, S., Pasternack, G.B., 2020. Machine
601 Learning Predicts Reach-Scale Channel Types From Coarse-Scale Geospatial Data in a Large
602 River Basin. Water Resour. Res. 56, e2019WR026691.
603 <https://doi.org/10.1029/2019WR026691>

604 Thank you for this insightful comment. We explicitly defined interpretability as structural
605 transparency, meaning that the relationship between predictors and predictands is
606 expressed in an analytical equation. Although tree-based models can be interpreted using
607 post-hoc tools such as SHAP values and partial dependence plots, they only provide
608 approximate, model-dependent interpretations. Moreover, they typically describe average
609 effects while overlooking higher-order interactions and may be sensitive to feature
610 correlation, potentially leading to unstable or misleading explanations (Apley and Zhu,
611 2020; Sundararajan and Najmi, 2020). In contrast, the SR-derived equation provide explicit
612 expression for N , making the model transparent and analyzable (Häfner et al., 2023;
613 Karpatne et al., 2024). For example, differentiating F_3 gives the marginal effects of each
614 attribute on N : $\frac{\partial N}{\partial A} = a_3 d_3 (A + c_3)^{a_3-1}$, $\frac{\partial N}{\partial K_{sat}} = e_3 K_{sat}^{e_3-1}$, and $\frac{\partial N}{\partial f_{SWE}} = f_3$. These derivatives
615 make it possible to quantify parameter sensitivity directly and to clarify how geomorphic
616 and climatic factors jointly influence baseflow separation. Overall, RF and SR should be
617 viewed as complementary rather than competing approaches: RF provides a benchmark
618 for achievable predictive performance, whereas SR offers structural transparency that
619 facilitates theoretical interpretation and model integration.

620 We revised the introduction and discussion sections to clarify these points.

621 [Section 1, lines 78-86]

622 *Despite higher prediction accuracy, the RF-based regionalization model does not provide an explicit*
623 *analytical expression linking catchment attributes to parameter N (Rudin, 2019). Although tree-based*
624 *models can be interpreted using post-hoc tools such as SHAP values and partial dependence plots, they*
625 *provide only approximate, model-dependent interpretations rather than explicit functional*
626 *relationships (Makke and Chawla, 2024; Rudin, 2019). Moreover, they typically describe average*
627 *effects while overlooking higher-order interactions and may be sensitive to feature correlation,*
628 *potentially leading to unstable or misleading explanations (Apley and Zhu, 2020; Sundararajan and*
629 *Najmi, 2020).*

630 [Section 1, lines 107-111]

631 *This study should not be viewed as an effort to assert a superior utility of SR over other machine*
632 *learning models in the regionalization of baseflow parameters. Instead, the SR formulas serve as post-*
633 *hoc interpretability tools to complement other black box models, enhancing the transparency of the*
634 *underlying relationship between hydrological signatures and catchment attributes (Rudin, 2019).*

635 [Section 5.1, lines 471-494]

636 *In this study, we used SR to derive mathematical expressions for the predictions of N using 9*
637 *catchment attributes. Across ten cross-validation iterations, the identified expressions*
638 *exhibited consistent structures, predictors, and nearly identical regression coefficients, indicating that*
639 *SR can yield stable functional relationships between catchment attributes and N. Compared to the RF-*
640 *based predictions reported by Lin et al. (2026), the SR-based approach showed lower predictive skill*
641 *($R^2 = 0.54$ vs. 0.80), reflecting the trade-off between predictive accuracy and interpretability. While RF*
642 *achieves superior predictive performance, it functions as a ‘black-box’ ensemble, offering no explicit*
643 *functional form to clarify whether environmental controls operate additively, multiplicatively, or*
644 *through nonlinear transformations. In contrast, SR provides structural transparency by yielding a*
645 *closed-form equation, facilitating direct analytical insights (Häfner et al., 2023; Karpatne et al., 2024).*
646 *This explicit representation enables rigorous sensitivity analysis via differentiation; for instance, the*
647 *marginal effects derived from equation F_3 quantify how geomorphic and climatic factors jointly govern*
648 *N. By trading a degree of predictive skill for parsimony, SR transforms the problem from simple*
649 *estimation into a hypothesis-generating exercise, providing compact transfer functions that are easily*
650 *integrated into regionalization frameworks (Feigl et al., 2022; Samaniego et al., 2010). Therefore, RF*
651 *and SR should be viewed as complementary rather than competing approaches: RF provides a*
652 *benchmark for predictive performance, while SR offers structural transparency that facilitates*
653 *theoretical interpretation and model integration.*

654

655 References

- 656 Apley, D.W., Zhu, J., 2020. Visualizing the Effects of Predictor Variables in Black Box
657 Supervised Learning Models. *Journal of the Royal Statistical Society Series B: Statistical*
658 *Methodology*, 82(4): 1059-1086. DOI:10.1111/rssb.12377
- 659 Chadalawada, J., Herath, H.M.V.V., Babovic, V., 2020. Hydrologically Informed Machine
660 Learning for Rainfall-Runoff Modeling: A Genetic Programming-Based Toolkit for
661 Automatic Model Induction. *Water Resources Research*, 56(4): e2019WR026933.
662 DOI:<https://doi.org/10.1029/2019WR026933>
- 663 Diebold, F.X., Mariano, R.S., 1995. Comparing Predictive Accuracy. *Journal of Business &*
664 *Economic Statistics*, 13(3): 253-263. DOI:10.1080/07350015.1995.10524599
- 665 Feigl, M., Herrnegger, M., Klotz, D., Schulz, K., 2020. Function Space Optimization: A
666 Symbolic Regression Method for Estimating Parameter Transfer Functions for Hydrological
667 Models. *Water Resour Res*, 56(10): e2020WR027385. DOI:10.1029/2020WR027385
- 668 Feigl, M. et al., 2022. Automatic Regionalization of Model Parameters for Hydrological
669 Models. *Water Resources Research*, 58(12): e2022WR031966.
670 DOI:<https://doi.org/10.1029/2022WR031966>
- 671 Häfner, D., Gemmrich, J., Jochum, M., 2023. Machine-guided discovery of a real-world rogue
672 wave model. *Proceedings of the National Academy of Sciences*, 120(48): e2306275120.
673 DOI:[doi:10.1073/pnas.2306275120](https://doi.org/10.1073/pnas.2306275120)
- 674 Karpatne, A., Jia, X., Kumar, V., 2024. Knowledge-guided machine learning: Current trends
675 and future prospects. arXiv preprint arXiv:2403.15989.

676 Klotz, D., Herrnegger, M., Schulz, K., 2017. Symbolic Regression for the Estimation of
677 Transfer Functions of Hydrological Models. *Water Resources Research*, 53(11): 9402-9423.
678 DOI:10.1002/2017wr021253

679 Kronberger, G., de Franca, F.O., Burlacu, B., Haider, C., Kommenda, M., 2022. Shape-
680 Constrained Symbolic Regression—Improving Extrapolation with Prior Knowledge.
681 *Evolutionary Computation*, 30(1): 75-98. DOI:10.1162/evco_a_00294

682 Li, Q. et al., 2024. Advancing symbolic regression for earth science with a focus on
683 evapotranspiration modeling. *npj Climate and Atmospheric Science*, 7(1).
684 DOI:10.1038/s41612-024-00861-5

685 Mei, Y. et al., 2024. Optimal Baseflow Separation Through Chemical Mass Balance:
686 Comparing the Usages of Two Tracers, Two Concentration Estimation Methods, and Four
687 Baseflow Filters. *Water Resources Research*, 60(7). DOI:10.1029/2023wr036386

688 Sheta, A., Abdel-raouf, A., Fraihat, K., Baareh, A.K., 2023. Evolutionary Design of a PSO-
689 Tuned Multigene Symbolic Regression Genetic Programming Model for River Flow
690 Forecasting. *International Journal of Advanced Computer Science and Applications*, 14:
691 2023. DOI:10.14569/IJACSA.2023.0140489

692 Sundararajan, M., Najmi, A., 2020. The Many Shapley Values for Model Explanation. In: Hal,
693 D., III, Aarti, S. (Eds.), *Proceedings of the 37th International Conference on Machine
694 Learning*. PMLR, *Proceedings of Machine Learning Research*, pp. 9269--9278.

695 Wilstrup, C., Kasak, J., 2021. Symbolic regression outperforms other models for small data
696 sets. arXiv preprint arXiv:2103.15147.

697

698

Reply to Reviewer #3

699 This manuscript presents an innovative approach to the regionalization of the baseflow
700 separation parameter N by bridging the gap between "black-box" machine learning and
701 traditional physical hydrology. By utilizing Symbolic Regression (SR) across 855
702 catchments in the US, the authors successfully identified explicit and interpretable
703 mathematical formulas that incorporate key catchment attributes such as drainage area,
704 soil conductivity, and snow processes. The study is well motivated, and the use of SEC for
705 validation adds significant scientific rigor. While the work is technically sound and highly
706 relevant to the community, there are several conceptual and methodological points that
707 require clarification and revision before publication.

708 Thank you for this positive and encouraging assessment of our work. We appreciate the
709 reviewer's recognition of the novelty of using symbolic regression to bridge data-driven
710 approaches and physical interpretability in hydrology, as well as the value of deriving
711 explicit mathematical relationships between the baseflow separation parameter and key
712 catchment attributes. We also thank the reviewer for the constructive comments. These
713 points have been carefully considered, and the manuscript has been revised accordingly.
714 Detailed responses to each comment are provided below.

715

716 1. A major concern is the mathematical structure of F2 and F3, where catchment area (A) is
717 directly added to saturated hydraulic conductivity (K_{sat}). Please clarify the physical or
718 mathematical justification for adding variables with different dimensions and how this
719 affects the transferability of the formulas across different unit systems

720 Thank you for this important comment. In the previous SR model training, the $\times 100$
721 scaling factor in K_{sat} was preserved, causing this dimensional inconsistency. To address
722 this issue, we retrained the SR models using the correct K_{sat} magnitude, together with
723 some additional structural constraints to improve the physical interpretability of the
724 resulting expressions. Specifically, multiplication, division, power-law, and logarithmic
725 operators were not allowed to be nested within operators of the same type. The internal
726 complexity of expressions within power-law and logarithmic operators was limited to a
727 maximum value of 3, and the overall expression complexity was capped at 20. In addition,
728 recursive formulations (i.e., expressions where the output variable appears as an input to
729 itself) were not permitted to ensure model interpretability and avoid trivial or ill-posed
730 solutions. The new SR expressions are consisted of single variable terms only, which are
731 more physically interpretable and avoid dimensionally inconsistent combinations.

732 [Section 3.2, Lines 191-206]

733 *To control the search space and ensure physically interpretable expressions, several structural*
734 *constraints were imposed in SR model training. Multiplication, division, power-law, and*
735 *logarithmic operators were not allowed to be nested within operators of the same type. The*
736 *internal complexity of expressions inside power-law and logarithmic operators was restricted*
737 *to a maximum value of 3. The maximum allowable total complexity was set to 20. Expression*
738 *complexity is defined as the sum of the complexity index assigned to each component in the*
739 *equation. Take $N = 1.6 \times A^{0.2}$ as an example, if multiplication and power-law operators are*

740 *each assigned a complexity of 2 and constants and input variables are assigned a complexity*
741 *of 1, the total complexity of the expression is calculated as $2 + 2 + 1 + 1 + 1 = 7$. In this*
742 *study, all operators were assigned a uniform complexity index of 1 to avoid bias toward*
743 *specific functional forms. Recursive formulations (i.e., expressions where the output variable*
744 *appears as an input to itself) were not permitted to ensure model interpretability and avoid*
745 *trivial or ill-posed solutions.*

746 [Table 2]

747 **Table 2. The calibrated F'_{PL} and the three stable forms of SR expressions across all ten iterations of the cross-validation. The**
 748 **numbers within the brackets of the first row are the complexity indices of the expressions.**

Iteration	F'_{LR} (5)	F_1 (5)	F_2 (13)	F_3 (17)
1	$3.00 \times A^{0.15}$	$A^{0.23} + 3.66$	$0.40(A + 529)^{0.31} + K_{sat}^{0.30} + 2.39$	$0.27(A + 640)^{0.33} + K_{sat}^{0.35} + 3.71f_{SWE} + 1.79$
2	$3.10 \times A^{0.15}$	$A^{0.23} + 3.76$	$0.62(A + 642)^{0.28} + K_{sat}^{0.32} + 1.37$	$0.40(A + 652)^{0.30} + K_{sat}^{0.36} + 3.38f_{SWE} + 1.26$
3	$3.10 \times A^{0.15}$	$A^{0.23} + 3.74$	$0.47(A + 526)^{0.30} + K_{sat}^{0.28} + 2.19$	$0.32(A + 575)^{0.32} + K_{sat}^{0.33} + 3.55f_{SWE} + 1.80$
4	$3.25 \times A^{0.14}$	$A^{0.23} + 3.86$	$0.39(A + 334)^{0.30} + K_{sat}^{0.31} + 2.78$	$0.26(A + 361)^{0.33} + K_{sat}^{0.35} + 3.50f_{SWE} + 2.32$
5	$3.17 \times A^{0.15}$	$A^{0.23} + 3.79$	$0.39(A + 371)^{0.31} + K_{sat}^{0.31} + 2.64$	$0.35(A + 547)^{0.31} + K_{sat}^{0.35} + 3.60f_{SWE} + 1.61$
6	$3.03 \times A^{0.15}$	$A^{0.23} + 3.76$	$0.56(A + 602)^{0.28} + K_{sat}^{0.29} + 1.77$	$0.41(A + 754)^{0.30} + K_{sat}^{0.34} + 3.63f_{SWE} + 1.21$
7	$3.14 \times A^{0.15}$	$A^{0.23} + 3.83$	$0.44(A + 401)^{0.30} + K_{sat}^{0.29} + 2.58$	$0.30(A + 443)^{0.32} + K_{sat}^{0.33} + 3.28f_{SWE} + 2.18$
8	$3.16 \times A^{0.15}$	$A^{0.23} + 3.85$	$0.44(A + 421)^{0.30} + K_{sat}^{0.31} + 2.49$	$0.31(A + 491)^{0.32} + K_{sat}^{0.35} + 3.54f_{SWE} + 1.95$
9	$3.13 \times A^{0.15}$	$A^{0.23} + 3.81$	$0.46(A + 481)^{0.30} + K_{sat}^{0.28} + 2.42$	$0.30(A + 525)^{0.32} + K_{sat}^{0.33} + 3.38f_{SWE} + 2.06$
10	$3.14 \times A^{0.15}$	$A^{0.23} + 3.79$	$0.43(A + 545)^{0.30} + K_{sat}^{0.31} + 2.29$	$0.29(A + 627)^{0.32} + K_{sat}^{0.35} + 3.57f_{SWE} + 1.86$

749

750 2. Why was a standard Genetic Programming-based SR approach chosen over more recent
751 continuous search methods or grammar-based optimization, such as those discussed by
752 Feigl et al. (2020)?

753 Thank you for this insightful comment. In our study, the search space was deliberately
754 restricted in both dimensionality and structural complexity to ensure tractability and
755 robustness. Prior to symbolic regression, we applied a mutual information-based feature
756 selection procedure, which reduced the candidate predictors to nine physically relevant
757 variables. In addition, the set of mathematical operators was explicitly constrained to a
758 small and interpretable subset, including addition, subtraction, multiplication, division,
759 power-law, and logarithmic functions. To further control structural complexity, we
760 imposed additional restrictions: a) multiplication, division, power-law, and logarithmic
761 operators were not allowed to be nested within operators of the same type; b) the internal
762 complexity of expressions within power-law and logarithmic operators was limited to a
763 maximum of 3; and c) the total expression complexity was capped at 20. These constraints
764 substantially reduce the search space and mitigate the known limitations of GP in handling
765 highly complex symbolic structures. Under these controlled conditions, the advantages of
766 more advanced approaches (e.g., their ability to efficiently explore high-dimensional and
767 flexible grammars) are less critical. Instead, our objective is not to exhaustively search for
768 highly complex symbolic forms, but to identify low-complexity, physically interpretable,
769 and structurally stable relationships that are consistent across cross-validation runs. For
770 this purpose, GP-based symbolic regression provides a good balance between search
771 efficiency, interpretability, and implementation simplicity, and remains well aligned with
772 our study goals. Similar GP-based SR approaches have also been successfully applied in
773 other studies (Dorgo et al., 2021; Pumo and Noto, 2023; Razaq et al., 2016).

774 To further assess whether the derived relationships depend on the choice of optimization
775 paradigm, we additionally implemented an alternative symbolic regression framework
776 based on deep learning with reinforcement learning-guided search. The resulting
777 expressions are: $A^{0.22} + 3.70$, $(A + 901)^{0.24} + K_{sat}^{0.28}$, $0.31 * (A + 549)^{0.31} + K_{sat}^{0.34} +$
778 $3.15f_{SWE} + 1.84$, with corresponding R^2 values of 0.47, 0.52, and 0.55, respectively. These
779 expressions exhibit highly similar functional forms and comparable predictive
780 performance to those obtained using the GP-based approach. This consistency suggests
781 that the identified relationships are robust across different symbolic regression
782 optimization paradigms.

783

784 3. There is a mismatch in the units of K_{sat} between Table 1 (10-2 mm/h) and Figure 3
785 (cm/day). Please unify the units throughout the manuscript and figures for consistency.

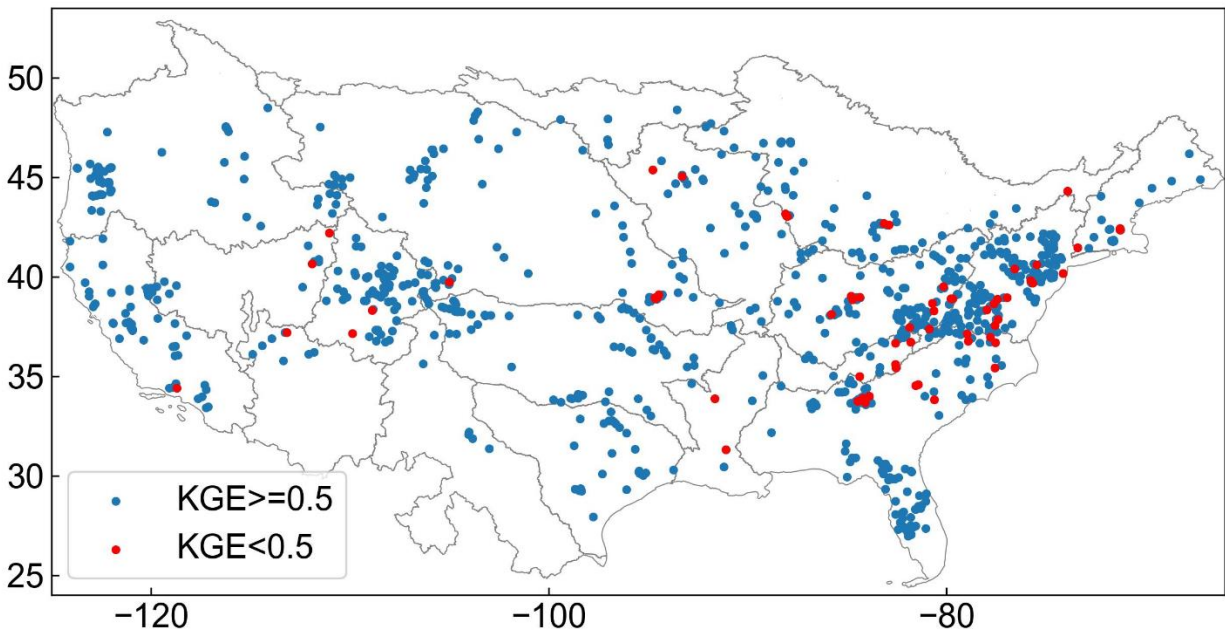
786 Thank you for pointing this out. The units were unified to mm/hr throughout the
787 manuscript.

788

789 4. Catchments with a KGE below 0.5 were excluded from the analysis. Please clarify if these
790 excluded catchments belong to a specific hydrologic regime (e.g., arid or ephemeral

791 streams), and discuss the potential limitations this imposes on the global applicability of
792 your formulas.

793 Thank you for this insightful comment. Mei et al. (2024) shows that the performance of
794 baseflow separation had no explicit spatial pattern. Therefore, the excluded gages with
795 $KGE < 0.5$ are unlikely from specific hydrologic regimes. We also examined the spatial
796 distribution of the excluded gages (63 in total) by plotting them with the others. The figure
797 reveal no explicit spatial clusters, which stands as a proof of our claim.



798 Spatial distribution of the 855 selected catchments with Kling Gupta Efficiency between estimated and
799 observed SEC below and above 0.5, respectively.
800

801

802 5. You report an R^2 of 0.54 for SR versus 0.80 for RF. This is a substantial loss in predictive
803 power. I request a more detailed discussion on whether the gain in interpretability justifies
804 such a high "accuracy cost," particularly for practical water resource management
805 applications.

806 Thank you for raising this important point. We clarify that the choice of SR is not intended
807 to replace high-accuracy models such as RF, but to complement them by providing
808 structural interpretability. While RF achieves higher accuracy, its ensemble structure does
809 not provide an explicit functional form linking predictors to N , such as whether controls
810 operate additively, multiplicatively, or through specific nonlinear transformations. In
811 contrast, the principal advantage of SR lies in its ability to produce an explicit closed-form
812 equation. This feature is particularly valuable in hydrology, where the goal is often not only
813 to predict a quantity of interest, but also to understand how catchment characteristics
814 jointly control that quantity. For instance, differentiating F_3 gives the marginal effects of
815 each attribute on N : $\frac{\partial N}{\partial A} = a_3 d_3 (A + c_3)^{a_3 - 1}$, $\frac{\partial N}{\partial K_{sat}} = e_3 K_{sat}^{e_3 - 1}$, and $\frac{\partial N}{\partial f_{SWE}} = f_3$. These
816 derivatives make it possible to quantify parameter sensitivity directly and to clarify how
817 geomorphic and climatic factors jointly influence baseflow separation. For practical water

818 resource management, understanding systematic controls and mechanistic relationships
819 can be as valuable as incremental improvements in predictive accuracy. An additional
820 advantage of SR is parsimony. The final model is a compact equation rather than a large
821 ensemble of trees, which makes it easier to communicate and potentially integrate into
822 regionalization frameworks. Therefore, RF and SR should be viewed as complementary
823 rather than competing approaches: RF provides a benchmark for achievable predictive
824 performance, whereas SR offers structural transparency that facilitates theoretical
825 interpretation and model integration.

826 We revised the introduction and discussion sections to clarify this point and to emphasize
827 that RF and SR serve complementary roles.

828 [Section 1, lines 107-111]

829 *This study should not be viewed as an effort to assert a superior utility of SR over other machine*
830 *learning models in the regionalization of baseflow parameters. Instead, the SR formulas serve as post-*
831 *hoc interpretability tools to complement other black box models, enhancing the transparency of the*
832 *underlying relationship between hydrological signatures and catchment attributes (Rudin, 2019).*

833 [Section 5.1, lines 471-494]

834 *In this study, we used SR to derive mathematical expressions for the predictions of N using 9*
835 *catchment attributes. Across ten cross-validation iterations, the identified expressions*
836 *exhibited consistent structures, predictors, and nearly identical regression coefficients, indicating that*
837 *SR can yield stable functional relationships between catchment attributes and N . Compared to the RF-*
838 *based predictions reported by Lin et al. (2026), the SR-based approach showed lower predictive skill*
839 *($R^2 = 0.54$ vs. 0.80), reflecting the trade-off between predictive accuracy and interpretability. While RF*
840 *achieves superior predictive performance, it functions as a 'black-box' ensemble, offering no explicit*
841 *functional form to clarify whether environmental controls operate additively, multiplicatively, or*
842 *through nonlinear transformations. In contrast, SR provides structural transparency by yielding a*
843 *closed-form equation, facilitating direct analytical insights (Häfner et al., 2023; Karpatne et al., 2024).*
844 *This explicit representation enables rigorous sensitivity analysis via differentiation; for instance, the*
845 *marginal effects derived from equation F_3 quantify how geomorphic and climatic factors jointly govern*
846 *N . By trading a degree of predictive skill for parsimony, SR transforms the problem from simple*
847 *estimation into a hypothesis-generating exercise, providing compact transfer functions that are easily*
848 *integrated into regionalization frameworks (Feigl et al., 2022; Samaniego et al., 2010). Therefore, RF*
849 *and SR should be viewed as complementary rather than competing approaches: RF provides a*
850 *benchmark for predictive performance, while SR offers structural transparency that facilitates*
851 *theoretical interpretation and model integration.*

852

853 6. In Section 5.2, the formulas indicate that N (flow event duration) increases with K_{sat} .
854 Physically, higher soil conductivity often implies faster drainage. While you attribute this to
855 a shift toward slower subsurface paths, please provide more quantitative evidence or
856 literature support to explain why this "slowing" effect dominates over the expected
857 increase in drainage efficiency.

858 Thank you for this insightful comment. It is important to clarify that the parameter N
859 represents a catchment-scale flow duration, reflecting the integrated hydrological response
860 rather than local-scale drainage velocity. Its relationship with K_{sat} therefore emerges from

861 the combined effects of runoff generation, flow partitioning, and subsurface transport
862 processes. In highly permeable catchments, enhanced infiltration reduces rapid overland
863 flow generation and increases the contribution of subsurface flow pathways. This shift in
864 runoff partitioning can modify the hydrograph shape and potentially extend the effective
865 flow duration.

866 To further examine this mechanism, we analyzed the statistical relationships between K_{sat} ,
867 the recession coefficient, and N across all study catchments. The Spearman correlation
868 between K_{sat} and the recession coefficient is positive ($r = 0.17$, $p < 0.01$), indicating that
869 catchments with higher K_{sat} tend to exhibit slower recession behavior (i.e., longer recession
870 times). In addition, the correlation between K_{sat} and N is $r = 0.24$ ($p < 0.01$), consistent with
871 the positive dependency identified in the derived formula. These correlations support the
872 interpretation that higher K_{sat} is associated with hydrological responses characterized by
873 relatively longer effective flow duration. This interpretation is also consistent with
874 previous studies showing that highly permeable geological settings can promote prolonged
875 recession behavior (Carlier et al., 2018; Krakauer and Temimi, 2011; Tague and Grant,
876 2004). Therefore, while higher K_{sat} enhances local drainage efficiency at the soil profile
877 scale, its integrated effect at the catchment scale appears to be mediated through runoff
878 partitioning and increased subsurface contributions, which can outweigh the expected
879 acceleration of depletion and result in longer effective flow durations.

880

881 7. Line 230, Performance is notably poor in HUC 12. Given that your input data includes
882 reservoir storage, please clarify if the SMM method's fundamental assumption of "natural"
883 baseflow is even applicable in these highly regulated and irrigated basins.

884 Thank you for this important comment. The relatively poor performance in HUC 12 may be
885 related to its distinct hydrological characteristics. Many catchments in this region are
886 relatively arid and exhibit weak storage effects, flashy hydrographs, and strong sensitivity
887 to short-duration, high-intensity rainfall events (Feng et al., 2020; Kratzert et al., 2019).
888 These fast hydrological responses may not be captured by SMM, which often reveals
889 smooth baseflow dynamics (Stewart, 2015). Similar challenges in accurately predictions of
890 baseflow and streamflow were also reported by other studies for the region (Feng et al.,
891 2020; Xie et al., 2022). We added this discussion to Section 4.3.

892 [Section 4.3, lines 367-371]

893 *In contrast, the lowest performance for all three formulas occurs in HUC 12, where median KGE*
894 *values fall below 0.65 and more than 25% of catchments show KGE values below 0.5. This may be*
895 *related to the relative arid climate and flashy hydrological response of HUC 12 (Feng et al., 2020;*
896 *Kratzert et al., 2019), which is difficult for SMM to capture. Note that SMM is more skillful for smooth*
897 *baseflow dynamics (Stewart, 2015).*

898 Regarding the concern about human regulation, we examined the relationships between
899 baseflow separation performance and two indicators of anthropogenic disturbance, namely
900 reservoir storage and the areal fraction of irrigation, across the 855 study catchments. The
901 correlation coefficients between R^2 of baseflow and the two attributes are both smaller
902 than 0.05 ($p > 0.1$), indicating no statistical significance. This suggests that the baseflow

903 separation performance of SMM does not appear to related to the reservoir storage and the
904 areal fraction of irrigation factors.

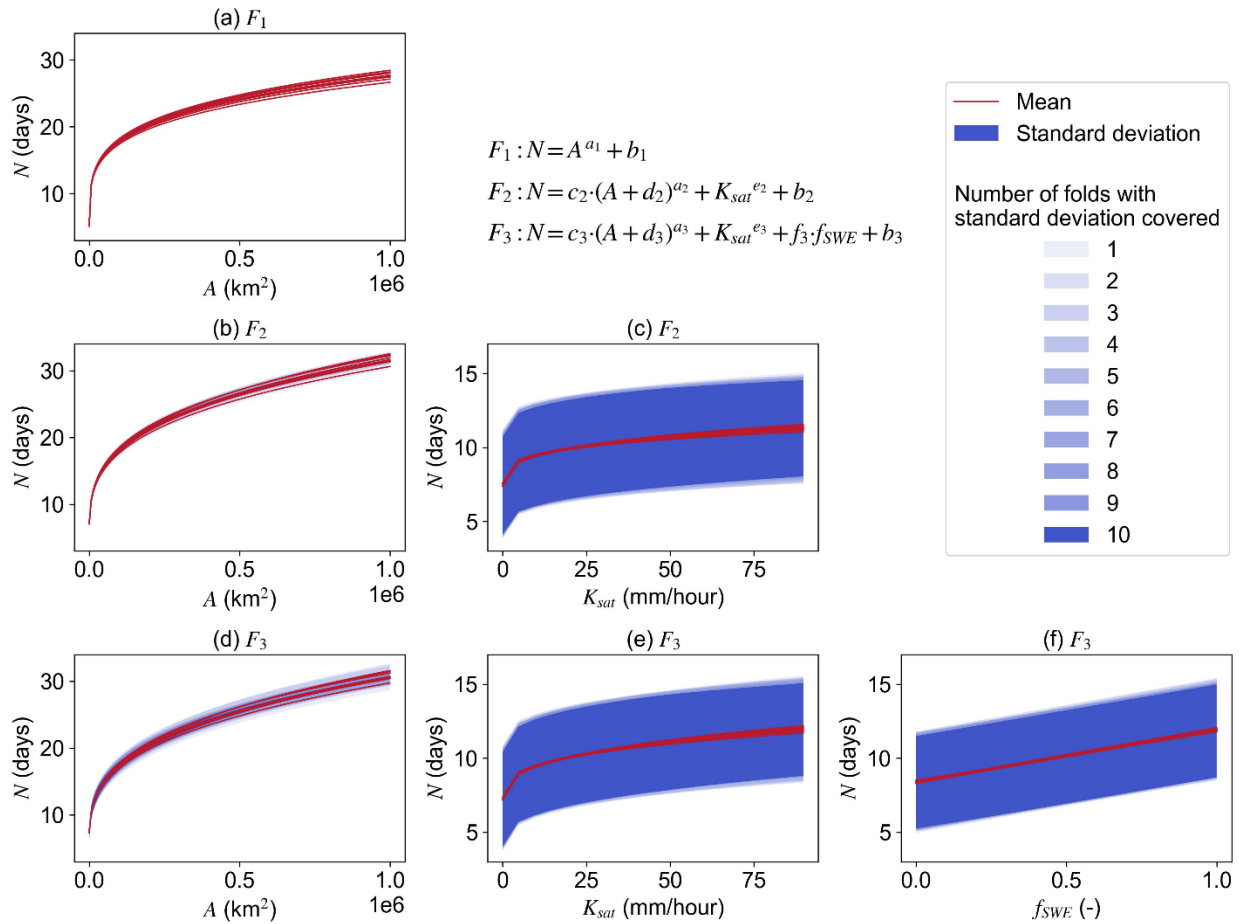
905

906 8. In Figure 3, the ten individual lines representing different cross-validation folds are
907 extremely difficult to distinguish. Please improve the visualization to more clearly show the
908 variance between folds.

909 Thank you for this comment. The individual lines represent the SR-derived formulas
910 obtained from the ten cross-validation folds. They appear closely clustered because the
911 corresponding coefficients are highly similar across folds. This behavior reflects the
912 robustness and stability of the SR-derived relationships, rather than a limitation of the
913 visualization. Therefore, we decide to retain the current presentation of this figure.
914 Additionally, we have added clarifying descriptions in the revised manuscript to better
915 explain this point.

916 [Section 4.1, lines 287-299]

917 *Figure 3 further illustrates the behavior of these formulas. For F_1 , the nearly identical exponents*
918 *(~ 0.23) and intercepts (3.66–3.86 days) result in almost overlapping curves (Figure 3a), indicating a*
919 *stable power-law relationship between N and A , with a diminishing rate of increase. For F_2 , the*
920 *similarly constrained exponents (0.28–0.32) and intercepts (1.37–2.78 days) produce tightly clustered*
921 *response curves (Figure 3b-c), showing that both A and K_{sat} contribute positively to N . F_3 extends F_2*
922 *by introducing f_{SWE} as a linear term. The replicated formulas still exhibit closely grouped slopes*
923 *(3.28–3.71 days) and intercepts (1.21–2.32 days), which explains the clustering of curves in Figure 3d–*
924 *f. The marginal relationships of N with A and K_{sat} in F_3 remain consistent with those in F_2 , whereas*
925 *increasing f_{SWE} leads to an approximately linear increase in N , at a rate of about 0.3–0.4 days per 0.1*
926 *increment in snow fraction. Overall, SR identifies A as the most influential factor in predicting N , as*
927 *evidenced by its presence in all SR-derived formulas. The narrower ranges of predicted N s in Figure*
928 *3b and d also suggest that A exerts greater influence than K_{sat} and f_{SWE} .*



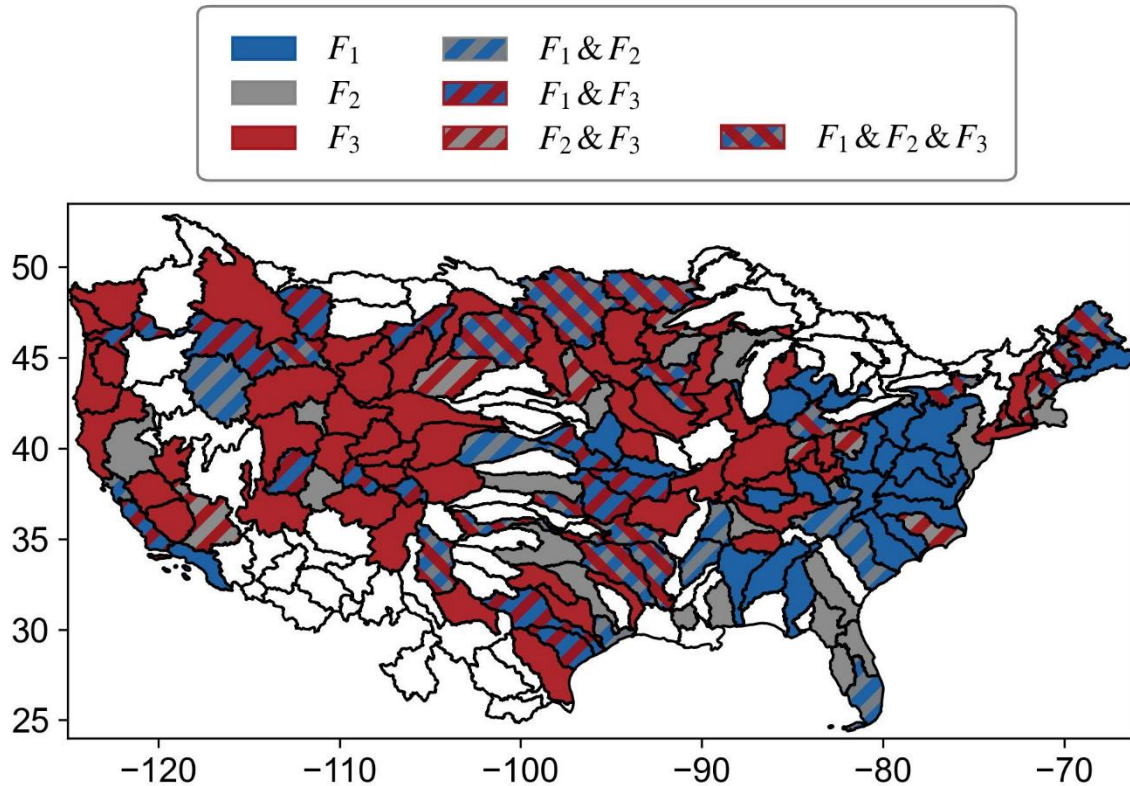
929

930 **Figure 3. Marginal relationship of N on different predictors (A , K_{sat} , and f_{SWE}) that consist of the SR**
 931 **expressions (F_1 , F_2 , and F_3). Each line represents one of the ten instances of F_1 , F_2 , and F_3 . Panels a,**
 932 **b, and d are for A ; panels c and e are for K_{sat} ; panel f is for f_{SWE} .**

933

934 **9. The use of multiple hatching patterns and overlapping colors in Figure 7b results in a**
 935 **cluttered visual presentation that is difficult for the reader to interpret.**

936 Thank you for pointing this out. The combined color–pattern encoding was intentionally
 937 designed to preserve compositional information, allowing readers to immediately identify
 938 which formula(s) participate in the best-performing category, rather than only indicating a
 939 single dominant formula. We explored alternative visualization schemes, including a
 940 triangular color spectrum to represent combinations of formulas. However, these
 941 approaches reduced interpretability, as it became difficult to visually trace the contribution
 942 of individual formulas across regions and to distinguish shared dominance (i.e., ties) from
 943 single-formula dominance. To improve readability, we have refined the figure in the
 944 revised manuscript by adjusting hatching density and enlarging legend. In addition, a
 945 higher-resolution version will be provided upon publication to allow readers to better
 946 examine regional details.



947

948 *Figure 7. Spatial distribution of the best-performing formulas for SEC estimations across different*
 949 *HUC4 regions. At the HUC4 scale, the best-performing formula is the one(s) appears most frequently as*
 950 *the best among all catchments within the region.*

951

952 References:

953 Carlier, C., Wirth, S.B., Cochand, F., Hunkeler, D., Brunner, P., 2018. Geology controls
 954 streamflow dynamics. *Journal of Hydrology*, 566: 756-769.

955 Dorgo, G., Kulcsar, T., Abonyi, J., 2021. Genetic programming-based symbolic regression for
 956 goal-oriented dimension reduction. *Chemical Engineering Science*, 244: 116769.
 957 DOI:<https://doi.org/10.1016/j.ces.2021.116769>

958 Feigl, M. et al., 2022. Automatic Regionalization of Model Parameters for Hydrological
 959 Models. *Water Resources Research*, 58(12): e2022WR031966.
 960 DOI:<https://doi.org/10.1029/2022WR031966>

961 Feng, D., Fang, K., Shen, C., 2020. Enhancing Streamflow Forecast and Extracting Insights
 962 Using Long - Short Term Memory Networks With Data Integration at Continental Scales.
 963 *Water Resources Research*, 56(9): e2019WR026793. DOI:10.1029/2019wr026793

964 Häfner, D., Gemrich, J., Jochum, M., 2023. Machine-guided discovery of a real-world rogue
 965 wave model. *Proceedings of the National Academy of Sciences*, 120(48): e2306275120.
 966 DOI:[doi:10.1073/pnas.2306275120](https://doi.org/10.1073/pnas.2306275120)

967 Karpatne, A., Jia, X., Kumar, V., 2024. Knowledge-guided machine learning: Current trends
968 and future prospects. arXiv preprint arXiv:2403.15989.

969 Krakauer, N.Y., Temimi, M., 2011. Stream recession curves and storage variability in small
970 watersheds. *Hydrology and Earth System Sciences*, 15(7): 2377-2389.

971 Kratzert, F. et al., 2019. Towards learning universal, regional, and local hydrological
972 behaviors via machine learning applied to large-sample datasets. *Hydrology and Earth
973 System Sciences*, 23(12): 5089-5110. DOI:10.5194/hess-23-5089-2019

974 Lin, Y. et al., 2026. Regionalization of Optimal Baseflow Separation using Catchment-scale
975 Characteristics. *Water Resources Research*, 62(3).

976 Pumo, D., Noto, L.V., 2023. Exploring the use of multi-gene genetic programming in regional
977 models for the simulation of monthly river runoff series. *Stochastic Environmental
978 Research and Risk Assessment*, 37(5): 1917-1941. DOI:10.1007/s00477-022-02373-1

979 Razaq, S.A. et al., 2016. Prediction of Flow Duration Curve in Ungauged Catchments Using
980 Genetic Expression Programming. *Procedia Engineering*, 154: 1431-1438.
981 DOI:<https://doi.org/10.1016/j.proeng.2016.07.516>

982 Rudin, C., 2019. Stop Explaining Black Box Machine Learning Models for High Stakes
983 Decisions and Use Interpretable Models Instead. *Nat Mach Intell*, 1(5): 206-215.
984 DOI:10.1038/s42256-019-0048-x

985 Samaniego, L., Kumar, R., Attinger, S., 2010. Multiscale parameter regionalization of a grid-
986 based hydrologic model at the mesoscale. *Water Resources Research*, 46(5).
987 DOI:<https://doi.org/10.1029/2008WR007327>

988 Stewart, M.K., 2015. Promising new baseflow separation and recession analysis methods
989 applied to streamflow at Glendhu Catchment, New Zealand. *Hydrol. Earth Syst. Sci.*, 19(6):
990 2587-2603. DOI:10.5194/hess-19-2587-2015

991 Tague, C.L., Grant, G.E., 2004. A geological framework for interpreting the low-flow regimes
992 of Cascade streams, Willamette River Basin, Oregon. *Water Resources Research*, 40.

993 Xie, J. et al., 2022. Estimating Gridded Monthly Baseflow From 1981 to 2020 for the
994 Contiguous US Using Long Short - Term Memory (LSTM) Networks. *Water Resources
995 Research*, 58(8). DOI:10.1029/2021wr031663

996