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

3

4  Pierre Nicolle[1,3], Vazken Andréassian[1,*], Paul Royer-Gaspard[1], Charles Perrin[1], Guillaume Thirel[1],
5  Laurent Coron[2], Léonard Santos[1]

6  [1]Université Paris-Saclay, INRAE, UR HYCAR, 92160, Antony, France

7  [2]EDF, DTG, Toulouse, France

8  [3]now at Université Gustave Eiffel, Nantes, France

9  [*]Corresponding author: Vazken Andréassian (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: it is therefore applicable to
17  uncalibrated models

18  •    the RAT method can be used to determine whether a hydrological model can be safely used
19  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.

# 1 Introduction

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

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

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

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

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

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

For model applications in a changing climate context, Klemeš's DSST procedure is of particular interest. Indeed, when calibration and evaluation are done over climatically-contrasted past periods, the model faces the difficulties it will have to deal with in the future. The power of DSST can be 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.

## 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 Donelly-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. (2015) put forward an SST-based protocol to investigate how hydrological models deal with
101 changing conditions, which was widely used during an IAHS workshop, both with physically-oriented
102 models (Gelfan et al., 2015; Magand et al, 2015), conceptual models (Brigode et al., 2015; Efstratiadis
103 et al., 2015; Hughes, 2015; Kling et al., 2015; Li et al., 2015; Yu and Zhu, 2015) or data-based models
104 (Tanaka and Tachikawa, 2015; Taver et al., 2015).

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

Hydrology and
Earth System
Sciences
Discussions

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

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

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

### 1.3  Scope of the technical note

134 This note presents a new generic diagnostic framework inspired by Klemeš's DSST procedure and by
135 our own previous attempts (Coron et al., 2012; Thirel et al., 2015a) to assess whether a hydrological
136 model can be considered "climate-proof". One of the problems of existing methods is the
137 requirement of multiple calibrations: these are relatively easy to implement with parsimonious
138 conceptual models but definitely not with complex models that require long interventions by
139 expert modellers and, obviously, not for those models with a once-for-all parameterisation.

140 Here, we propose a framework that is applicable with only one long period for which a model
141 simulation is available. Thus, the proposed test is even applicable to those models that do not
142 require calibration (or to those for which only a single calibration exists).
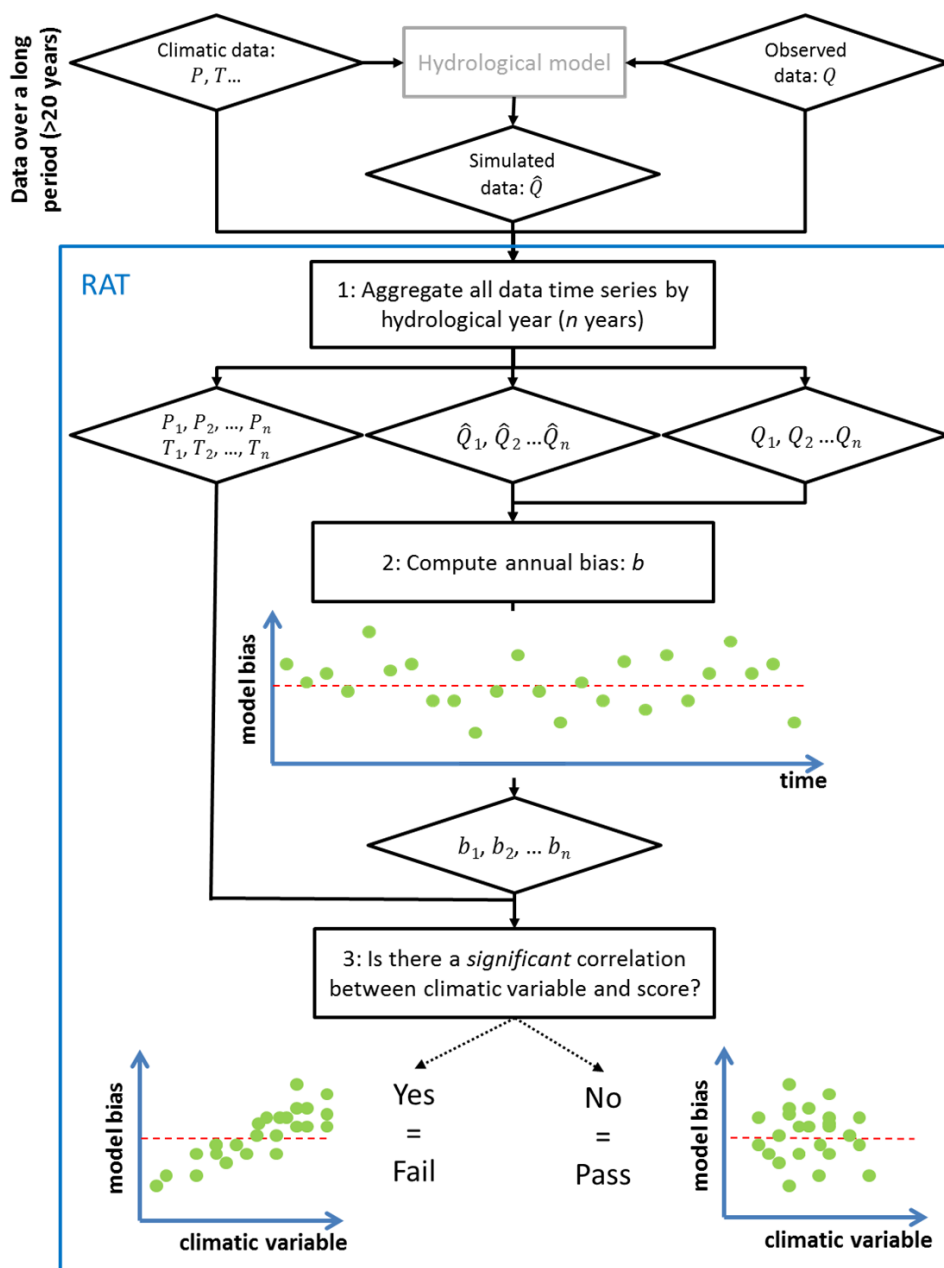
143 Section 2 presents and discusses the concept of the proposed test, section 3 presents the catchment
144 set and the evaluation method, and section 4 illustrates the application of the test on a set of French
145 catchments, with a comparison to a reference procedure.

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

147 The robustness assessment test (RAT) proposed in this note is inspired by the work of Coron et al.
148 (2014). The specificity of the RAT is that it requires only one calibration (or one parameterisation)
149 covering a sufficiently-long period (at least 30 years) with as much climatic variability as possible.
150 Thus, it applies at the same time to simple conceptual models that can be calibrated automatically,
151 to more complex models requiring expert calibration, and to uncalibrated models for which

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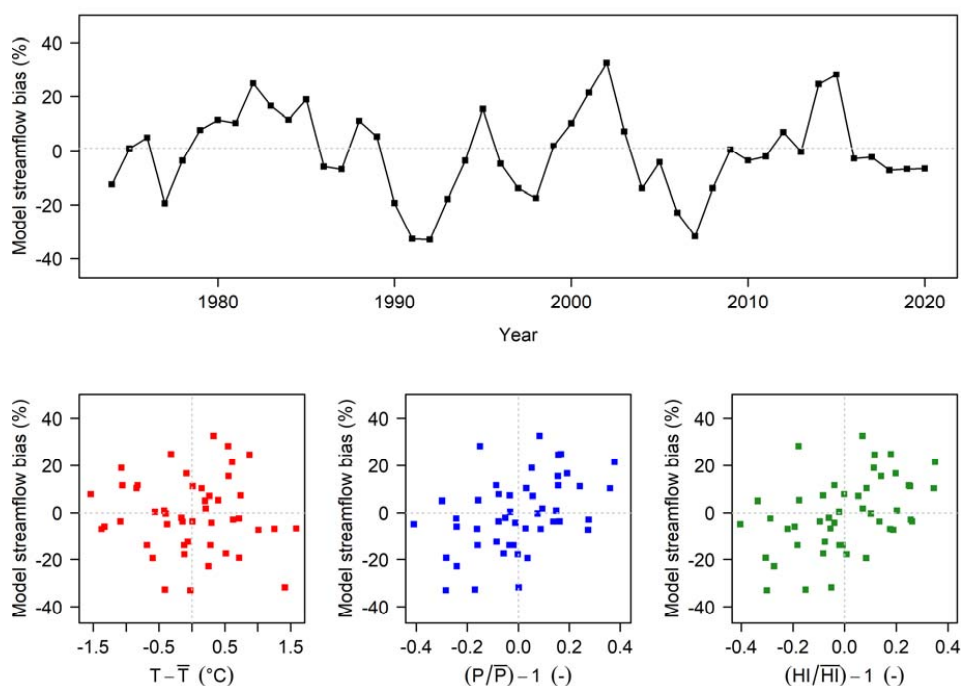158 the concept.

159

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

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

Hydrology and
Earth System
Sciences
Discussions

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



172

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

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

185   Last, we would like to mention that the RAT procedure is different from the Proxy metric for Model
186   Robustness (PMR) presented by Royer-Gaspard et al. (2021), even if both methods aim to evaluate
187   hydrological model robustness without employing a multiple calibrations process: the PMR is a
188   simple metric to estimate the robustness of a hydrological model, while the RAT is a method to
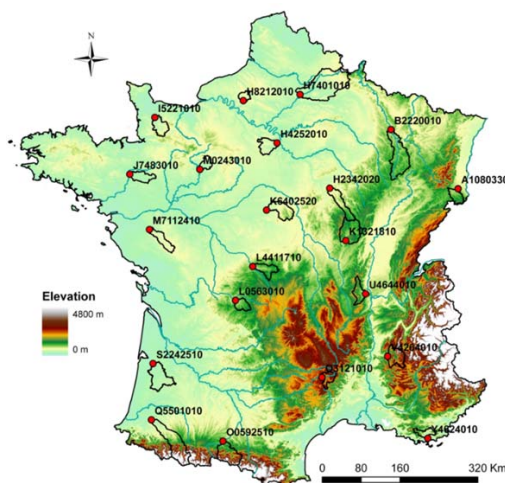
189 diagnose the dependencies of model errors to certain types of climatic changes. Thus, the RAT and
190 the PMR may be seen as complementary tools to assess a variety of aspects about model robustness.

## 3   Material and methods

### 3.1   Catchment set

193 We employed the dataset previously used by Nicolle et al. (2014), comprising 21 French catchments
194 (Figure 3), with complementary data until 2020. Catchments were chosen to represent a large range
195 of physical and climatic conditions in France, with sufficiently-long observation time series (daily
196 streamflow from 1974 to 2020) in order to provide a diverse representation of past hydroclimatic
197 conditions. Streamflow data come from the French HYDRO database (Leleu et al., 2014) and with
198 quality control performed by the operational hydrometric services. Catchment size ranges from 380
199 to 4,300 km² and median elevation from 70 to 1020 m.

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



204

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

### 3.2   Hydrological model

207 The RAT diagnostic framework is generic and can be applied to any type of model. Here daily
208 streamflow was simulated using the daily lumped GR4J rainfall–runoff model (Perrin et al., 2003).
209 The objective function used for calibration is the KGE criterion (Gupta et al., 2009) computed on
210 square-root-transformed flows. Model implementation was done with the airGR R package (Coron et
211 al., 2017, 2018).

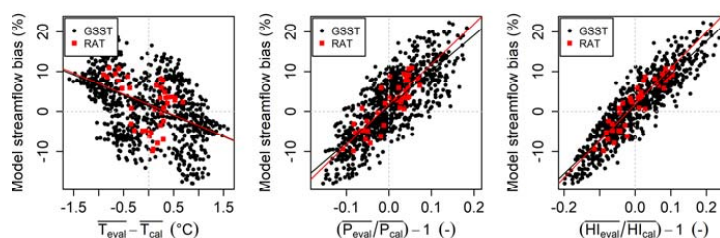212 **3.3 Evaluation of the RAT framework**

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

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

227 # 4 Results

228 ## 4.1 Comparison between the RAT and the GSST procedure

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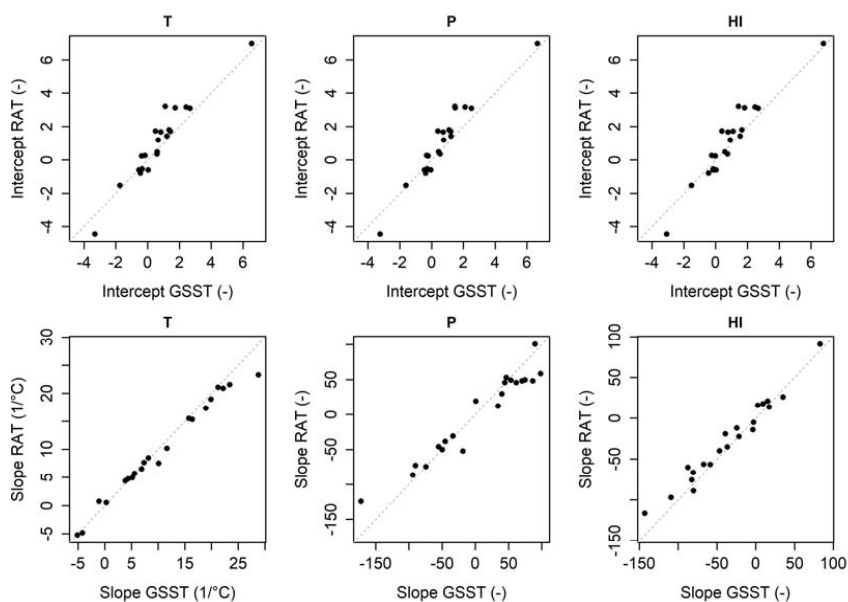236 catchments are provided as supplementary material.

237


238 **Figure 4. Streamflow bias obtained with the RAT (red squares) and the GSST (black dots), as a function of**
239 **temperature, precipitation and humidity index anomalies, for the Orge River at Morsang-sur-Orge**
240 **(H4252010) (934 km²).**

241 To summarize the results on the 21 catchments, we present on Figure 5 the slope and intercept of a
242 linear regression computed between model streamflow bias and climatic variable anomaly, for the

243    GSST and the RAT over the 21 catchments: the slope of the regressions obtained for both methods
244    are very similar and the intercept also exhibits a good match (although somewhat larger differences).
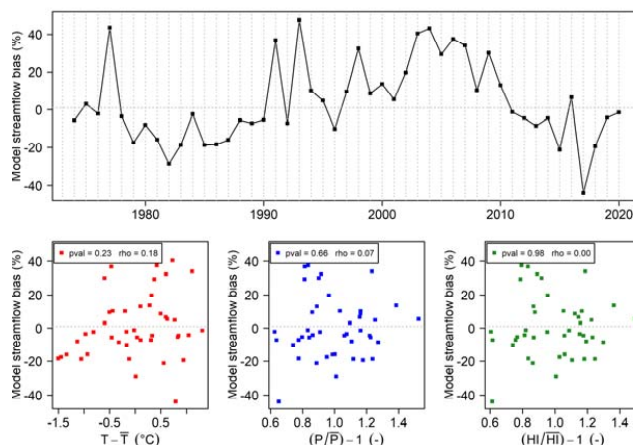


245

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

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

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

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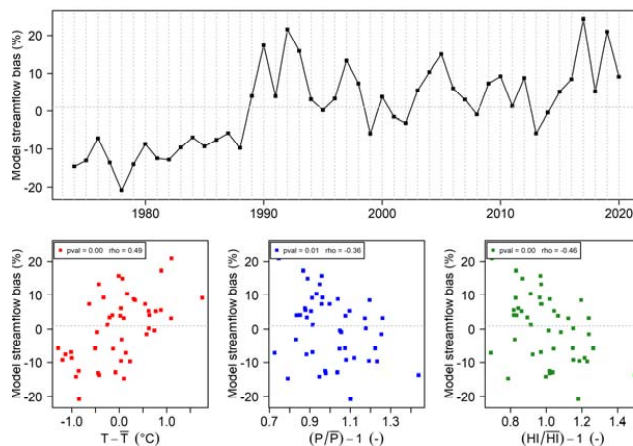258    Spearman correlation and a 5% threshold. Five cases were identified:

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

**Figure 6. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute anomalies (bottom left), and precipitation P (bottom centre) and humidity index P/E$_0$ (bottom right) anomalies, for the Orne Saosnoise River at Montbizot (M0243010) (510 km²).**

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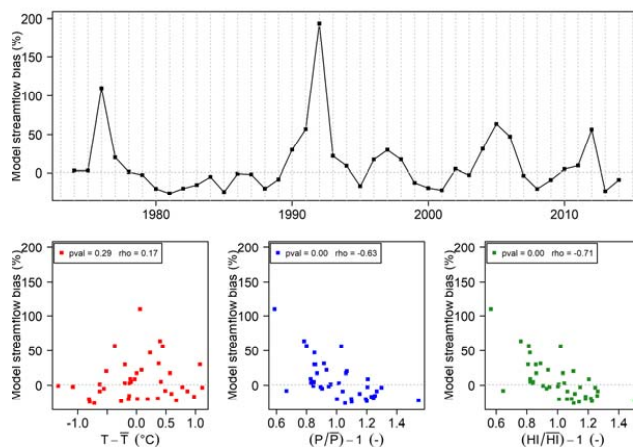


**Figure 7. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute anomalies (bottom left), and precipitation P (bottom center) and humidity index P/E$_0$ (bottom right) anomalies, for the Arroux River at Etang-sur-Arroux (K1321810) (1790 km²)**
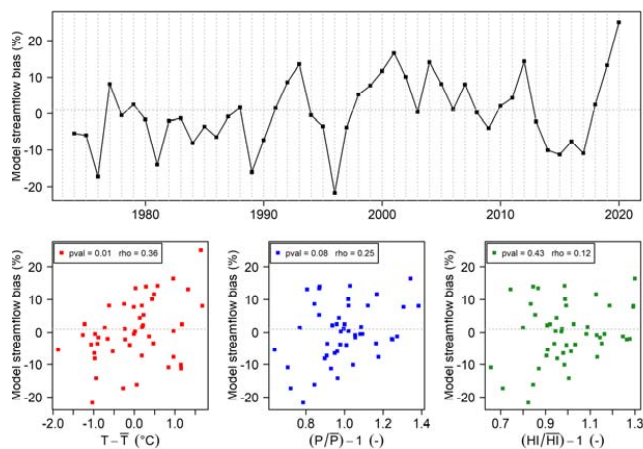
3. **Significant climate dependency on precipitation and humidity index but not on temperature** (Figure 8). This case happens on 5 of the 21 catchments.



**Figure 8.** Streamflow annual bias obtained with the RAT function of time (top), temperature absolute anomalies (bottom left), and precipitation P (bottom center) and humidity index $P/E_0$ (bottom right) anomalies, for the Seiche River at Bruz (J7483010) (810 km²)
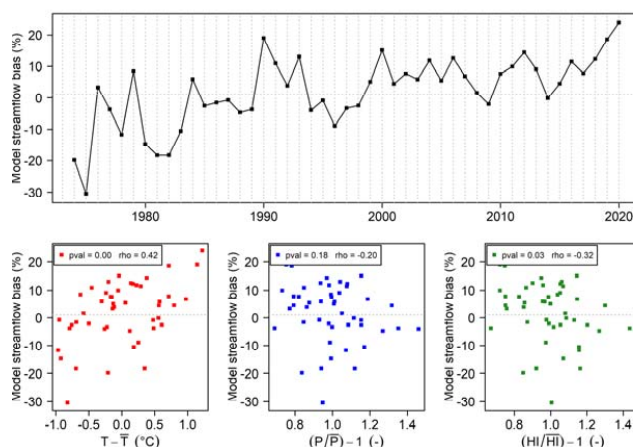
4. **Significant climate dependency on temperature but not on precipitation and humidity index** (Figure 9). This case happens on 3 of the 21 catchments.



**Figure 9.** Streamflow annual bias obtained with the RAT function of time (top), temperature absolute changes (bottom left), and precipitation P (bottom center) and humidity index $P/E_0$ (bottom right) anomalies, for the Ill at Didenheim (A1080330) (670 km²)

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

**Figure 10. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute changes (bottom left), and precipitation P (bottom center) and humidity index $P/E_0$ (bottom right) anomalies, for the Briance River at Condat-sur-Vienne (L0563010) (597 km²)**

## 5   Conclusion

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

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

1.   In this note, the RAT concept was illustrated with a rank-based test (Spearman correlation) and a significance threshold of 0.05. Like all thresholds, this one is arbitrary. Moreover, other non-parametric tests could be used and would probably yield slightly different results (we also tested the Kendall tau test, with very similar results, but do not show the results here);
2.   Detecting a relationship between model bias and a climate variable using the RAT does not allow to directly conclude on a lack of model robustness. Indeed, changes in the precipitation monitoring network or in the hydrometric rating curves can also give the false impression that the hydrological model lacks robustness. Such an erroneous conclusion could also be due to

323      widespread changes in land use, construction of an unaccounted storage reservoir or the
324      evolution of water uses. Some of the lacks of robustness detected among the 21 catchments
325      presented here could be in fact due to metrological causes;

326   3.   Also, because of the ongoing rise of temperatures (over the last 40 years at least), we have a
327      correlation between temperature and time since the beginning of streamgaging. If for any
328      reason, time is having an impact on model bias, this may cause an artefact in the RAT in the
329      form of a dependency between model bias and temperature;

330   4.   Similarly to the Differential Split Sample Test, the diagnostic of model climatic robustness is
331      limited to the climatic variable against which the bias is compared. As such, the RAT should not
332      be seen as an *absolute* test, but rather as a *necessary but not sufficient* condition to use a
333      model for climate change studies: because the climatic variability present in the past
334      observations is limited to the historic range, so is the extrapolation test. With Popper's words
335      (Popper, 1959), the RAT can only allow falsifying a hydrological model… but not proving it true;

336   5.   Although it would be tempting to transform the RAT into a post-processing method, we do not
337      recommend it. Indeed, detecting a relationship between model bias and a climate variable
338      using the RAT does not necessarily mean that a simple (linear) debiasing solution can be
339      proposed to solve the issue (see e.g. the paper by Bellprat et al. (2013) on this topic). What we
340      do recommend is to work as much as possible on the model structure, to turn it less climate
341      dependent;

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

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

## 6   Acknowledgments

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

# 7    References

366    Andréassian, V., Le Moine, N., Perrin, C., Ramos, M.-H., Oudin, L., Mathevet, T., et al. (2012). All that
367         glitters is not gold: the case of calibrating hydrological models. Hydrological Processes, 26(14),
368         2206–2210. https://doi.org/10.1002/hyp.9264

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

371    Bellprat, O., Kotlarski, S., Lüthi, D., & Schär, C. (2013). Physical constraints for temperature biases in
372         climate models: limits of temperature biases. Geophysical Research Letters, 40(15), 4042–
373         4047. https://doi.org/10.1002/grl.50737

374    Beven, K. (2016). Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood,
375         hypothesis testing, and communication. Hydrological Sciences Journal, 61(9), 1652–1665.
376         https://doi.org/10.1080/02626667.2015.1031761

377    Bisselink, B., Zambrano-Bigiarini, M., Burek, P., & de Roo, A. (2016). Assessing the role of uncertain
378         precipitation estimates on the robustness of hydrological model parameters under highly
379         variable climate conditions. Journal of Hydrology: Regional Studies, 8, 112–129.
380         https://doi.org/10.1016/j.ejrh.2016.09.003

381    Blöschl, G., et al. 2019. Twenty-three Unsolved Problems in Hydrology – a community perspective.
382         Hydrological Sciences Journal, https://doi.org/10.1080/02626667.2019.1620507

383    Brigode, P., et al. (2015). Dependence of model-based extreme flood estimation on the calibration
384         period: case study of the Kamp River (Austria). Hydrological Sciences Journal, 60 (7–8).
385         https://doi.org/10.1080/02626667.2015.1006632

386    Broderick, C., Matthews, T., Wilby, R. L., Bastola, S., & Murphy, C. (2016). Transferability of
387         hydrological models and ensemble averaging methods between contrasting climatic periods.
388         Water Resources Research, 52(10), 8343–8373. https://doi.org/10.1002/2016WR018850

389    Coron, L., Andréassian, V., Bourqui, M., Perrin, C., & Hendrickx, F. (2011). Pathologies of hydrological
390         models used in changing climatic conditions: a review. In S. Franks, E. Boegh, E. Blyth, D.
391         Hannah, & K. K. Yilmaz (Eds.), Hydro-climatology: Variability and Change. IAHS Red Books
392         Series 344 (pp. 39–44). Wallingford: IAHS.

393    Coron, L., Andréassian, V., Perrin, C., Lerat, J., Vaze, J., Bourqui, M., & Hendrickx, F. (2012). Crash
394         testing hydrological models in contrasted climate conditions: An experiment on 216 Australian
395         catchments. Water Resources Research, 48, W05552. https://doi.org/10.1029/2011WR011721

396    Coron, L., Andréassian, V., Perrin, C., Bourqui, M., & Hendrickx, F. (2014). On the lack of robustness of
397         hydrologic models regarding water balance simulation: a diagnostic approach applied to three
398         models of increasing complexity on 20 mountainous catchments. Hydrol. Earth Syst. Sci., 18(2),
399         727–746.

400    Coron, L., Thirel, G., Delaigue, O., Perrin, C., & Andréassian, V. (2017). The Suite of Lumped GR
401         Hydrological Models in an R package. Environmental Modelling and Software, 94, 337.
402         https://doi.org/10.1016/j.envsoft.2017.05.002

403    Coron, L., Delaigue, O., Thirel, G., Perrin, C. and Michel, C. (2020). airGR: Suite of GR Hydrological
404         Models for Precipitation-Runoff Modelling. R package version 1.4.3.65. DOI:
405         10.15454/EX11NA. URL: https://CRAN.R-project.org/package=airGR.

406  Dakhlaoui, H., Ruelland, D., Tramblay, Y., & Bargaoui, Z. (2017). Evaluating the robustness of
407      conceptual rainfall-runoff models under climate variability in northern Tunisia. Journal of
408      Hydrology, 550, 201–217. https://doi.org/10.1016/j.jhydrol.2017.04.032
409  Dakhlaoui, H, Ruelland, D., & Tramblay, Y. (2019). A bootstrap-based differential split-sample test to
410      assess the transferability of conceptual rainfall-runoff models under past and future climate
411      variability. Journal of Hydrology, 575, 470–486. https://doi.org/10.1016/j.jhydrol.2019.05.056
412  Delaigue, O., Génot, Lebecherel, L., Brigode, P., & Bourgin, P. Y. (2018). Base de données
413      hydroclimatiques observées à l'échelle de la France. IRSTEA. IRSTEA, UR HYCAR, Équipe
414      Hydrologie    des    bassins    versants,    Antony.    Retrieved    from
415      https://webgr.inrae.fr/en/activities/database-1-2/
416  Donelly-Makowecki, L. M., & Moore, R. D. (1999). Hierarchical testing of three rainfall-runoff models
417      in small forested catchments. Journal of Hydrology, 219(3–4), 136–152.
418  Efstratiadis, A., Nalbantis, I., and Koutsoyiannis, D. (2015). Hydrological modelling of temporally-
419      varying catchments: facets of change and the value of information. Hydrological Sciences
420      Journal, 60 (7–8). https://doi.org/10.1080/02626667.2014.982123
421  Fowler, K., Coxon, G., Freer, J., Peel, M., Wagener, T., Western, A., et al. (2018). Simulating Runoff
422      Under Changing Climatic Conditions: A Framework for Model Improvement. Water Resources
423      Research, 54(12), 9812–9832. https://doi.org/10.1029/2018WR023989
424  Gaborit, É., Ricard, S., Lachance-Cloutier, S., Anctil, F., & Turcotte, R. (2015). Comparing global and
425      local calibration schemes from a differential split-sample test perspective. Canadian Journal of
426      Earth Sciences, 52(11), 990–999. https://doi.org/10.1139/cjes-2015-0015
427  Gelfan, A.N., & Millionshchikova, T.D. (2018). Validation of a Hydrological Model Intended for Impact
428      Study: Problem Statement and Solution Example for Selenga River Basin. Water Resour 45, 90–
429      101. https://doi.org/10.1134/S0097807818050354
430  Gelfan, A., et al. (2015). Testing robustness of the physically-based ECOMAG model with respect to
431      changing    conditions.    Hydrological    Sciences    Journal,    60    (7–8).
432      https://doi.org/10.1080/02626667.2014.935780
433  Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). Decomposition of the mean squared
434      error and NSE performance criteria: Implications for improving hydrological modelling. Journal
435      of Hydrology, 377(1–2), 80–91. https://doi.org/10.1016/j.jhydrol.2009.08.003
436  Klemeš, V. (1986). Operational testing of hydrologic simulation models. Hydrological Sciences
437      Journal, 31(1), 13–24. https://doi.org/10.1080/02626668609491024
438  Krysanova, V., Donnelly, C., Gelfan, A., Gerten, D., Arheimer, B., Hattermann, F., & Kundzewicz, Z. W.
439      (2018). How the performance of hydrological models relates to credibility of projections under
440      climate    change.    Hydrological    Sciences    Journal,    63(5),    696–720.
441      https://doi.org/10.1080/02626667.2018.1446214
442  Hughes, D.A. (2015). Simulating temporal variability in catchment response using a monthly rainfall–
443      runoff    model.    Hydrological    Sciences    Journal,    60    (7–8).
444      https://doi.org/10.1080/02626667.2014.909598
445  Kling, H., et al. (2015). Performance of the COSERO precipitation–runoff model under non-stationary
446      conditions in basins with different climates. Hydrological Sciences Journal, 60 (7–8).
447      https://doi.org/10.1080/02626667.2014.959956
448  Larson, S. C. (1931). The shrinkage of the coefficient of multiple correlation. Journal of Educational
449      Psychology, 22(1), 45–55. https://doi.org/10.1037/h0072400

450    Leleu, I., Tonnelier, I., Puechberty, R., Gouin, P., Viquendi, I., Cobos, L., et al. (2014). La refonte du
451          système d'information national pour la gestion et la mise à disposition des données
452          hydrométriques. La Houille Blanche, (1), 25–32. https://doi.org/10.1051/lhb/2014004
453    Li, C. Z., Zhang, L., Wang, H., Zhang, Y. Q., Yu, F. L., and Yan, D. H. (2012). The transferability of
454          hydrological models under nonstationary climatic conditions, Hydrol. Earth Syst. Sci., 16, 1239–
455          1254, https://doi.org/10.5194/hess-16-1239-2012
456    Li, H., Beldring, S., and Xu, C.-Y. (2015). Stability of model performance and parameter values on two
457          catchments facing changes in climatic conditions. Hydrological Sciences Journal, 60 (7–8).
458          https://doi.org/10.1080/02626667.2014.978333
459    Magand, C., et al. (2015). Parameter transferability under changing climate: case study with a land
460          surface model in the Durance watershed, France. Hydrological Science Journal, 60 (7–8).
461          https://doi.org/10.1080/02626667.2014.993643
462    Montanari, A., Young, G., Savenije, H. H. G., Hughes, D., Wagener, T., Ren, L. L., et al. (2013). "Panta
463          Rhei—Everything Flows": Change in hydrology and society—The IAHS Scientific Decade 2013–
464          2022.        Hydrological         Sciences         Journal,        58(6),        1256–1275.
465          https://doi.org/10.1080/02626667.2013.809088
466    Mosteller, F., & Tukey, J. W. (1968). Data Analysis, Including Statistics. The Collected Works of John
467          W. Tukey (1988) Graphics pp. 1965-1985, vol. 5 (123)
468    Motavita, D.F., R. Chow, A. Guthke & W. Nowak. (2019). The comprehensive differential split-sample
469          test: A stress-test for hydrological model robustness under climate variability, Journal of
470          Hydrology, 573: 501-515, https://doi.org/10.1016/j.jhydrol.2019.03.054
471    Nicolle, P., Pushpalatha, R., Perrin, C., François, D., Thiéry, D., Mathevet, T., et al. (2014).
472          Benchmarking hydrological models for low-flow simulation and forecasting on French
473          catchments. Hydrol. Earth Syst. Sci., 18(8), 2829–2857. https://doi.org/10.5194/hess-18-2829-
474          2014
475    Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andréassian, V., Anctil, F., & Loumagne, C. (2005). Which
476          potential evapotranspiration input for a lumped rainfall–runoff model?: Part 2-Towards a
477          simple and efficient potential evapotranspiration model for rainfall–runoff modelling. Journal
478          of Hydrology, 303(1–4), 290–306. https://doi.org/10.1016/j.jhydrol.2004.08.026
479    Perrin, C., Michel, C., & Andréassian, V. (2003). Improvement of a parsimonious model for
480          streamflow        simulation.        Journal        of        Hydrology,        279(1–4),        275–289.
481          https://doi.org/10.1016/S0022-1694(03)00225-7
482    Popper, K. (1959). The logic of scientific discovery. London: Routldege.
483    Rau, P., Bourrel, L., Labat, D., Ruelland, D., Frappart, F., Lavado, W., et al. (2019). Assessing
484          multidecadal runoff (1970–2010) using regional hydrological modelling under data and water
485          scarcity conditions in Peruvian Pacific catchments. Hydrological Processes, 33(1), 20–35.
486          https://doi.org/10.1002/hyp.13318
487    Refsgaard, J. C., & Henriksen, H. J. (2004). Modelling guidelines—terminology and guiding principles.
488          Advances in Water Resources, 27, 71–82. https://doi.org/10.1016/j.advwatres.2003.08.006
489    Refsgaard, J. C., & Knudsen, J. (1996). Operational validation and intercomparison of different types
490          of        hydrological        models.        Water        Resources        Research,        32(7),        2189–2202.
491          https://doi.org/10.1029/96WR00896
492    Refsgaard, J. C., Madsen, H., Andréassian, V., Arnbjerg-Nielsen, K., Davidson, T. A., Drews, M., et al.
493          (2013). A framework for testing the ability of models to project climate change and its impacts.
494          Climatic Change, 122(1–2), 271–282. https://doi.org/10.1007/s10584-013-0990-2

495    Roberts, D. R., Bahn, V., Ciuti, S., Boyce, M. S., Elith, J., Guillera-Arroita, G., et al. (2017). Cross-
496        validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure.
497        Ecography, 40(8), 913–929. https://doi.org/10.1111/ecog.02881
498    Royer-Gaspard, P., Andréassian, V., and Thirel, G. (2021) Technical note: PMR − a proxy metric to
499        assess hydrological model robustness in a changing climate. In review for Hydrol. Earth Syst.
500        Sci. Discuss., https://doi.org/10.5194/hess-2021-58
501    Seibert, J. (2003). Reliability of model predictions outside calibration conditions. Nordic Hydrology,
502        34(5), 477–492. https://doi.org/10.2166/nh.2003.0019
503    Seifert, D., Sonnenborg, T. O., Refsgaard, J. C., Højberg, A. L., and Troldborg, L. (2012). Assessment of
504        hydrological model predictive ability given multiple conceptual geological models, Water
505        Resour. Res., 48, W06503, doi:10.1029/2011WR011149.
506    Seiller, G., Anctil, F., & Perrin, C. (2012). Multimodel evaluation of twenty lumped hydrological
507        models under contrasted climate conditions. Hydrology and Earth System Sciences, 16(4),
508        1171–1189. https://doi.org/10.5194/hess-16-1171-2012
509    Tanaka, T. and Tachikawa, Y. (2015). Testing the applicability of a kinematic wave-based distributed
510        hydrological model in two climatically contrasting catchments. Hydrological Sciences Journal,
511        60 (7–8). https://doi.org/10.1080/02626667.2014.967693
512    Taver, V., et al. (2015). Feed-forward vs recurrent neural network models for non-stationarity
513        modelling using data assimilation and adaptivity. Hydrological Sciences Journal, 60 (7–8).
514        https://doi.org/10.1080/02626667.2014.967696
515    Teutschbein, C. & Seibert, J. 2013. Is bias correction of regional climate model (RCM) simulations
516        possible for non-stationary conditions? Hydrology and Earth System Sciences, 17, 5061–5077.,
517        https://doi.org/10.5194/hess-17-5061-2013
518    Thirel, G., Andréassian, V., Perrin, C., Audouy, J.-N., Berthet, L., Edwards, P., et al. (2015a). Hydrology
519        under change: an evaluation protocol to investigate how hydrological models deal with
520        changing catchments. Hydrological Sciences Journal, 60(7–8), 1184–1199.
521        https://doi.org/10.1080/02626667.2014.967248
522    Thirel, G., Andréassian, V., & Perrin, C. (2015b). On the need to test hydrological models under
523        changing conditions. Hydrological Sciences Journal, 60(7–8), 1165–1173.
524        https://doi.org/10.1080/02626667.2015.1050027
525    Vaze, J., Post, D. A., Chiew, F. H. S., Perraud, J. M., Viney, N. R., & Teng, J. (2010). Climate non-
526        stationarity - Validity of calibrated rainfall-runoff models for use in climate change studies.
527        Journal of Hydrology, 394(3–4), 447–457. https://doi.org/10.1016/j.jhydrol.2010.09.018
528    Vidal, J.-P., Martin, E., Franchistéguy, L., Baillon, M., & Soubeyroux, J.-M. (2010). A 50-year high-
529        resolution atmospheric reanalysis over France with the Safran system. International Journal of
530        Climatology, 30(11), P. 1627–1644. DOI: 10.1002/joc.2003. https://doi.org/10.1002/joc.2003
531    Vormoor, K., Heistermann, M., Bronstert, A., & Lawrence, D. (2018). Hydrological model parameter
532        (in)stability – "crash testing" the HBV model under contrasting flood seasonality conditions.
533        Hydrological Sciences Journal, 63(7), 991–1007.
534        https://doi.org/10.1080/02626667.2018.1466056
535    Wilby, R. L. (2019). A global hydrology research agenda fit for the 2030s. Hydrology Research, 50(6):
536        1464-1480. https://doi.org/10.2166/nh.2019.100
537    Xu, C. (1999). Climate Change and Hydrologic Models: A Review of Existing Gaps and Recent Research
538        Developments. Water Resources Management, 13(5), 369–382.
539        https://doi.org/10.1023/A:1008190900459

540    Yu, B. and Zhu, Z. (2015). A comparative assessment of AWBM and SimHyd for forested watersheds.
541         Hydrological Sciences Journal, 60 (7–8). https://doi.org/10.1080/02626667.2014.961924

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

In this appendix, we deal with calibrated models, for which we verify that the main hypothesis underlying the RAT is reasonable, i.e. that when considering a long calibration period, the weight of each individual year in the overall calibration process is almost negligible. We then explore the limits of this hypothesis when reducing the length of the overall calibration period.
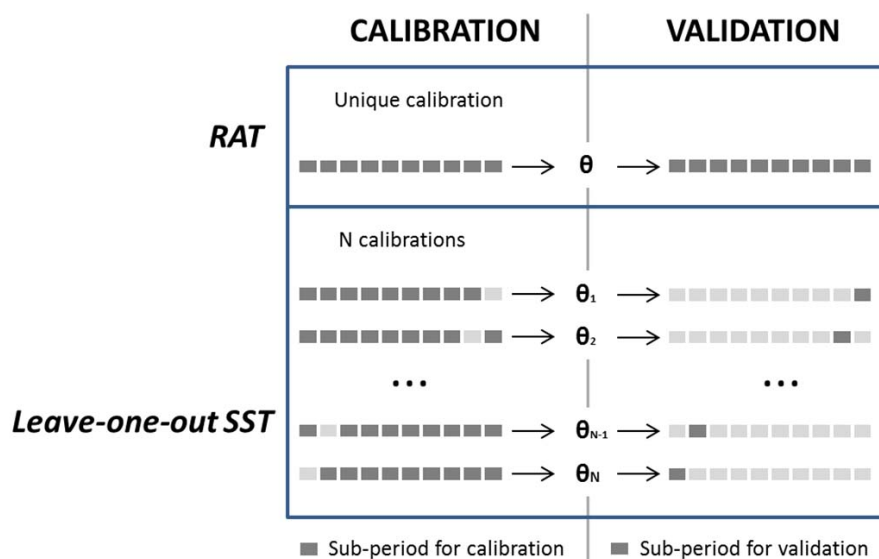
- **Evaluation method**

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

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

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

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

The difference between the two approaches is schematized in Figure 11: the leave-one-out procedure consists in performing $N$ calibrations over ($N$-$1$)-year-long periods followed by an independent evaluation on the remaining 1-year-long period. As shown in Figure 11, the two procedures result in the same number of validation points ($N$). Eq. 1 provides a way to quantify whether both methods differ, i.e. whether the partial overlap between calibration and validation periods in the RAT makes a difference.
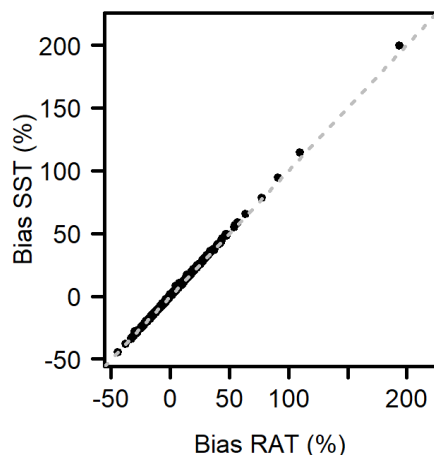
567

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

571

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

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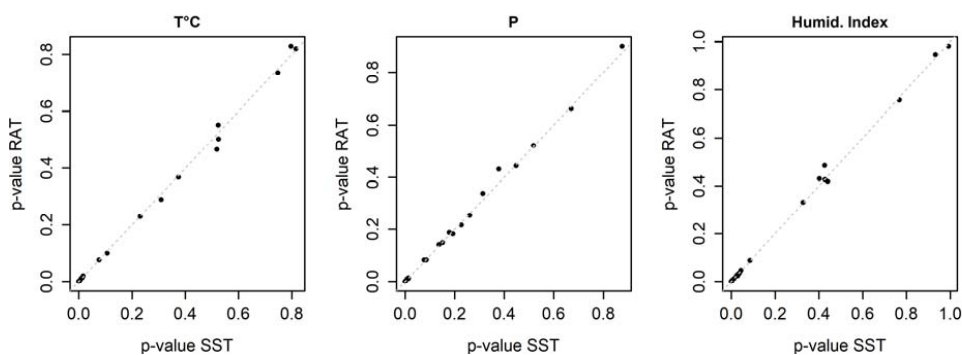576 catchment set).

577

578

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

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

586



587

**Figure 13. Spearman correlation *p*-value from the correlation for annual bias and annual temperature, precipitation, and humidity index (P/E$_0$). Comparison between RAT and SST (one point per catchment)**
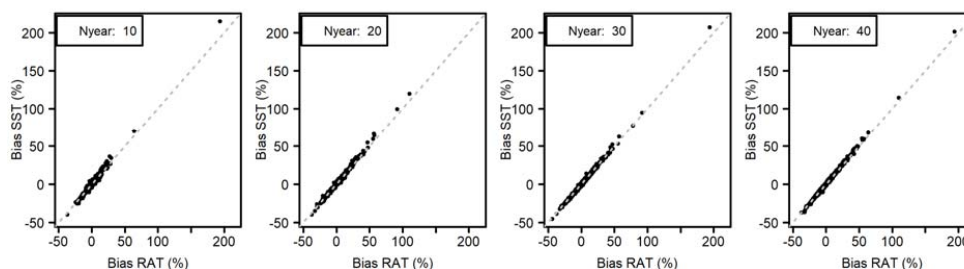
590

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

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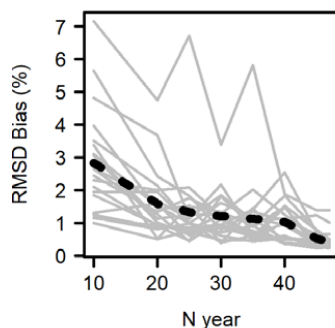600 We would not advise to use RAT with time series of less than 20 years.

601



602

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

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



610

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