

# 1 **Technical Note – RAT: a Robustness Assessment Test for calibrated** 2 **and uncalibrated hydrological models**

3  
4 Pierre Nicolle<sup>1,3</sup>, Vazken Andréassian<sup>1\*</sup>, Paul Royer-Gaspard<sup>1</sup>, Charles Perrin<sup>1</sup>, Guillaume Thirel<sup>1</sup>,  
5 Laurent Coron<sup>2</sup>, Léonard Santos<sup>1</sup>

6 <sup>1</sup>Université Paris-Saclay, INRAE, UR HYCAR, 92160, Antony, France

7 <sup>2</sup>EDF, DTG, Toulouse, France

8 <sup>3</sup>now at Université Gustave Eiffel, Nantes, France

9 \*Corresponding author: Vazken Andréassian ([vazken.andreassian@inrae.fr](mailto:vazken.andreassian@inrae.fr))

## 10 **Key Words**

11 hydrological modelling, split-sample test, differential split-sample test, model evaluation, robustness,  
12 climate change

## 13 **Key Points**

- 14 • a new method (RAT) is proposed to assess the robustness of hydrological models, as an  
15 alternative to the classical split-sample test
- 16 • the RAT method does not require multiple calibrations [of hydrological models](#): it is therefore  
17 applicable to uncalibrated models
- 18 • the RAT method can be used to determine whether a hydrological model cannot be safely  
19 used for climate change impact studies
- 20 • success at the RAT test is a necessary (but not sufficient) condition of model robustness

## 21 **Abstract**

22 Prior to their use under future changing climate conditions, all hydrological models should be  
23 thoroughly evaluated regarding their temporal transferability (application in different time periods)  
24 and extrapolation capacity (application beyond the range of known past conditions). This note  
25 presents a straightforward evaluation framework aimed at detecting potential undesirable climate  
26 dependencies in hydrological models: the robustness assessment test (RAT). Although it is  
27 conceptually inspired by the classic differential split-sample test of Klemeš (1986), the RAT presents  
28 the advantage to be applicable to all types of models, be they calibrated or not (i.e. regionalized or  
29 physically based). In this note, we present the RAT, illustrate its application on a set of 21  
30 catchments, verify its applicability hypotheses and compare it to previously published tests. Results  
31 show that the RAT is an efficient evaluation approach, passing it successfully can be considered a  
32 prerequisite for any hydrological model to be used for climate change impact studies.

## 33 1 Introduction

### 34 1.1 All hydrological models should be evaluated for their robustness

35 Hydrologists are increasingly requested to provide predictions of the impact of climate change  
36 (Wilby, 2019). Given the expected evolution of climate conditions, the actual ability of models to  
37 predict the corresponding evolution of hydrological variables should be verified (Beven, 2016).  
38 Indeed, when using a hydrological model for climate change impact assessment, we make two  
39 implicit hypotheses concerning:

40 • the **capacity of extrapolation beyond known hydroclimatic conditions**: we assume that the  
41 hydrological model used is able to extrapolate catchment behaviour under conditions not or rarely  
42 seen in the past. While we do not expect hydrological models to be able to simulate a behaviour  
43 which would result from a modification of catchment physical characteristics, we do expect them to  
44 be able to represent the catchment response to extreme climatic conditions (and possibly to  
45 conditions more extreme than those observed in the past);

46 • the **independence of the model set-up period**: we assume that the model functioning is  
47 independent of the climate it experienced during its set-up/calibration period. For those models  
48 which are calibrated, we assume that the parameters are generic and not specific to the calibration  
49 period, i.e. they do not suffer from overcalibration on this period (Andréassian et al., 2012).

50 Hydrologists make the hypothesis that model structure and parameters are well-identified over the  
51 calibration period and that parameters remain relevant over the future period, when climate  
52 conditions will be different. Unfortunately, the majority of hydrological models are not entirely  
53 independent of climate conditions (Refsgaard et al., 2013; Thirel et al., 2015b). When run under  
54 changing climate conditions, they sometimes reveal an unwanted sensitivity to the data used to  
55 conceive or calibrate them (Coron et al., 2011).

56 The diagnostic tool most widely used to assess the robustness of hydrological models is the split-  
57 sample test (SST) (Klemeš, 1986), which is considered by ~~all~~ most hydrologists as a “good modelling  
58 practice” (Refsgaard & Henriksen, 2004). The SST stipulates that when a model requires calibration  
59 (i.e. when its parameters cannot be deduced directly from physical measurements or catchment  
60 descriptors), it should be evaluated twice: once on the data used for calibration and once on an  
61 independent dataset. This practice has been promoted in hydrology by Klemeš (1986), who did not  
62 invent the concept (Arlot & Celisse, 2010; Larson, 1931; Mosteller & Tukey, 1968), but who  
63 formalized it for hydrological modelling. Klemeš proposed initially a four-level testing scheme for  
64 evaluating model transposability in time and space: (i) split-sample test on two independent periods,  
65 (ii) proxy-basin test on neighbouring catchments, (iii) differential split-sample test on contrasted  
66 independent periods (DSST), and (iv) proxy-basin differential split-sample test on neighbouring  
67 catchments and contrasted periods.

68 For model applications in a changing climate context, Klemeš’s DSST procedure is of particular  
69 interest. Indeed, when calibration and evaluation are done over climatically-contrasted past periods,  
70 the model faces the difficulties it will have to deal with in the future. The power of DSST can be  
71 limited by the climatic variability observed in the past, which may be far below the drastic changes

72 expected in the future. However, a satisfactory behaviour during the DSST can be seen as a  
73 prerequisite of model robustness.

## 74 1.2 Past applications of the DSST method

75 The DSST received limited attention up to the 2010s, with only a few studies which applied it. The  
76 studies by Refsgaard & Knudsen (1996) and Donnelly-Makowecki & Moore (1999) investigated to  
77 which extent Klemeš's hierarchical testing scheme could be used to improve the conclusions of  
78 model intercomparisons. ~~Though the authors of the first study did not find large differences between  
79 the SST and DSST when comparing conceptual and physically-oriented models, the authors of the  
80 second study found that the DSST was more powerful than the SST to discriminate between four  
81 event-based models.~~The study by Xu (1999) questioned the applicability of models in nonstationary  
82 conditions and was one of the early attempts to apply the Klemeš's testing scheme in this  
83 perspective. Similarly, tests carried out by Seibert (2003) explicitly intended to test the ability of a  
84 model to extrapolate beyond calibration range and showed limitations of the tested model, stressing  
85 the need for improved calibration strategies. Last, Vaze et al. (2010) also investigated the behaviour  
86 of four rainfall-runoff models under contrasting conditions, using wet and dry periods on catchments  
87 in Australia that experienced a prolonged drought period. They observed different model behaviours  
88 when going from wet to dry or dry to wet conditions.

89 More recently, Coron et al. (2012) proposed a generalized SST (GSST) allowing for an exhaustive DSST  
90 to evaluate model transposability over time under various climate conditions. The concept of GSST  
91 consists in testing "the model in as many and as varied climatic configurations as possible, including  
92 similar and contrasted conditions between calibration and validation. [...] ~~The GSST procedure simply  
93 consists of a series of calibration-validation tests on subperiods of equal length, considering all  
94 possible configurations~~". Seifert et al. (2012) used a differential split-sample approach to test a  
95 hydrogeological model (differential being understood with respect to differences in groundwater  
96 abstractions). Li et al. (2012) identified two dry and two wet periods in long hydroclimatic series to  
97 understand how a model should be parameterised to work under nonstationary climatic conditions.  
98 Teutschbein and Seibert (2013) performed differential split-sample tests by dividing the data series  
99 into cold and warm as well as dry and wet years, in order to evaluate bias correction methods. Thirel  
100 et al. (2015a) put forward an SST-based protocol to investigate how hydrological models deal with  
101 changing conditions, which was widely used during an [IAHS-workshop of the International  
102 Association of Hydrological Sciences \(IAHS\)](#), both with physically-oriented models (Gelfan et al., 2015;  
103 Magand et al, 2015), conceptual models (Brigode et al., 2015; Efstratiadis et al., 2015; Hughes, 2015;  
104 Kling et al., 2015; Li et al., 2015; Yu and Zhu, 2015) or data-based models (Tanaka and Tachikawa,  
105 2015; Taver et al., 2015).

106 Recently, with the growing concern on model robustness in link with the Panta Rhei decade of the  
107 [International Association of Hydrological Sciences \(IAHS\)](#) (Montanari et al., 2013), a slow but steadily  
108 increasing interest is noticeable for procedures inspired by Klemeš's DSST (see e.g. the Unsolved  
109 Hydrological Problem n° 19 in the paper by Blöschl et al., 2019: *How can hydrological models be  
110 adapted to be able to extrapolate to changing conditions?*). A few studies used the original DSST or  
111 GSST to implement more demanding model tests (Bisselink et al., 2016; Gelfan and Millionshchikova,  
112 2018; Rau et al., 2019; Vormoor et al., 2018). For example, based on an ensemble approach using six  
113 hydrological models, Broderick et al. (2016) investigated under DSST conditions how the robustness

114 can be improved by multi-model combinations. ~~They recommend selecting the best available~~  
115 ~~analogues of expected annual mean and seasonal climate conditions.~~

116 A few authors also tried to propose improved implementations of these testing schemes. Seiller et al.  
117 (2012) used non-continuous periods or years selected on mean temperature and precipitation to  
118 enhance the contrast between testing periods. This idea to jointly use these two climate variables to  
119 select periods was further investigated by Gaborit et al. (2015), who assessed how the temporal  
120 model robustness can be improved by advanced calibration schemes. They showed that the  
121 robustness of the tested model was improved when going from humid-cold to dry-warm or from dry-  
122 cold to humid-warm conditions when using regional calibration instead of local calibration. Dakhlaoui  
123 et al. (2017) investigated the impact of DSST on model robustness by selecting dry/wet and cold/hot  
124 hydrological years to increase the contrast in climate conditions between calibration and validation  
125 periods. These authors later proposed a bootstrap technique to widen the testing conditions  
126 (Dakhlaoui et al. 2019). The investigations of Fowler et al. (2018) identified some limits of the DSST  
127 procedure and concluded that “model evaluation based solely on the DSST is hampered due to  
128 contingency on the chosen calibration method, and it is difficult to distinguish which cases of DSST  
129 failure are truly caused by model structural inadequacy”. Last, Motavita et al. (2019) combined DSST  
130 with periods of variable length, and conclude that parameters obtained on dry periods may be more  
131 robust.

132 All these past studies show that there is still methodological work needed on the issue of model  
133 testing and robustness assessment. This note is a further step in that direction.

### 134 1.3 Scope of the technical note

135 This note presents a new generic diagnostic framework inspired by Klemeš’s DSST procedure and by  
136 our own previous attempts (Coron et al., 2012; Thirel et al., 2015a) to assess [the relative confidence](#)  
137 ~~one may have whether with~~ a hydrological model ~~can be considered “climate proof” to be used in a~~  
138 ~~changing climate context~~. One of the problems of existing methods is the requirement of multiple  
139 [calibrations of hydrological models](#): these are relatively easy to implement with parsimonious  
140 conceptual models but definitively not with complex models that require long interventions by  
141 expert modellers and, obviously, not for those models with a once-for-all parameterisation.

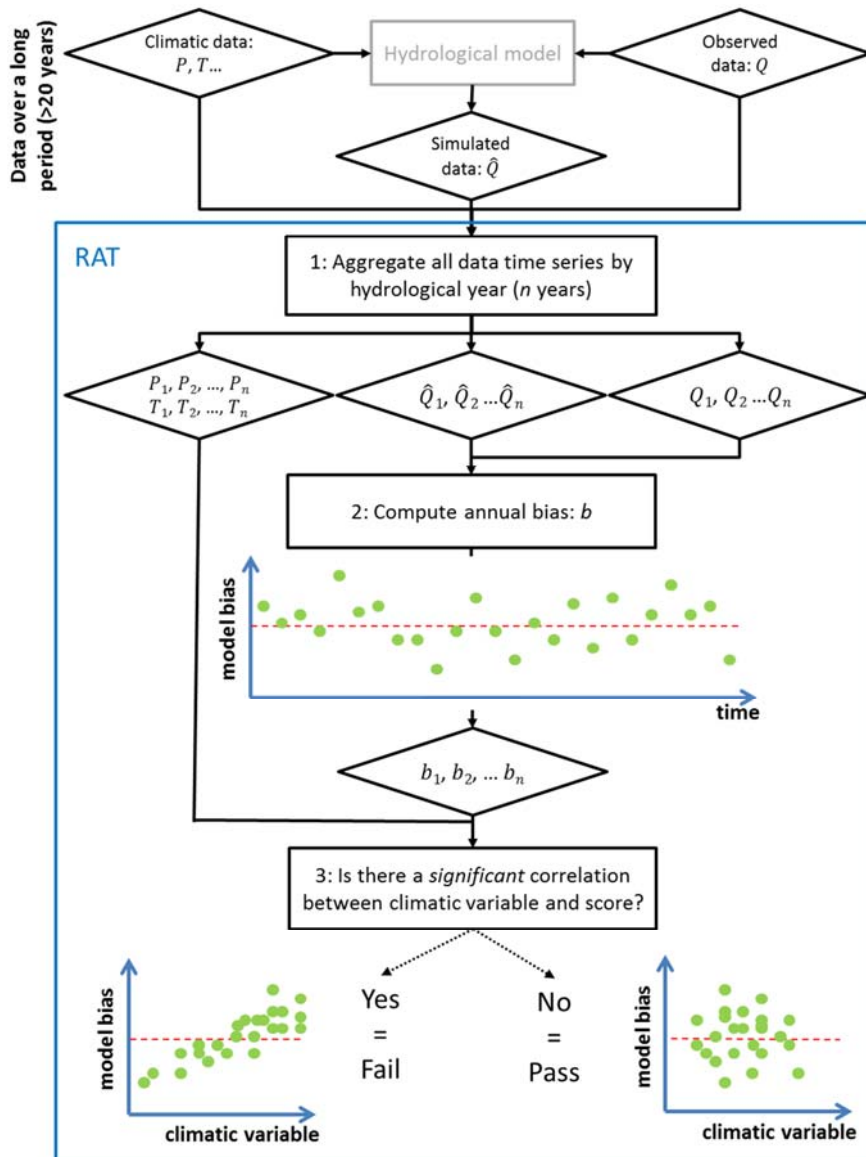
142 Here, we propose a framework that is applicable with only one long period for which a model  
143 simulation is available. Thus, the proposed test is even applicable to those models that do not  
144 require calibration (or to those for which only a single calibration exists).

145 Section 2 presents and discusses the concept of the proposed test, section 3 presents the catchment  
146 set and the evaluation method, and section 4 illustrates the application of the test on a set of French  
147 catchments, with a comparison to a reference procedure.

## 148 2 The robustness assessment test (RAT) concept

149 The robustness assessment test (RAT) proposed in this note is inspired by the work of Coron et al.  
150 (2014). The specificity of the RAT is that it requires only one ~~calibration (or one~~  
151 ~~parameterisation) simulation~~ covering a sufficiently-long period (at least ~~30-20~~ years) with as much  
152 climatic variability as possible. Thus, it applies at the same time to simple conceptual models that can

153 be calibrated automatically, to more complex models requiring expert calibration, and to  
154 uncalibrated models for which parameters are derived from the measurement of certain physical  
155 properties. The RAT consists in computing a relevant numeric [bias](#) criterion repeatedly each year and  
156 then exploring its correlation with a climatic factor deemed meaningful, in order to identify  
157 undesirable dependencies and thus to assess the extrapolation capacity (Roberts et al., 2017) of any  
158 hydrological model. Indeed, if the performances of a model are shown to be dependent on a given  
159 climate variable, this can be an issue when the model is used on a period with a changing climate.  
160 The flowchart in Figure 1 summarizes the concept.

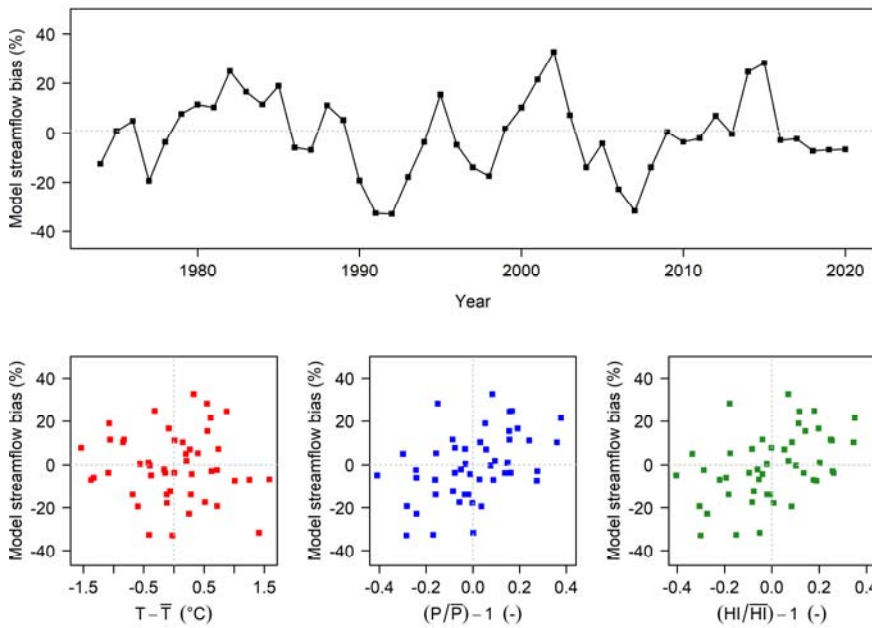


161

162 **Figure 1. Flow chart of the Robustness Assessment Test**

163 An example is shown in Figure 2, with a daily time step hydrological model calibrated on a 47-year  
 164 streamflow record. Note that this plot could be obtained from any hydrological model calibrated or  
 165 not. The relative streamflow bias ( $(\overline{Q_{sim}}/\overline{Q_{obs}} - 1)$ , with  $\overline{Q_{sim}}$  and  $\overline{Q_{obs}}$  being the mean simulated  
 166 and observed streamflows respectively) is calculated on an annual basis (47 values in total). Then,

167 the annual bias values are plotted against climate descriptors, typically the annual temperature  
 168 absolute anomaly ( $T - \bar{T}$ , where  $T$  is the annual mean and  $\bar{T}$  is the long-term mean annual  
 169 temperature), the annual precipitation relative anomaly  $P/\bar{P} - 1$  and the humidity index relative  
 170 anomaly  $HI/\bar{HI} - 1$ , where  $HI = P/E_0$ ,  $E_0$  being the potential evaporation). Note that the mean  
 171 annual values are computed on hydrological years (here from August 1<sup>st</sup> of year  $n-1$  to July 31<sup>st</sup> of  
 172 year  $n$ ). In this example, there is a slight dependency of model bias on precipitation and humidity  
 173 index. Clearly, this could be a problem if we were to use this model in an extrapolation mode.



174  
 175 **Figure 2. Robustness Assessment Test (RAT) applied to a hydrological model: the upper graph presents the**  
 176 **evolution in time (year by year) of model streamflow bias; the lower scatterplots present the relationship**  
 177 **between model bias and climatic variables (temperature  $T$ , precipitation  $P$  and humidity index  $HI$ , from left**  
 178 **to right)**

179 Whereas the methods based on the split-sample test (i.e. Coron et al, 2012 and Thirel et al., 2015b)  
 180 evaluate model robustness on periods that are independent of the calibration period, it is not the  
 181 case for the RAT. Consequently, one could fear that the results of the RAT evaluation may be  
 182 influenced by the calibration process. However, because the RAT uses a very long period for  
 183 calibration, we hypothesize that the weight of each individual year in the overall calibration process  
 184 is small, almost negligible. This assumption can be checked by comparing the RAT with a leave-one-  
 185 out SST (see Appendix). The analysis showed that this hypothesis is reasonable for long time series,  
 186 but that the RAT is not applicable when the available time period is too short (less than 20 years).

187 Last, we would like to mention that the RAT procedure is different from the Proxy metric for Model  
 188 Robustness (PMR) presented by Royer-Gaspard et al. (2021), even if both methods aim to evaluate

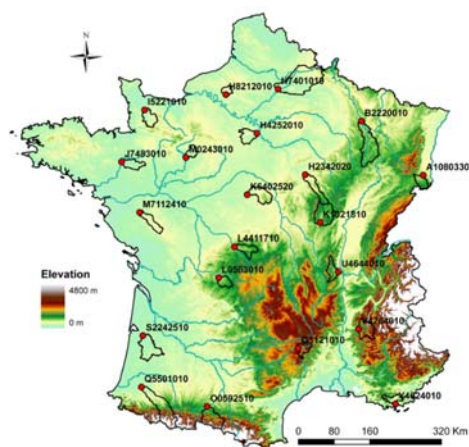
189 hydrological model robustness without employing a multiple calibrations process: the PMR is a  
190 simple metric to estimate the robustness of a hydrological model, while the RAT is a method to  
191 diagnose the dependencies of model errors to certain types of climatic changes. Thus, the RAT and  
192 the PMR may be seen as complementary tools to assess a variety of aspects about model robustness.

### 193 3 Material and methods

#### 194 3.1 Catchment set

195 We employed the dataset previously used by Nicolle et al. (2014), comprising 21 French catchments  
196 (Figure 3), with complementary data until 2020. Catchments were chosen to represent a large range  
197 of physical and climatic conditions in France, with sufficiently-long observation time series (daily  
198 streamflow from 1974 to 2020) in order to provide a diverse representation of past hydroclimatic  
199 conditions. Streamflow data come from the French HYDRO database (Leleu et al., 2014) and with  
200 quality control performed by the operational hydrometric services. Catchment size ranges from 380  
201 to 4,300 km<sup>2</sup> and median elevation from 70 to 1020 m.

202 The daily precipitation and temperature data originate from the gridded SAFRAN climate reanalysis  
203 (Vidal et al., 2010) over the 1959–2020 period. More information about the catchment set can be  
204 found in Nicolle et al. (2014). Aggregated catchment files and computation of Oudin potential  
205 evaporation (Oudin et al., 2005) was made as described in Delaigue et al. (2018).



206  
207 **Figure 3. Location of the 21 catchments in France. Red dots represent the catchment outlets**

#### 208 3.2 Hydrological model

209 The RAT diagnostic framework is generic and can be applied to any type of model. Here daily  
210 streamflow was simulated using the daily lumped GR4J rainfall–runoff model (Perrin et al., 2003).  
211 The objective function used for calibration is the KGE criterion (Gupta et al., 2009) computed on



212 square-root-transformed flows. Model implementation was done with the airGR R package (Coron et  
213 al., 2017, 2018).

### 214 3.3 Evaluation of the RAT framework

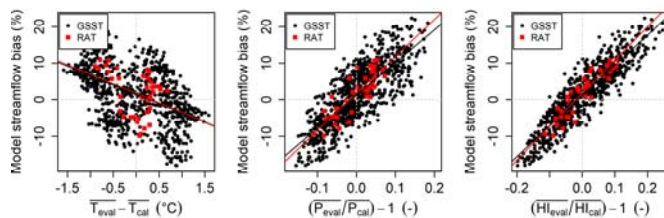
215 The RAT was evaluated against the GSST of Coron et al. (2012) used as a benchmark, in order to  
216 check whether it yields similar results. The GSST procedure was applied to each catchment using a  
217 10-year period to calibrate the model. For each calibration, each 10-year sliding period over the  
218 remaining available period, strictly independent of the calibration one, was used to evaluate the  
219 model. The results of the two approaches were compared by plotting on the same graph the annual  
220 streamflow bias obtained from the unique simulation period for the RAT, and the average  
221 streamflow bias over the sliding calibration-validation time periods for GSST, as a function of  
222 temperature, precipitation and humidity anomalies as in Figure 2. The similarity of the trends  
223 (between streamflow bias and climatic anomaly) obtained by the two methods was evaluated on the  
224 catchment set by comparing the slope and intercept of the linear regressions obtained in each case.

225 We then identified the catchments where the RAT procedure detected a ~~lack of~~ dependency of  
226 streamflow bias to ~~one or several~~ climate variables, ~~or a dependency to one or several variables~~. The  
227 Spearman correlation between model bias and climate variables was computed and a significance  
228 threshold of 5% was used (p-value 0.05).

## 229 4 Results

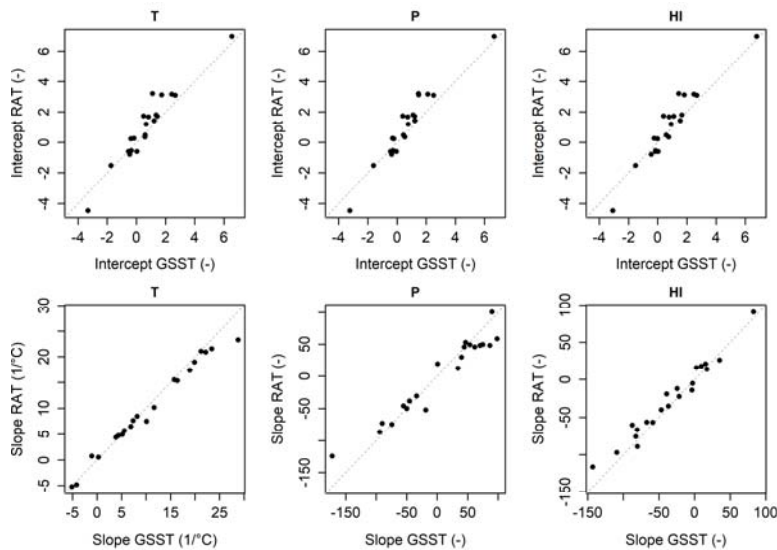
### 230 4.1 Comparison between the RAT and the GSST procedure

231 Figure 4 presents an example for the Orge River at Morsang-sur-Orge: GSST points are represented  
232 by black dots and RAT points by red squares. Let us first note that since red points represent only  
233 each of the  $N$  years of the period for the RAT and black points represent all GSST possible  
234 independent calibration-validation pairs (a number close to  $N(N-1)$ ), black points are much more  
235 numerous. We can observe that the amplitude of both streamflow bias and climatic variable change  
236 is larger for the GSST than for the RAT as there are more calibration periods, whatever the climatic  
237 variable (P, T or HI). However, the trends in the scatterplot are quite similar. Graphs for all  
238 catchments are provided as supplementary material.



239  
240 **Figure 4. Streamflow bias obtained with the RAT (red squares) and the GSST (black dots), as a function of**  
241 **temperature, precipitation and humidity index anomalies, for the Orge River at Morsang-sur-Orge**  
242 **(H4252010) (934 km<sup>2</sup>).**

243 To summarize the results on the 21 catchments, we present on Figure 5 the slope and intercept of a  
 244 linear regression computed between model streamflow bias and climatic variable anomaly, for the  
 245 GSST and the RAT over the 21 catchments: the slope of the regressions obtained for both methods  
 246 are very similar and the intercept also exhibits a good match (although somewhat larger differences).



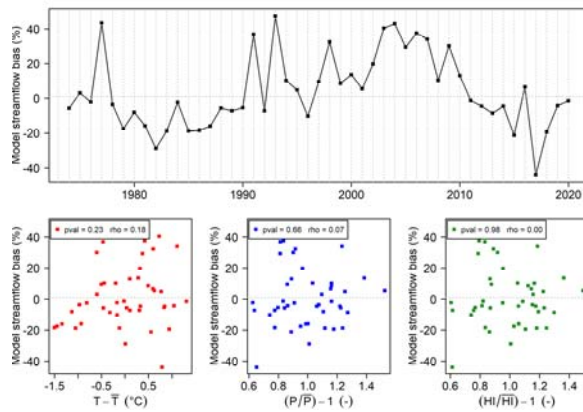
247  
 248 **Figure 5. Comparison of slopes and intercept of linear regressions between streamflow bias and temperature**  
 249 **(T), precipitation (P) and humidity index (HI) anomalies (from left to right) obtained by the GSST and the RAT**  
 250 **procedures (each point represents one of the 21 test catchments)**

251 We can thus conclude that the RAT reproduces the results of GSST, but at a much lower  
 252 computational [price/cost](#), and this is what we were aiming at. One should however acknowledge that  
 253 switching from the GSST to the RAT unavoidably reduces the severity of the climate anomalies we  
 254 can expose the hydrological models to: indeed, the climate anomalies with the RAT are computed  
 255 with respect to the mean over the whole period, whilst with the GSST they are computed between  
 256 two shorter (and hence potentially more different) periods.

#### 257 4.2 Application of the RAT procedure to the detection of climate dependencies

258 We now illustrate the different behaviours found among the 21 catchments when applying the RAT  
 259 procedure. The significance of the link between model bias and climate anomalies was based on the  
 260 Spearman correlation and a 5% threshold. Five cases were identified:

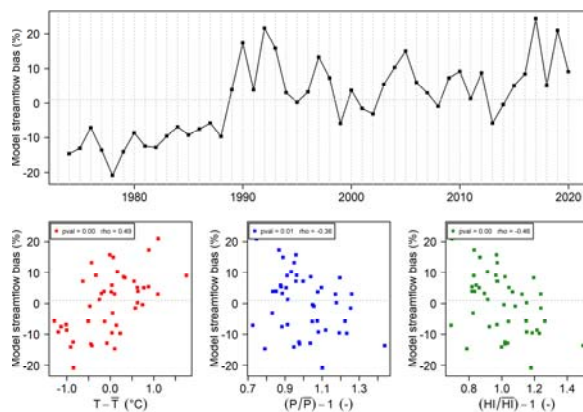
- 261 1. **No climate dependency** (Figure 6): This is the case for 6 catchments out of 21 and the  
 262 expected situation of a “robust” model. The different plots show a lack of dependence,  
 263 for temperature, precipitation and humidity index alike. For the catchment of Figure 6,  
 264 the p-value of the Spearman correlation is [quite](#) high (between 0.23 and 0.98) and thus  
 265 not significant.



266

267 **Figure 6. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute**  
 268 **anomalies (bottom left), and precipitation P (bottom centre) and humidity index  $P/E_0$  (bottom right)**  
 269 **anomalies, for the Orne Saosnoise River at Montbizot (M0243010) (510 km<sup>2</sup>).**

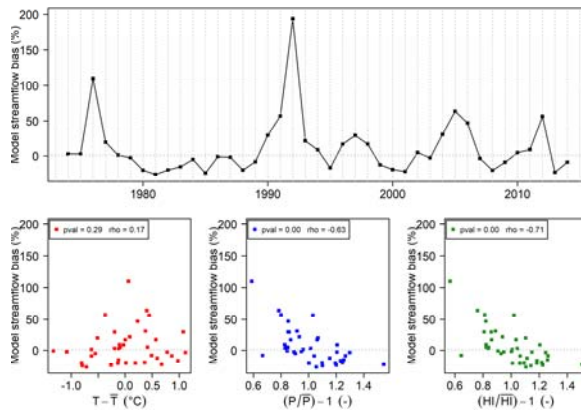
270 2. **Significant dependency on annual temperature, precipitation and humidity index**  
 271 (Figure 7): This is a clearly undesirable situation illustrating a lack of robustness of the  
 272 hydrological model. It happens on only two catchments out of 21. The Spearman  
 273 correlation between model bias and temperature, precipitation and humidity index  
 274 anomalies (respectively 0.49, -0.36 and -0.46) is significant (i.e. below the classic  
 275 significance threshold of 5%). In Figure 7, the annual bias shows an increasing trend with  
 276 annual temperature and a decreasing trend with annual precipitation and humidity  
 277 index.



278

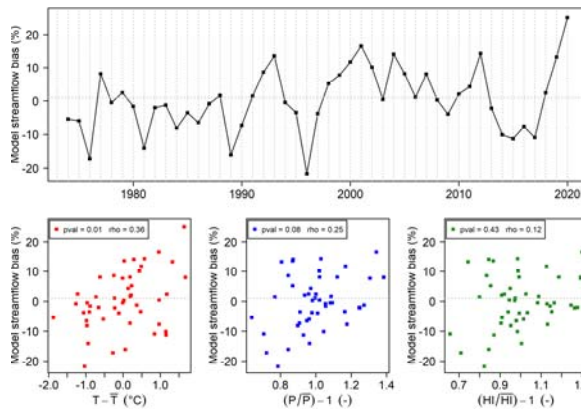
279 **Figure 7. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute**  
 280 **anomalies (bottom left), and precipitation P (bottom center) and humidity index  $P/E_0$  (bottom right)**  
 281 **anomalies, for the Arroux River at Etang-sur-Arroux (K1321810) (1790 km<sup>2</sup>)**

282 3. Significant climate dependency on precipitation and humidity index but not on  
 283 temperature (Figure 8). This case happens on 5 of the 21 catchments.



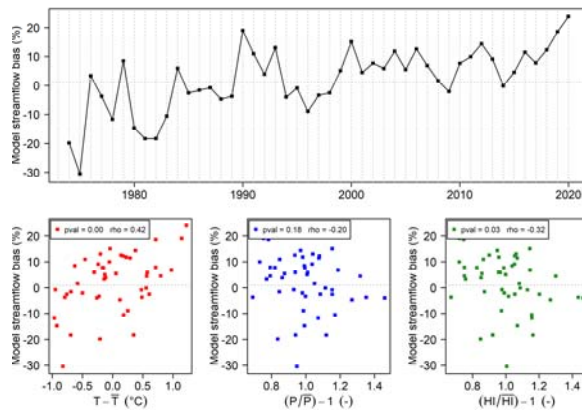
284  
 285 Figure 8. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute  
 286 anomalies (bottom left), and precipitation P (bottom center) and humidity index  $P/E_0$  (bottom right)  
 287 anomalies, for the Seiche River at Bruz (J7483010) (810 km<sup>2</sup>)

288 4. Significant climate dependency on temperature but not on precipitation and humidity  
 289 index (Figure 9). This case happens on 3 of the 21 catchments.



290  
 291 Figure 9. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute  
 292 changes (bottom left), and precipitation P (bottom center) and humidity index  $P/E_0$  (bottom right)  
 293 anomalies, for the Ill at Didenheim (A1080330) (670 km<sup>2</sup>)

294 5. Significant climate dependency on temperature and humidity index but not on  
 295 precipitation (Figure 10). This case happens on 5 of the 21 catchments.



296

297 **Figure 10.** Streamflow annual bias obtained with the RAT function of time (top), temperature absolute  
 298 changes (bottom left), and precipitation P (bottom center) and humidity index  $P/E_0$  (bottom right)  
 299 anomalies, for the Briance River at Condat-sur-Vienne (L0563010) (597 km<sup>2</sup>)

300 **4.3 How to use RAT results?**

301 A question that many modelers may ask us is *what can be done when different types of model failure*  
 302 *are identified? Some of the authors of this paper have long be fond of the concept of Crash test*  
 303 *(Andréassian et al., 2009), and we would like to argue here that the RAT too can be seen as a kind of*  
 304 *crash-test. As all crash tests, it will end up identifying failures. But the fact that a car may be*  
 305 *destroyed when projected against a wall does not mean that it is entirely unsafe, it rather means that*  
 306 *it is not entirely safe. Although we are conscious of this, we keep driving cars... but, we are also*  
 307 *willing to pay (invest) more for a safer car (even if this safer-and-more-expensive toy did also*  
 308 *ultimately fail the crash test). We believe that the same will occur with hydrological models: The RAT*  
 309 *may help identify safer models, or safer ways to parameterize models. If applied on large datasets, it*  
 310 *may help identify model flaws, and thus help us work to eliminate them. It will not however help*  
 311 *identify perfect models: these do not exist.*

312

313 **5 Conclusion**

314 The proposed robustness assessment test (RAT) is an easy-to-implement evaluation framework that  
 315 allows robustness evaluation from all types of hydrological models to be compared, by using only  
 316 one long period for which model simulations are available. The RAT consists in identifying undesired  
 317 dependencies of model errors to the variations of some climate variables over time. Such  
 318 dependencies can indeed be detrimental for model performance in a changing climate context. This  
 319 test can be particularly useful for climate change impact studies where the robustness of hydrological  
 320 models is often not evaluated at all: as such, our test can help users to discriminate alternative  
 321 models and select the most reliable models for climate change studies, which ultimately should  
 322 reduce uncertainties on climate change impact predictions (Krysanova et al., 2018).

Mis en forme : Titre 2

Mis en forme : Police :Italique

Mis en forme : Police :Italique

323 The proposed test has obviously its limits, and a first difficulty that we see in using the RAT is that it is  
324 only applicable in cases where the hypothesis of independence between the 1-year subperiods and  
325 the whole period is sufficient. This is the case when long series are available (at least 20 years, see  
326 last graph in appendix). If it is not the case, the RAT procedure should not be used. Therefore, we  
327 would indeed recommend its use in cases where modellers cannot “afford” multiple calibrations, or  
328 where the parameterisation strategy is considered (by the modeller) as ‘calibration free’ (i.e.  
329 physically-based models). A few other limitations should be mentioned:

- 330 1. In this note, the RAT concept was illustrated with a rank-based test (Spearman correlation) and  
331 a significance threshold of 0.05. Like all thresholds, this one is arbitrary. Moreover, other non-  
332 parametric tests could be used and would probably yield slightly different results (we also  
333 tested the Kendall tau test, with very similar results, but do not show the results here);
- 334 2. Detecting a relationship between model bias and a climate variable using the RAT does not  
335 allow to directly conclude on a lack of model robustness, because even a robust model will be  
336 affected by a trend in input data. Indeed, changes in the precipitation monitoring network or  
337 in the hydrometric rating curves can also give the false, yielding the impression that the  
338 hydrological model lacks robustness. Such an erroneous conclusion could also be due to  
339 widespread changes in land use, construction of an unaccounted storage reservoir or the  
340 evolution of water uses. Some of the lacks of robustness detected among the 21 catchments  
341 presented here could be in fact due to metrological causes;
- 342 3. Also, because of the ongoing rise of temperatures (over the last 40 years at least), we have a  
343 correlation between temperature and time since the beginning of streamgaging. If for any  
344 reason, time is having an impact on model bias, this may cause an artefact in the RAT in the  
345 form of a dependency between model bias and temperature;
- 346 4. Similarly to the Differential Split Sample Test, the diagnostic of model climatic robustness is  
347 limited to the climatic variable against which the bias is compared. As such, the RAT should not  
348 be seen as an *absolute* test, but rather as a *necessary but not sufficient* condition to use a  
349 model for climate change studies: because the climatic variability present in the past  
350 observations is limited to the historic range, so is the extrapolation test. With Popper’s words  
351 (Popper, 1959), the RAT can only allow falsifying a hydrological model... but not proving it  
352 true/right;
- 353 5. Although it would be tempting to transform the RAT into a post-processing method, we do not  
354 recommend it. Indeed, detecting a relationship between model bias and a climate variable  
355 using the RAT does not necessarily mean that a simple (linear) debiasing solution can be  
356 proposed to solve the issue (see e.g. the paper by Bellprat et al. (2013) on this topic). What we  
357 do recommend is to work as much as possible on the model structure, to turn it less climate  
358 dependent;
- 359 6. Some of the modalities of the RAT, that we initially thought of importance, are not really  
360 important: this is for example the case with the use of hydrological years. We tested the  
361 twelve possible annual aggregations schemes (see [https://doi.org/10.5194/hess-2021-147-](https://doi.org/10.5194/hess-2021-147-AC6)  
362 AC6) and found no significant impact;
- 363 5-7. Upon recommendation by one of the reviewers, we tried to assess the possible impact of the  
364 quality of the precipitation forcing on RAT results (see [https://doi.org/10.5194/hess-2021-147-](https://doi.org/10.5194/hess-2021-147-AC5)  
365 AC5) and found that the type of forcing used does have an impact on RAT results (interestingly,  
366 the climatic dataset yielding the best simulation results was also the dataset yielding the less

Mis en forme : Retrait : Gauche : 0 cm, Suspendu : 1 cm

367 [catchments failing the robustness test](#)). It seems unavoidable that forcing data quality will  
368 [impact the results of RAT, but we would argue that it would similarly have an impact on the](#)  
369 [results of a Differential Split Sample Test. We believe that there is no way to avoid entirely this](#)  
370 [dependency, and that evaluating the quality of input data should be done before looking at](#)  
371 [model robustness](#);

372 6-8. Last, we could mention that a model showing a small overall annual bias (but linked to a  
373 climate variable) could still be preferred to one showing a large overall annual bias (but  
374 independent of the tested climate variables): the RAT should not be seen as the only basis for  
375 model choice.

376 Beyond the limitations, we also see the perspective for further development of the method: although  
377 this note only considered overall model bias (as the most basic requirement for a model to be used  
378 to predict the impact of a future climate), we think that this methodology could be applied to bias in  
379 different flow ranges (low or high flows) or to statistical indicators describing low-flow characteristics  
380 or maximum annual streamflow. And characteristics other than bias could be tested, e.g. ratios  
381 pertaining to the variability of flows. Further, while we only tested the dependency to mean annual  
382 temperature, precipitation and humidity index, other characteristics, such as precipitation intensity  
383 or fraction of snowfall, could be considered in this framework.

## 384 6 Acknowledgments

385 This work was funded by the project AQUACLEW, which is part of ERA4CS, an ERA-NET initiated by JPI  
386 Climate, and funded by FORMAS (SE), DLR (DE), BMWFW (AT), IFD (DK), MINECO (ES), ANR (FR) with  
387 co-funding by the European Commission [Grant 690462].

388 The authors gratefully acknowledge the comments of Prof. Jens-Christian Refsgaard and Dr Nans  
389 Addor on a preliminary version of the note, [and the reviews by Dr Bettina Schäfli and by two](#)  
390 [anonymous reviewers](#).

391 The gridded SAFRAN climate reanalysis data can be ordered from Météo-France.

392 Observed streamflow data are available on the French HYDRO database  
393 (<http://www.hydro.eaufrance.fr/>).

394 The GR models, including GR4J, are available from the airGR R package.

## 395 7 References

396 Andréassian, V., Le Moine, N., Perrin, C., Ramos, M.-H., Oudin, L., Mathevet, T., Lerat, J., Berthet, L  
397 (2012). All that glitters is not gold: the case of calibrating hydrological models. *Hydrological*  
398 *Processes*, 26(14), 2206–2210. <https://doi.org/10.1002/hyp.9264>

399 [Andréassian, V., Perrin, C., Berthet, L., Le Moine, N., Lerat, J., Loumagne, C., Oudin, L., Mathevet, T.,](#)  
400 [Ramos, M.H., Valéry, A. \(2009\). \*Crash tests for a standardized evaluation of hydrological\*](#)  
401 [models. \*Hydrology and Earth System Sciences\*, 13, 1757-1764. \[https://doi.org/10.5194/hess-\]\(https://doi.org/10.5194/hess-13-1757-2009\)](#)  
402 [13-1757-2009](#)

403 Arlot, S., & Celisse, A. (2010). A survey of cross-validation procedures for model selection. *Statistics*  
404 *Surveys*, 4(0), 40–79. <https://doi.org/10.1214/09-SS054>

Mis en forme : Français (France)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)

- 405 Bellprat, O., Kotlarski, S., Lüthi, D., & Schär, C. (2013). Physical constraints for temperature biases in  
406 climate models: limits of temperature biases. Geophysical Research Letters, 40(15), 4042–  
407 4047. <https://doi.org/10.1002/grl.50737>
- 408 Beven, K. (2016). Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood,  
409 hypothesis testing, and communication. *Hydrological Sciences Journal*, 61(9), 1652–1665.  
410 <https://doi.org/10.1080/02626667.2015.1031761>
- 411 Bisselink, B., Zambrano-Bigiarini, M., Burek, P., & de Roo, A. (2016). Assessing the role of uncertain  
412 precipitation estimates on the robustness of hydrological model parameters under highly  
413 variable climate conditions. Journal of Hydrology: Regional Studies, 8, 112–129.  
414 <https://doi.org/10.1016/j.ejrh.2016.09.003>
- 415 Blöschl, G., et al. 2019. Twenty-three Unsolved Problems in Hydrology – a community perspective.  
416 *Hydrological Sciences Journal*, <https://doi.org/10.1080/02626667.2019.1620507>
- 417 Brigode, P., et al. (2015). Dependence of model-based extreme flood estimation on the calibration  
418 period: case study of the Kamp River (Austria). *Hydrological Sciences Journal*, 60 (7–8).  
419 <https://doi.org/10.1080/02626667.2015.1006632>
- 420 Broderick, C., Matthews, T., Wilby, R. L., Bastola, S., & Murphy, C. (2016). Transferability of  
421 hydrological models and ensemble averaging methods between contrasting climatic periods.  
422 *Water Resources Research*, 52(10), 8343–8373. <https://doi.org/10.1002/2016WR018850>
- 423 Coron, L., Andréassian, V., Bourqui, M., Perrin, C., & Hendrickx, F. (2011). Pathologies of hydrological  
424 models used in changing climatic conditions: a review. In S. Franks, E. Boegh, E. Blyth, D.  
425 Hannah, & K. K. Yilmaz (Eds.), Hydro-climatology: Variability and Change. IAHS Red Books  
426 Series 344 (pp. 39–44). Wallingford: IAHS.
- 427 Coron, L., Andréassian, V., Perrin, C., Lerat, J., Vaze, J., Bourqui, M., & Hendrickx, F. (2012). Crash  
428 testing hydrological models in contrasted climate conditions: An experiment on 216 Australian  
429 catchments. Water Resources Research, 48, W05552. <https://doi.org/10.1029/2011WR011721>
- 430 Coron, L., Andréassian, V., Perrin, C., Bourqui, M., & Hendrickx, F. (2014). On the lack of robustness of  
431 hydrologic models regarding water balance simulation: a diagnostic approach applied to three  
432 models of increasing complexity on 20 mountainous catchments. *Hydrol. Earth Syst. Sci.*, 18(2),  
433 727–746.
- 434 Coron, L., Thirel, G., Delaigue, O., Perrin, C., & Andréassian, V. (2017). The Suite of Lumped GR  
435 Hydrological Models in an R package. *Environmental Modelling and Software*, 94, 337.  
436 <https://doi.org/10.1016/j.envsoft.2017.05.002>
- 437 Coron, L., Delaigue, O., Thirel, G., Perrin, C. and Michel, C. (2020). airGR: Suite of GR Hydrological  
438 Models for Precipitation-Runoff Modelling. R package version 1.4.3.65. DOI:  
439 10.15454/EX11NA. URL: <https://CRAN.R-project.org/package=airGR>.
- 440 Dakhlaoui, H., Ruelland, D., Trambly, Y., & Bargaoui, Z. (2017). Evaluating the robustness of  
441 conceptual rainfall-runoff models under climate variability in northern Tunisia. *Journal of*  
442 *Hydrology*, 550, 201–217. <https://doi.org/10.1016/j.jhydrol.2017.04.032>
- 443 Dakhlaoui, H., Ruelland, D., & Trambly, Y. (2019). A bootstrap-based differential split-sample test to  
444 assess the transferability of conceptual rainfall-runoff models under past and future climate  
445 variability. Journal of Hydrology, 575, 470–486. <https://doi.org/10.1016/j.jhydrol.2019.05.056>
- 446 Delaigue, O., Génot, Lebecherel, L., Brigode, P., & Bourgin, P. Y. (2018). Base de données  
447 hydroclimatiques observées à l'échelle de la France. IRSTEA. IRSTEA, UR HYCAR, Équipe  
448 Hydrologie des bassins versants, Antony. Retrieved from  
449 <https://webgr.inrae.fr/en/activities/database-1-2/>

Mis en forme : Anglais (États-Unis)

Mis en forme : Français (France)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)



450 Donnelly-Makowecki, L. M., & Moore, R. D. (1999). Hierarchical testing of three rainfall-runoff models  
 451 in small forested catchments. *Journal of Hydrology*, 219(3–4), 136–152.

452 Efstratiadis, A., Nalbantis, I., and Koutsoyiannis, D. (2015). Hydrological modelling of temporally-  
 453 varying catchments: facets of change and the value of information. *Hydrological Sciences*  
 454 *Journal*, 60 (7–8). <https://doi.org/10.1080/02626667.2014.982123>

455 Fowler, K., Coxon, G., Freer, J., Peel, M., Wagener, T., Western, A., et al. (2018). Simulating Runoff  
 456 Under Changing Climatic Conditions: A Framework for Model Improvement. *Water Resources*  
 457 *Research*, 54(12), 9812–9832. <https://doi.org/10.1029/2018WR023989>

458 Gaborit, É., Ricard, S., Lachance-Cloutier, S., Anctil, F., & Turcotte, R. (2015). Comparing global and  
 459 local calibration schemes from a differential split-sample test perspective. *Canadian Journal of*  
 460 *Earth Sciences*, 52(11), 990–999. <https://doi.org/10.1139/cjes-2015-0015>

461 Gelfan, A.N., & Millionshchikova, T.D. (2018). Validation of a Hydrological Model Intended for Impact  
 462 Study: Problem Statement and Solution Example for Selenga River Basin. *Water Resour* 45, 90–  
 463 101. <https://doi.org/10.1134/S0097807818050354>

464 Gelfan, A., et al. (2015). Testing robustness of the physically-based ECOMAG model with respect to  
 465 changing conditions. *Hydrological Sciences Journal*, 60 (7–8).  
 466 <https://doi.org/10.1080/02626667.2014.935780>

467 Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). Decomposition of the mean squared  
 468 error and NSE performance criteria: Implications for improving hydrological modelling. *Journal*  
 469 *of Hydrology*, 377(1–2), 80–91. <https://doi.org/10.1016/j.jhydrol.2009.08.003>

470 Klemeš, V. (1986). Operational testing of hydrologic simulation models. *Hydrological Sciences*  
 471 *Journal*, 31(1), 13–24. <https://doi.org/10.1080/02626668609491024>

472 Krysanova, V., Donnelly, C., Gelfan, A., Gerten, D., Arheimer, B., Hattermann, F., & Kundzewicz, Z. W.  
 473 (2018). How the performance of hydrological models relates to credibility of projections under  
 474 climate change. *Hydrological Sciences Journal*, 63(5), 696–720.  
 475 <https://doi.org/10.1080/02626667.2018.1446214>

476 Hughes, D.A. (2015). Simulating temporal variability in catchment response using a monthly rainfall-  
 477 runoff model. *Hydrological Sciences Journal*, 60 (7–8).  
 478 <https://doi.org/10.1080/02626667.2014.909598>

479 Kling, H., et al. (2015). Performance of the COSERO precipitation-runoff model under non-stationary  
 480 conditions in basins with different climates. *Hydrological Sciences Journal*, 60 (7–8).  
 481 <https://doi.org/10.1080/02626667.2014.959956>

482 Larson, S. C. (1931). The shrinkage of the coefficient of multiple correlation. *Journal of Educational*  
 483 *Psychology*, 22(1), 45–55. <https://doi.org/10.1037/h0072400>

484 Leleu, I., Tonnelier, I., Puechberty, R., Gouin, P., Viquendi, I., Cobos, L., et al. (2014). La refonte du  
 485 système d'information national pour la gestion et la mise à disposition des données  
 486 hydrométriques. *La Houille Blanche*, (1), 25–32. <https://doi.org/10.1051/lhb/2014004>

487 Li, C. Z., Zhang, L., Wang, H., Zhang, Y. Q., Yu, F. L., and Yan, D. H. (2012). The transferability of  
 488 hydrological models under nonstationary climatic conditions, *Hydrol. Earth Syst. Sci.*, 16, 1239–  
 489 1254, <https://doi.org/10.5194/hess-16-1239-2012>

490 Li, H., Beldring, S., and Xu, C.-Y. (2015). Stability of model performance and parameter values on two  
 491 catchments facing changes in climatic conditions. *Hydrological Sciences Journal*, 60 (7–8).  
 492 <https://doi.org/10.1080/02626667.2014.978333>

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (Royaume-Uni)

Mis en forme : Français (France)

- 493 Magand, C., et al. (2015). Parameter transferability under changing climate: case study with a land  
 494 surface model in the Durance watershed, France. *Hydrological Science Journal*, 60 (7–8).  
 495 <https://doi.org/10.1080/02626667.2014.993643>
- 496 Montanari, A., Young, G., Savenije, H. H. G., Hughes, D., Wagener, T., Ren, L. L., et al. (2013). “Panta  
 497 Rhei—Everything Flows”: Change in hydrology and society—The IAHS Scientific Decade 2013–  
 498 2022. *Hydrological Sciences Journal*, 58(6), 1256–1275.  
 499 <https://doi.org/10.1080/02626667.2013.809088>
- 500 Mosteller, F., & Tukey, J. W. (1968). *Data Analysis, Including Statistics. The Collected Works of John*  
 501 *W. Tukey (1988) Graphics pp. 1965-1985, vol. 5 (123)*
- 502 Motavita, D.F., R. Chow, A. Guthke & W. Nowak. (2019). The comprehensive differential split-sample  
 503 test: A stress-test for hydrological model robustness under climate variability, *Journal of*  
 504 *Hydrology*, 573: 501-515, <https://doi.org/10.1016/j.jhydrol.2019.03.054>
- 505 Nicolle, P., Pushpalatha, R., Perrin, C., François, D., Thiéry, D., Mathevet, T., et al. (2014).  
 506 Benchmarking hydrological models for low-flow simulation and forecasting on French  
 507 catchments. *Hydrol. Earth Syst. Sci.*, 18(8), 2829–2857. [https://doi.org/10.5194/hess-18-2829-](https://doi.org/10.5194/hess-18-2829-508-2014)  
 508 2014
- 509 Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andréassian, V., Anctil, F., & Loumagne, C. (2005). Which  
 510 potential evapotranspiration input for a lumped rainfall–runoff model?: Part 2-Towards a  
 511 simple and efficient potential evapotranspiration model for rainfall–runoff modelling. *Journal*  
 512 *of Hydrology*, 303(1–4), 290–306. <https://doi.org/10.1016/j.jhydrol.2004.08.026>
- 513 Perrin, C., Michel, C., & Andréassian, V. (2003). Improvement of a parsimonious model for  
 514 streamflow simulation. *Journal of Hydrology*, 279(1–4), 275–289.  
 515 [https://doi.org/10.1016/S0022-1694\(03\)00225-7](https://doi.org/10.1016/S0022-1694(03)00225-7)
- 516 Popper, K. (1959). *The logic of scientific discovery*. London: Routledge.
- 517 Rau, P., Bourrel, L., Labat, D., Ruelland, D., Frappart, F., Lavado, W., et al. (2019). Assessing  
 518 multidecadal runoff (1970–2010) using regional hydrological modelling under data and water  
 519 scarcity conditions in Peruvian Pacific catchments. *Hydrological Processes*, 33(1), 20–35.  
 520 <https://doi.org/10.1002/hyp.13318>
- 521 Refsgaard, J. C., & Henriksen, H. J. (2004). *Modelling guidelines—terminology and guiding principles*.  
 522 *Advances in Water Resources*, 27, 71–82. <https://doi.org/10.1016/j.advwatres.2003.08.006>
- 523 Refsgaard, J. C., & Knudsen, J. (1996). Operational validation and intercomparison of different types  
 524 of hydrological models. *Water Resources Research*, 32(7), 2189–2202.  
 525 <https://doi.org/10.1029/96WR00896>
- 526 Refsgaard, J. C., Madsen, H., Andréassian, V., Arnbjerg-Nielsen, K., Davidson, T. A., Drews, M., et al.  
 527 (2013). A framework for testing the ability of models to project climate change and its impacts.  
 528 *Climatic Change*, 122(1–2), 271–282. <https://doi.org/10.1007/s10584-013-0990-2>
- 529 Roberts, D. R., Bahn, V., Ciuti, S., Boyce, M. S., Elith, J., Guillerá-Arroita, G., et al. (2017). Cross-  
 530 validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure.  
 531 *Ecography*, 40(8), 913–929. <https://doi.org/10.1111/ecog.02881>
- 532 Royer-Gaspard, P., Andréassian, V., and Thirel, G. (2021) Technical note: PMR – a proxy metric to  
 533 assess hydrological model robustness in a changing climate. In review for *Hydrol. Earth Syst.*  
 534 *Sci. Discuss.*, <https://doi.org/10.5194/hess-2021-58>
- 535 Seibert, J. (2003). Reliability of model predictions outside calibration conditions. *Nordic Hydrology*,  
 536 34(5), 477–492. <https://doi.org/10.2166/nh.2003.0019>

Mis en forme : Anglais (Royaume-Uni)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (Royaume-Uni)

- 537 Seifert, D., Sonnenborg, T. O., Refsgaard, J. C., Højberg, A. L., and Trolborg, L. (2012). Assessment of  
538 hydrological model predictive ability given multiple conceptual geological models, *Water*  
539 *Resour. Res.*, 48, W06503, doi:10.1029/2011WR011149.
- 540 Seiller, G., Anctil, F., & Perrin, C. (2012). Multimodel evaluation of twenty lumped hydrological  
541 models under contrasted climate conditions. *Hydrology and Earth System Sciences*, 16(4),  
542 1171–1189. <https://doi.org/10.5194/hess-16-1171-2012>
- 543 Tanaka, T. and Tachikawa, Y. (2015). Testing the applicability of a kinematic wave-based distributed  
544 hydrological model in two climatically contrasting catchments. *Hydrological Sciences Journal*,  
545 60 (7–8). <https://doi.org/10.1080/02626667.2014.967693>
- 546 Taver, V., et al. (2015). Feed-forward vs recurrent neural network models for non-stationarity  
547 modelling using data assimilation and adaptivity. *Hydrological Sciences Journal*, 60 (7–8).  
548 <https://doi.org/10.1080/02626667.2014.967696>
- 549 Teutschbein, C. & Seibert, J. 2013. Is bias correction of regional climate model (RCM) simulations  
550 possible for non-stationary conditions? *Hydrology and Earth System Sciences*, 17, 5061–5077.,  
551 <https://doi.org/10.5194/hess-17-5061-2013>
- 552 Thirel, G., Andréassian, V., Perrin, C., Audouy, J.-N., Berthet, L., Edwards, P., et al. (2015a). Hydrology  
553 under change: an evaluation protocol to investigate how hydrological models deal with  
554 changing catchments. *Hydrological Sciences Journal*, 60(7–8), 1184–1199.  
555 <https://doi.org/10.1080/02626667.2014.967248>
- 556 Thirel, G., Andréassian, V., & Perrin, C. (2015b). On the need to test hydrological models under  
557 changing conditions. *Hydrological Sciences Journal*, 60(7–8), 1165–1173.  
558 <https://doi.org/10.1080/02626667.2015.1050027>
- 559 Vaze, J., Post, D. A., Chiew, F. H. S., Perraud, J. M., Viney, N. R., & Teng, J. (2010). Climate non-  
560 stationarity - Validity of calibrated rainfall-runoff models for use in climate change studies.  
561 *Journal of Hydrology*, 394(3–4), 447–457. <https://doi.org/10.1016/j.jhydrol.2010.09.018>
- 562 Vidal, J.-P., Martin, E., Franchistéguy, L., Baillon, M., & Soubeyroux, J.-M. (2010). A 50-year high-  
563 resolution atmospheric reanalysis over France with the Safran system. *International Journal of*  
564 *Climatology*, 30(11), P. 1627–1644. DOI: 10.1002/joc.2003. <https://doi.org/10.1002/joc.2003>
- 565 Vormoor, K., Heistermann, M., Bronstert, A., & Lawrence, D. (2018). Hydrological model parameter  
566 (in)stability – “crash testing” the HBV model under contrasting flood seasonality conditions.  
567 *Hydrological Sciences Journal*, 63(7), 991–1007.  
568 <https://doi.org/10.1080/02626667.2018.1466056>
- 569 Wilby, R. L. (2019). A global hydrology research agenda fit for the 2030s. *Hydrology Research*, 50(6):  
570 1464-1480. <https://doi.org/10.2166/nh.2019.100>
- 571 Xu, C. (1999). Climate Change and Hydrologic Models: A Review of Existing Gaps and Recent Research  
572 Developments. *Water Resources Management*, 13(5), 369–382.  
573 <https://doi.org/10.1023/A:1008190900459>
- 574 Yu, B. and Zhu, Z. (2015). A comparative assessment of AWBM and SimHyd for forested watersheds.  
575 *Hydrological Sciences Journal*, 60 (7–8). <https://doi.org/10.1080/02626667.2014.961924>

Mis en forme : Anglais (États-Unis)

Mis en forme : Anglais (Royaume-Uni)

## 576 **8 Appendix – Checking the impact of the partial overlap between** 577 **calibration and validation periods in the RAT**

578 In this appendix, we deal with calibrated models, for which we verify that the main hypothesis  
579 underlying the RAT is reasonable, i.e. that when considering a long calibration period, the weight of

580 each individual year in the overall calibration process is almost negligible. We then explore the limits  
581 of this hypothesis when reducing the length of the overall calibration period.

582

583 • **Evaluation method**

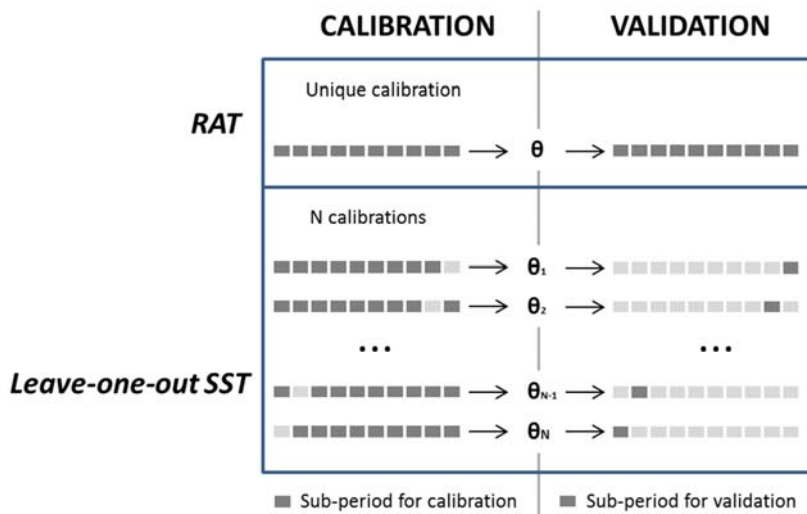
584 In order to check the impact of the partial overlap between calibration and validation periods in the  
585 RAT, it is possible (provided one works with a calibrated model) to compare the RAT with a “leave-  
586 one out” version of it, which is a classical variant of the Split Sample Test (SST): instead of computing  
587 the annual bias after a single calibration encompassing the whole period (RAT), we compute the  
588 annual bias with a different calibration each time, encompassing the whole period minus the year in  
589 question (“leave-one-out SST”).

590 The comparison between the RAT and the SST can be quantified using the root mean square  
591 difference (RMSD) of annual biases:

$$RMSD_{Bias} = \sqrt{(Bias_{RAT} - Bias_{SST})^2} \quad \text{Eq.1}$$

592 where  $Bias_{RAT}$  is the bias of validation year  $n$  when calibrating the model over the entire  
593 period (RAT procedure), and  $Bias_{SST}$  the bias of validation year  $n$  when calibrating the model  
594 over the entire period minus year  $n$  (leave-one-out SST procedure).

595 The difference between the two approaches is schematized in Figure 11: the leave-one-out  
596 procedure consists in performing  $N$  calibrations over  $(N-1)$ -year-long periods followed by an  
597 independent evaluation on the remaining 1-year-long period. As shown in Figure 11, the two  
598 procedures result in the same number of validation points ( $N$ ). Eq. 1 provides a way to quantify  
599 whether both methods differ, i.e. whether the partial overlap between calibration and validation  
600 periods in the RAT makes a difference.



601

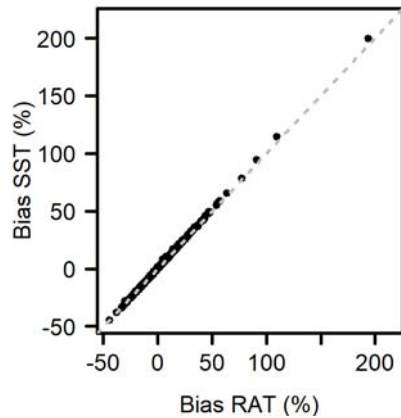
602 **Figure 11. Comparison of the RAT procedure with a leave-one-out split-sample test (SST). Both methods have**  
 603 **N validation periods (one per year). The RAT needs only one calibration, whereas the SST requires N**  
 604 **calibrations. Dark grey squares represent the years used for calibration or validation**

605

606 • **Comparison between the RAT and the leave-one-out SST**

607 Figure 12 plots the annual bias values obtained with the RAT versus the annual bias obtained with  
 608 the leave-one-out SST for the 21 test catchments, showing a total of 21x47 points. The almost perfect  
 609 alignment confirms that our underlying "negligibility" hypothesis is reasonable (at least on our  
 610 catchment set).

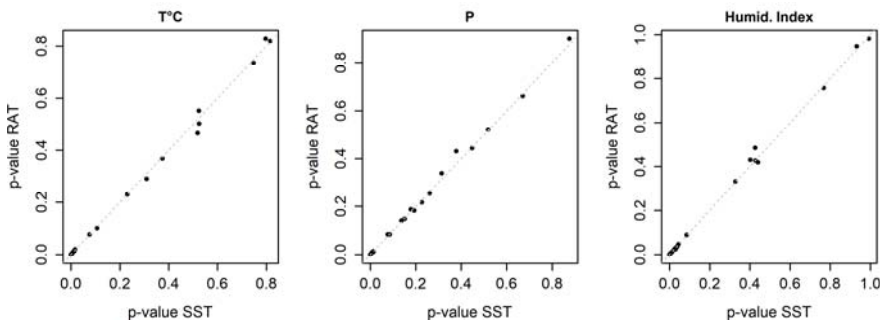
611



612  
 613 **Figure 12. Comparison of the annual bias obtained with the RAT with the annual bias obtained with the**  
 614 **leave-one-out SST. Each of the 21 catchments is represented with annual bias values (47 points by**  
 615 **catchment, 21x47 points in total)**

616 Figure 13 presents the Spearman correlation  $p$ -values for the correlation between annual bias and  
 617 changes in annual temperature, precipitation, and humidity index ( $P/E_0$ ), for the RAT and the leave-  
 618 one-out SST. The results from the RAT and the SST show the same dependencies on climate variables  
 619 (similar  $p$ -values).

620



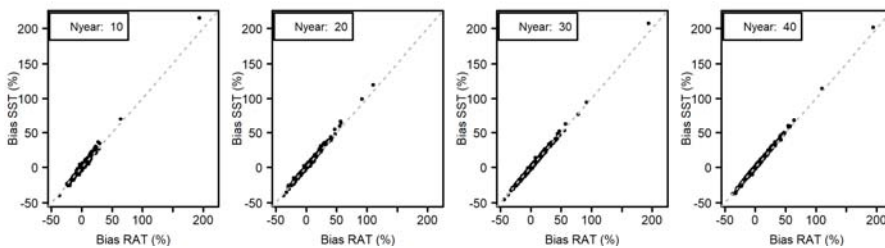
621  
 622 **Figure 13. Spearman correlation  $p$ -value from the correlation for annual bias and annual temperature,**  
 623 **precipitation, and humidity index ( $P/E_0$ ). Comparison between RAT and SST (one point per catchment)**

624

625 • **Sensitivity of the RAT procedure to the period length**

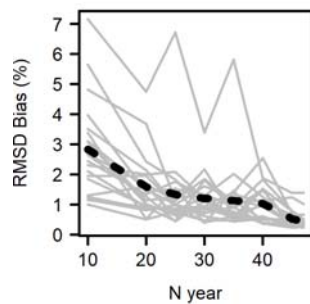
626 It is also interesting to investigate the limit of our hypothesis (i.e. that the relative weight of one year  
 627 within a long time series is very small) by progressively reducing the period length: indeed, the  
 628 shorter the data series available to calibrate the model, the more important the relative weight of  
 629 each individual year. Figure 14 compares the annual bias obtained with the RAT procedure with the  
 630 annual bias obtained with the leave-one-out SST, for 10-, 20-, 30-, and 40-year period lengths  
 631 (selection of the shorter periods was realized by sampling 10, 20, 30, and 40 years regularly among  
 632 the complete time series). The shorter the calibration period, the larger the differences between  
 633 both approaches (wider points scatter): there, we reach the limit of the single calibration procedure.  
 634 We would not advise to use RAT with time series of less than 20 years.

635



636 **Figure 14. Annual bias obtained with the RAT procedure vs. annual bias obtained with leave-one-out SST. Shorter time periods are obtained by sampling 10, 20, 30, and 40 years regularly among the complete time series. Each of the 21 catchments is represented with annual bias values**

640 These differences can be quantitatively measured by computing the RMSD (see Eq.1) between the  
 641 annual bias obtained with the RAT procedure and with the SST for different calibration period lengths  
 642 (see Figure 15). The RMSD tends to increase when the number of years available to calibrate the  
 643 model decreases, but it seems to be stable for periods longer than 20 years.



644

645 **Figure 15. RMSD between annual bias obtained with the RAT procedure and with the leave-one-out SST for**  
646 **different calibration period lengths for each catchment. The dotted line represents the mean RMSD for all**  
647 **catchments. Each grey line represents one of the 21 catchments.**