

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

3
4 Pierre Nicolle^{1,3}, Vazken Andréassian^{1,*}, Paul Royer-Gaspard¹, Charles Perrin¹, Guillaume Thirel¹,
5 Laurent Coron², Léonard Santos¹

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

7 ²EDF, DTG, Toulouse, France

8 ³now 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 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 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 presents
25 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 catchments,
30 verify its applicability hypotheses and compare it to previously published tests. Results show that the
31 RAT is an efficient evaluation approach, passing it successfully can be considered a prerequisite for any
32 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 (Wilby,
36 2019). Given the expected evolution of climate conditions, the actual ability of models to predict the
37 corresponding evolution of hydrological variables should be verified (Beven, 2016). Indeed, when using
38 a hydrological model for climate change impact assessment, we make two implicit hypotheses
39 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 which
43 would result from a modification of catchment physical characteristics, we do expect them to be able
44 to represent the catchment response to extreme climatic conditions (and possibly to conditions more
45 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 which
48 are calibrated, we assume that the parameters are generic and not specific to the calibration period,
49 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 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 formalized
63 it for hydrological modelling. Klemeš proposed initially a four-level testing scheme for evaluating
64 model transposability in time and space: (i) split-sample test on two independent periods, (ii) proxy-
65 basin test on neighbouring catchments, (iii) differential split-sample test on contrasted independent
66 periods (DSST), and (iv) proxy-basin differential split-sample test on neighbouring catchments and
67 contrasted periods.

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

72 future. However, a satisfactory behaviour during the DSST can be seen as a prerequisite of model
73 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 which
77 extent Klemeš's hierarchical testing scheme could be used to improve the conclusions of model
78 intercomparisons. The study by Xu (1999) questioned the applicability of models in nonstationary
79 conditions and was one of the early attempts to apply the Klemeš's testing scheme in this perspective.
80 Similarly, tests carried out by Seibert (2003) explicitly intended to test the ability of a model to
81 extrapolate beyond calibration range and showed limitations of the tested model, stressing the need
82 for improved calibration strategies. Last, Vaze et al. (2010) also investigated the behaviour of four
83 rainfall-runoff models under contrasting conditions, using wet and dry periods on catchments in
84 Australia that experienced a prolonged drought period. They observed different model behaviours
85 when going from wet to dry or dry to wet conditions.

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

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

109 A few authors also tried to propose improved implementations of these testing schemes. Seiller et al.
110 (2012) used non-continuous periods or years selected on mean temperature and precipitation to
111 enhance the contrast between testing periods. This idea to jointly use these two climate variables to
112 select periods was further investigated by Gaborit et al. (2015), who assessed how the temporal model

113 robustness can be improved by advanced calibration schemes. They showed that the robustness of
114 the tested model was improved when going from humid-cold to dry-warm or from dry-cold to humid-
115 warm conditions when using regional calibration instead of local calibration. Dakhlaoui et al. (2017)
116 investigated the impact of DSST on model robustness by selecting dry/wet and cold/hot hydrological
117 years to increase the contrast in climate conditions between calibration and validation periods. These
118 authors later proposed a bootstrap technique to widen the testing conditions (Dakhlaoui et al. 2019).
119 The investigations of Fowler et al. (2018) identified some limits of the DSST procedure and concluded
120 that “model evaluation based solely on the DSST is hampered due to contingency on the chosen
121 calibration method, and it is difficult to distinguish which cases of DSST failure are truly caused by
122 model structural inadequacy”. Last, Motavita et al. (2019) combined DSST with periods of variable
123 length, and conclude that parameters obtained on dry periods may be more robust.

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

126 **1.3 Scope of the technical note**

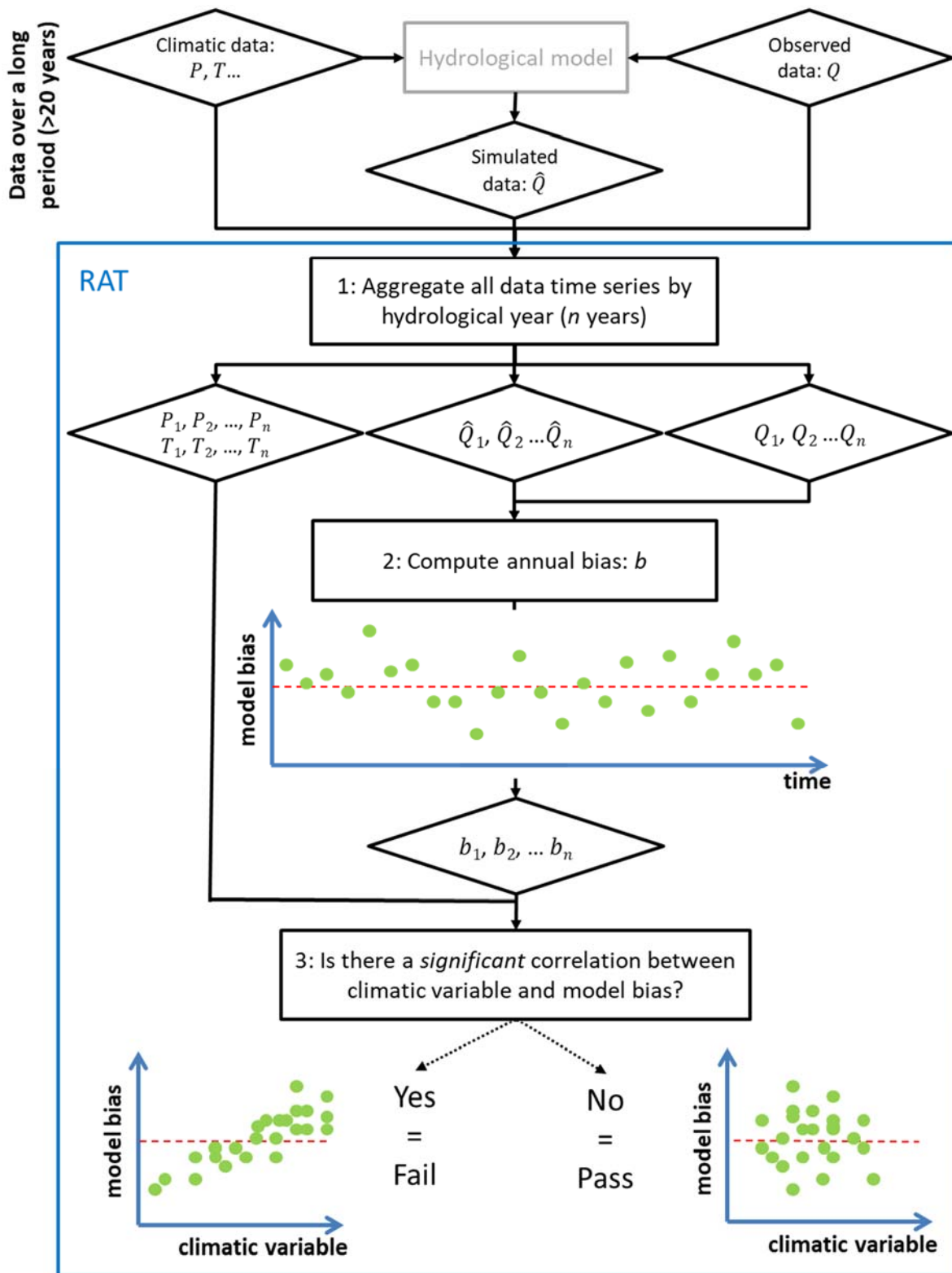
127 This note presents a new generic diagnostic framework inspired by Klemeš’s DSST procedure and by
128 our own previous attempts (Coron et al., 2012; Thirel et al., 2015a) to assess the relative confidence
129 one may have with a hydrological model to be used in a changing climate context. One of the problems
130 of existing methods is the requirement of multiple calibrations of hydrological models: these are
131 relatively easy to implement with parsimonious conceptual models but definitively not with complex
132 models that require long interventions by expert modellers and, obviously, not for those models with
133 a once-for-all parameterisation.

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

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

140 **2 The robustness assessment test (RAT) concept**

141 The robustness assessment test (RAT) proposed in this note is inspired by the work of Coron et al.
142 (2014). The specificity of the RAT is that it requires only one simulation covering a sufficiently-long
143 period (at least 20 years) with as much climatic variability as possible. Thus, it applies at the same time
144 to simple conceptual models that can be calibrated automatically, to more complex models requiring
145 expert calibration, and to uncalibrated models for which parameters are derived from the
146 measurement of certain physical properties. The RAT consists in computing a relevant numeric bias
147 criterion repeatedly each year and then exploring its correlation with a climatic factor deemed
148 meaningful, in order to identify undesirable dependencies and thus to assess the extrapolation
149 capacity (Roberts et al., 2017) of any hydrological model. Indeed, if the performances of a model are
150 shown to be dependent on a given climate variable, this can be an issue when the model is used on a
151 period with a changing climate. The flowchart in Figure 1 summarizes the concept.

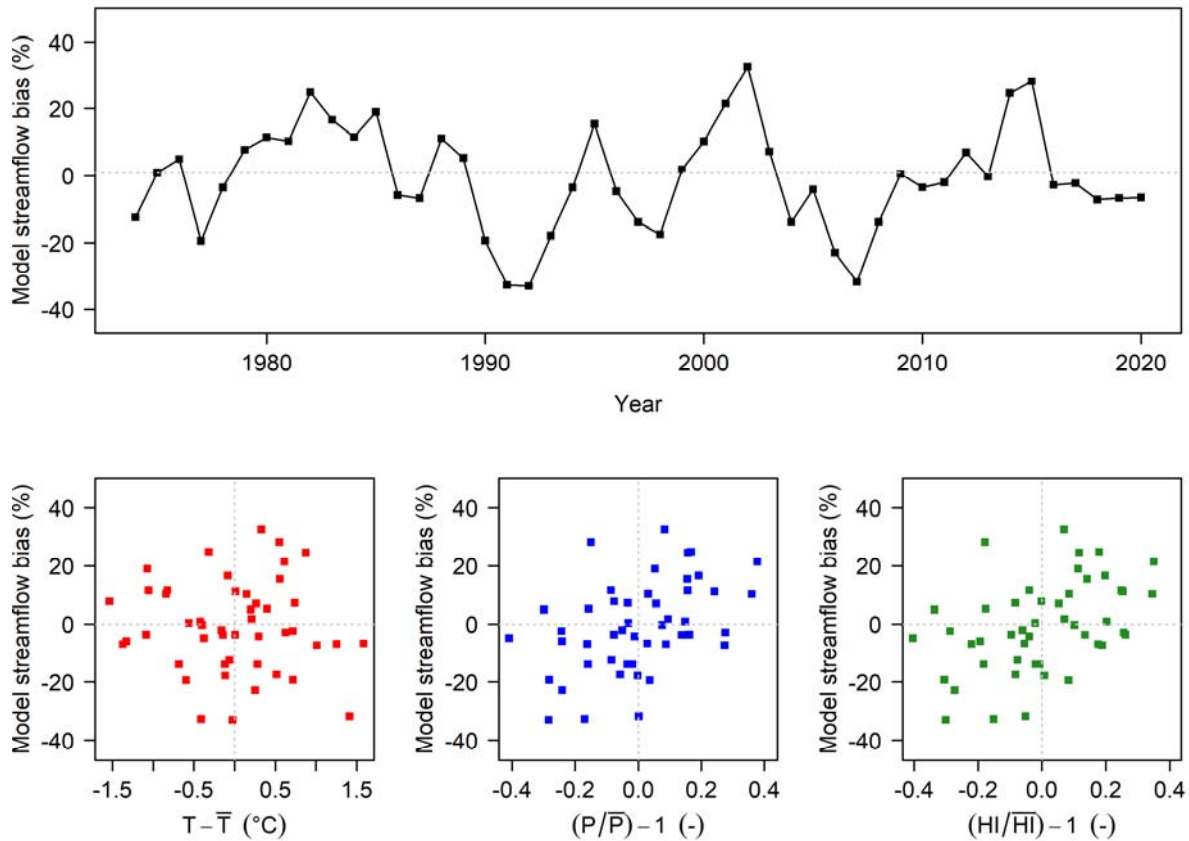


152

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

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

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



165
 166 **Figure 2. Robustness Assessment Test (RAT) applied to a hydrological model: the upper graph presents the**
 167 **evolution in time (year by year) of model streamflow bias; the lower scatterplots present the relationship**
 168 **between model bias and climatic variables (temperature T , precipitation P and humidity index HI , from left to**
 169 **right)**

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

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

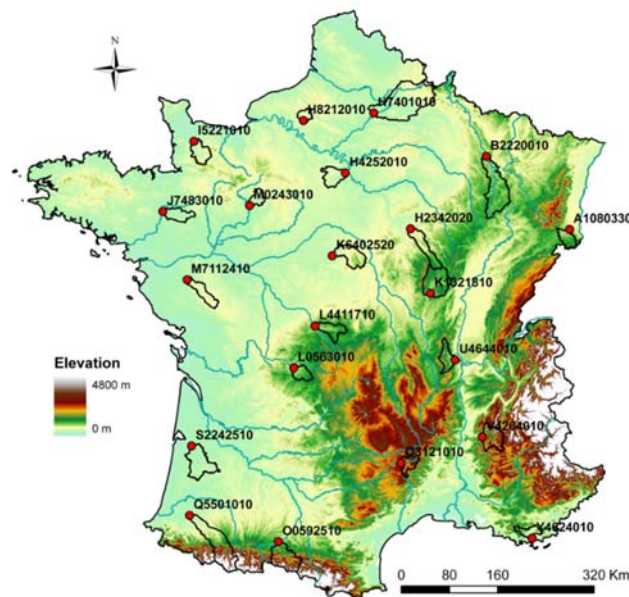
180 hydrological model robustness without employing a multiple calibrations process: the PMR is a simple
181 metric to estimate the robustness of a hydrological model, while the RAT is a method to diagnose the
182 dependencies of model errors to certain types of climatic changes. Thus, the RAT and the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

184 3 Material and methods

185 3.1 Catchment set

186 We employed the dataset previously used by Nicolle et al. (2014), comprising 21 French catchments
187 (Figure 3), extended up to 2020. Catchments were chosen to represent a large range of physical and
188 climatic conditions in France, with sufficiently-long observation time series (daily streamflow from
189 1974 to 2020) in order to provide a diverse representation of past hydroclimatic conditions.
190 Streamflow data come from the French HYDRO database (Leleu et al., 2014) and with quality control
191 performed by the operational hydrometric services. Catchment size ranges from 380 to 4,300 km² and
192 median elevation from 70 to 1020 m.

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



197
198 **Figure 3. Location of the 21 catchments in France. Red dots represent the catchment outlets**

199 3.2 Hydrological model

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

203 root-transformed flows. Model implementation was done with the airGR R package (Coron et al., 2017,
204 2018).

205 3.3 Evaluation of the RAT framework

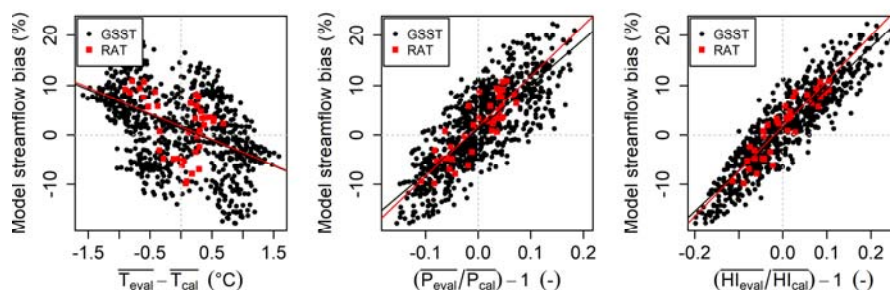
206 The RAT was evaluated against the GSST of Coron et al. (2012) used as a benchmark, in order to check
207 whether it yields similar results. The GSST procedure was applied to each catchment using a 10-year
208 period to calibrate the model. For each calibration, each 10-year sliding period over the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

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

219 4 Results

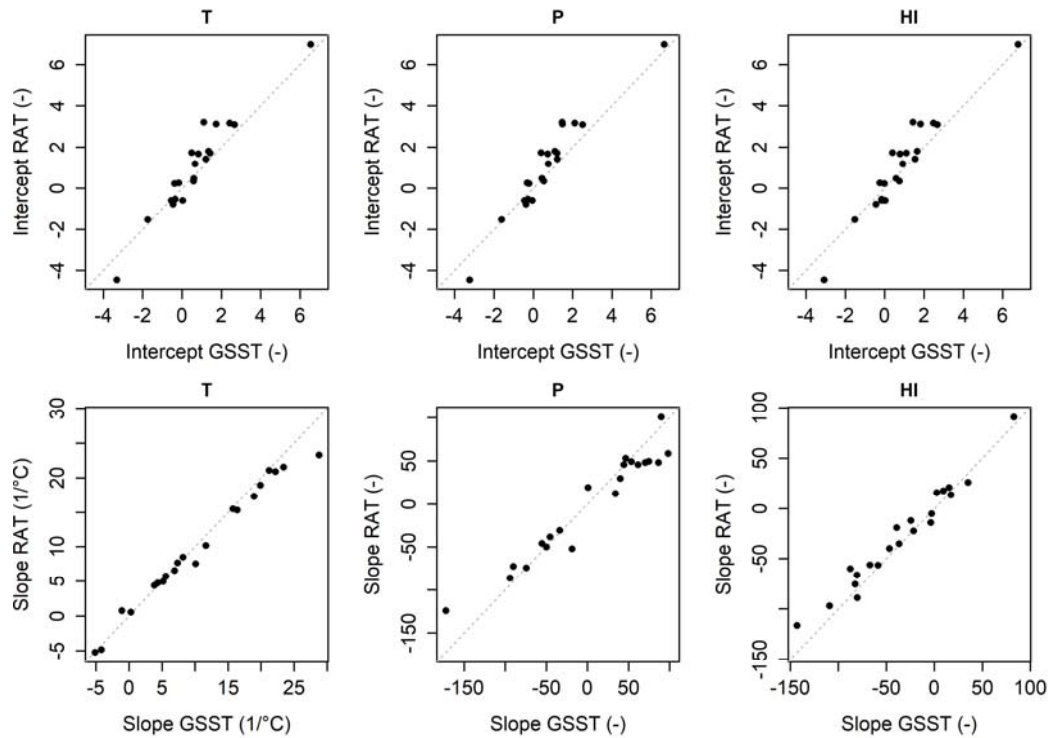
220 4.1 Comparison between the RAT and the GSST procedure

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



229
230 **Figure 4. Streamflow bias obtained with the RAT (red squares) and the GSST (black dots), as a function of**
231 **temperature, precipitation and humidity index anomalies, for the Orge River at Morsang-sur-Orge (H4252010)**
232 **(934 km²).**

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



237

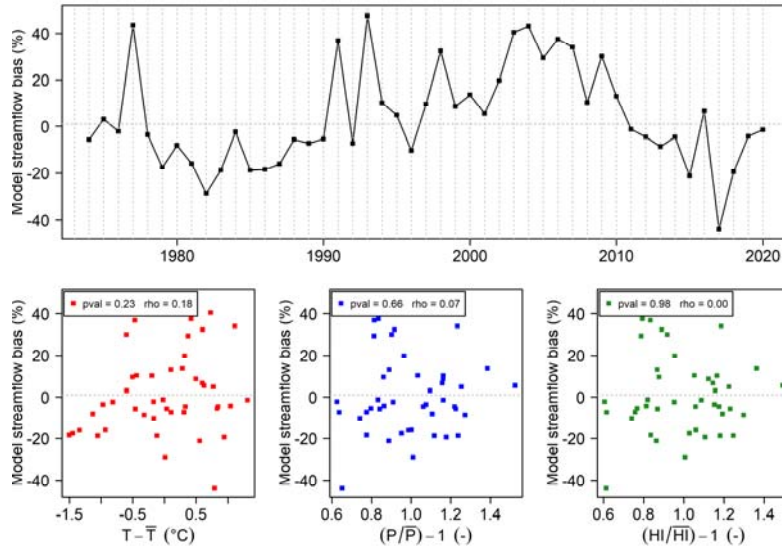
238 **Figure 5. Comparison of slopes and intercept of linear regressions between streamflow bias and temperature**
 239 **(T), precipitation (P) and humidity index (HI) anomalies (from left to right) obtained by the GSST and the RAT**
 240 **procedures (each point represents one of the 21 test catchments)**

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

247 4.2 Application of the RAT procedure to the detection of climate dependencies

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

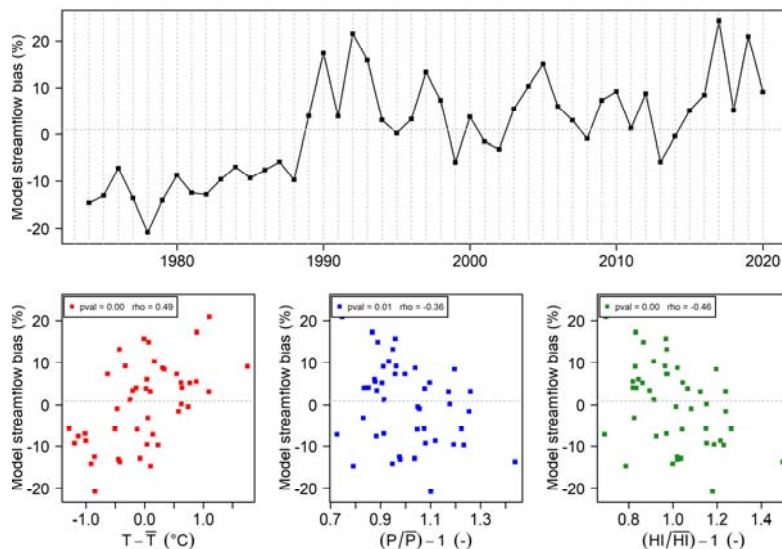
- 251 1. **No climate dependency** (Figure 6): This is the case for 6 catchments out of 21 and the
 252 expected situation of a “robust” model. The different plots show a lack of dependence, for
 253 temperature, precipitation and humidity index alike. For the catchment of Figure 6, the p-
 254 value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



255

256 **Figure 6. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute**
 257 **anomalies (bottom left), and precipitation P (bottom centre) and humidity index P/E_0 (bottom right)**
 258 **anomalies, for the Orne Saosnoise River at Montbizot (M0243010) (510 km²).**

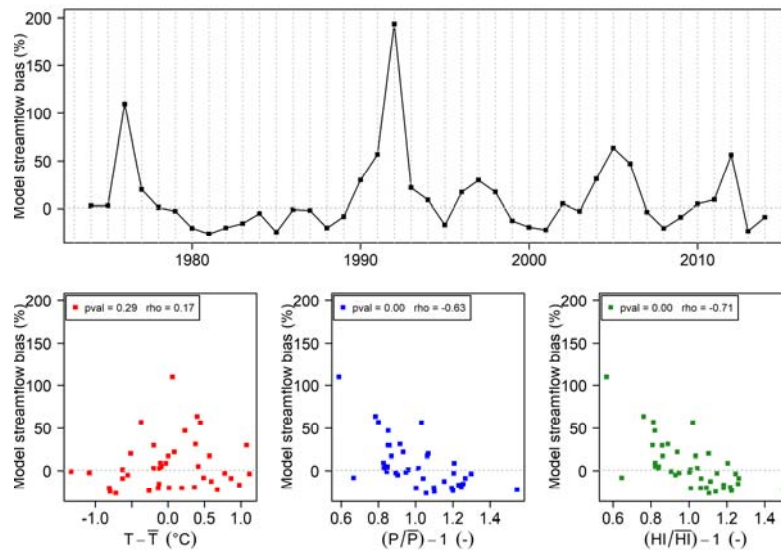
259 2. **Significant dependency on annual temperature, precipitation and humidity index** (Figure
 260 7): This is a clearly undesirable situation illustrating a lack of robustness of the hydrological
 261 model. It happens on only two catchments out of 21. The Spearman correlation between
 262 model bias and temperature, precipitation and humidity index anomalies (respectively
 263 0.49, -0.36 and -0.46) is significant (i.e. below the classic significance threshold of 5%). In
 264 Figure 7, the annual streamflow bias shows an increasing trend with annual temperature
 265 and a decreasing trend with annual precipitation and humidity index.



266

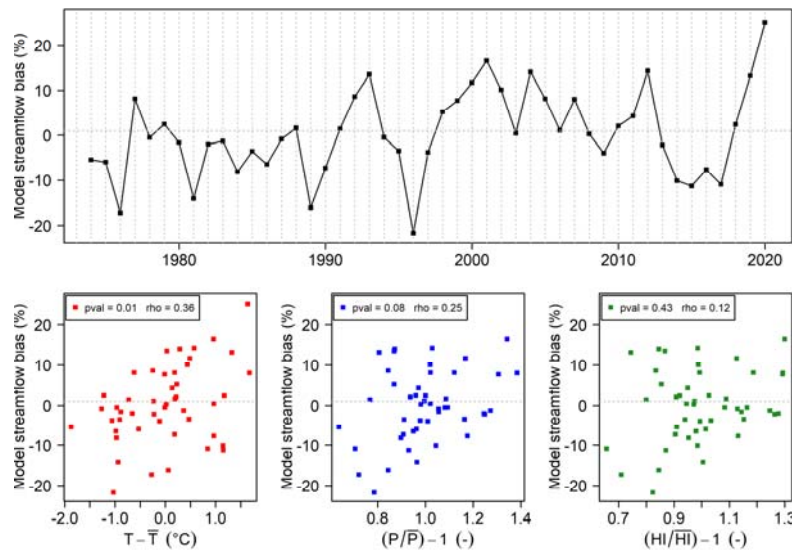
267 **Figure 7. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute**
 268 **anomalies (bottom left), and precipitation P (bottom center) and humidity index P/E_0 (bottom right)**
 269 **anomalies, for the Arroux River at Etang-sur-Arroux (K1321810) (1790 km²).**

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



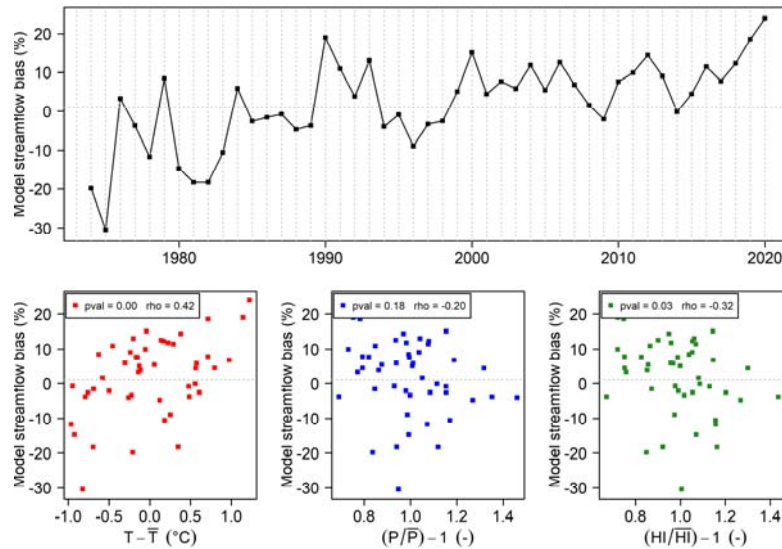
272
 273 **Figure 8. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute**
 274 **anomalies (bottom left), and precipitation P (bottom center) and humidity index P/E_0 (bottom right)**
 275 **anomalies, for the Seiche River at Bruz (J7483010) (810 km²)**

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



278
 279 **Figure 9. Streamflow annual bias obtained with the RAT function of time (top), temperature absolute changes**
 280 **(bottom left), and precipitation P (bottom center) and humidity index P/E_0 (bottom right) anomalies, for the**
 281 **Ill at Didenheim (A1080330) (670 km²)**

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



284

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

288 4.3 How to use RAT results?

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

300

301 5 Conclusion

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

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

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

355 dependency, and that evaluating the quality of input data should be done before looking at
356 model robustness;

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

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

368 **6 Acknowledgments**

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

372 The authors gratefully acknowledge the comments of Prof. Jens-Christian Refsgaard and Dr Nans Addor
373 on a preliminary version of the note, and the reviews by Dr Bettina Schäfli and by two anonymous
374 reviewers.

375 **7 Code/Data availability**

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

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

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

380 **8 Author contribution**

381 VA proposed the RAT concept based on discussions that had been going on in INRAE's HYCAR research
382 unit for more than a decade and whose origin can be traced back to the blessed era when Claude
383 Michel was providing his hydrological teaching in Antony. PN performed the computations and wrote
384 the paper with the help of all co-authors. LS proposed the summary flow chart.

385 **9 Competing interests**

386 None

387 **10 References**

- 388 Andréassian, V., Le Moine, N., Perrin, C., Ramos, M.-H., Oudin, L., Mathevet, T., Lerat, J., Berthet, L.
389 (2012). All that glitters is not gold: the case of calibrating hydrological models. *Hydrological*
390 *Processes*, 26(14), 2206–2210. <https://doi.org/10.1002/hyp.9264>
- 391 Andréassian, V., Perrin, C., Berthet, L., Le Moine, N., Lerat, J., Loumagne, C., Oudin, L., Mathevet, T.,
392 Ramos, M.H., Valéry, A. (2009). Crash tests for a standardized evaluation of hydrological models.
393 *Hydrology and Earth System Sciences*, 13, 1757-1764. [https://doi.org/10.5194/hess-13-1757-](https://doi.org/10.5194/hess-13-1757-2009)
394 2009
- 395 Arlot, S., & Celisse, A. (2010). A survey of cross-validation procedures for model selection. *Statistics*
396 *Surveys*, 4(0), 40–79. <https://doi.org/10.1214/09-SS054>
- 397 Bellprat, O., Kotlarski, S., Lüthi, D., & Schär, C. (2013). Physical constraints for temperature biases in
398 climate models: limits of temperature biases. *Geophysical Research Letters*, 40(15), 4042–4047.
399 <https://doi.org/10.1002/grl.50737>
- 400 Beven, K. (2016). Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood, hypothesis
401 testing, and communication. *Hydrological Sciences Journal*, 61(9), 1652–1665.
402 <https://doi.org/10.1080/02626667.2015.1031761>
- 403 Bisselink, B., Zambrano-Bigiarini, M., Burek, P., & de Roo, A. (2016). Assessing the role of uncertain
404 precipitation estimates on the robustness of hydrological model parameters under highly
405 variable climate conditions. *Journal of Hydrology: Regional Studies*, 8, 112–129.
406 <https://doi.org/10.1016/j.ejrh.2016.09.003>
- 407 Blöschl, G., et al. 2019. Twenty-three Unsolved Problems in Hydrology – a community perspective.
408 *Hydrological Sciences Journal*, <https://doi.org/10.1080/02626667.2019.1620507>
- 409 Brigode, P., et al. (2015). Dependence of model-based extreme flood estimation on the calibration
410 period: case study of the Kamp River (Austria). *Hydrological Sciences Journal*, 60 (7–8).
411 <https://doi.org/10.1080/02626667.2015.1006632>
- 412 Broderick, C., Matthews, T., Wilby, R. L., Bastola, S., & Murphy, C. (2016). Transferability of hydrological
413 models and ensemble averaging methods between contrasting climatic periods. *Water*
414 *Resources Research*, 52(10), 8343–8373. <https://doi.org/10.1002/2016WR018850>
- 415 Coron, L., Andréassian, V., Bourqui, M., Perrin, C., & Hendrickx, F. (2011). Pathologies of hydrological
416 models used in changing climatic conditions: a review. In S. Franks, E. Boegh, E. Blyth, D. Hannah,
417 & K. K. Yilmaz (Eds.), *Hydro-climatology: Variability and Change*. IAHS Red Books Series 344 (pp.
418 39–44). Wallingford: IAHS.
- 419 Coron, L., Andréassian, V., Perrin, C., Bourqui, M., & Hendrickx, F. (2014). On the lack of robustness of
420 hydrologic models regarding water balance simulation: a diagnostic approach applied to three
421 models of increasing complexity on 20 mountainous catchments. *Hydrol. Earth Syst. Sci.*, 18(2),
422 727–746.
- 423 Coron, L., Andréassian, V., Perrin, C., Lerat, J., Vaze, J., Bourqui, M., & Hendrickx, F. (2012). Crash testing
424 hydrological models in contrasted climate conditions: An experiment on 216 Australian
425 catchments. *Water Resources Research*, 48, W05552. <https://doi.org/10.1029/2011WR011721>
- 426 Coron, L., Delaigue, O., Thirel, G., Perrin, C. and Michel, C. (2020). airGR: Suite of GR Hydrological
427 Models for Precipitation-Runoff Modelling. R package version 1.4.3.65. DOI: 10.15454/EX11NA.
428 URL: <https://CRAN.R-project.org/package=airGR>.

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

432 Dakhlaoui, H., Ruelland, D., & Tramblay, Y. (2019). A bootstrap-based differential split-sample test to
433 assess the transferability of conceptual rainfall-runoff models under past and future climate
434 variability. *Journal of Hydrology*, 575, 470–486. <https://doi.org/10.1016/j.jhydrol.2019.05.056>

435 Dakhlaoui, H., Ruelland, D., Tramblay, Y., & Bargaoui, Z. (2017). Evaluating the robustness of
436 conceptual rainfall-runoff models under climate variability in northern Tunisia. *Journal of*
437 *Hydrology*, 550, 201–217. <https://doi.org/10.1016/j.jhydrol.2017.04.032>

438 Delaigue, O., Génot, Lebecherel, L., Brigode, P., & Bourgin, P. Y. (2018). Base de données
439 hydroclimatiques observées à l'échelle de la France. IRSTEA. IRSTEA, UR HYCAR, Équipe
440 Hydrologie des bassins versants, Antony. Retrieved from
441 <https://webgr.inrae.fr/en/activities/database-1-2/>

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

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

447 Fowler, K., et al. (2018). Simulating Runoff Under Changing Climatic Conditions: A Framework for
448 Model Improvement. *Water Resources Research*, 54(12), 9812–9832.
449 <https://doi.org/10.1029/2018WR023989>

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

453 Gelfan, A., Motovilov, Y., Krylenko, I., Moreido, V., & Zakharova, E. (2015). Testing robustness of the
454 physically-based ECOMAG model with respect to changing conditions. *Hydrological Sciences*
455 *Journal*, 60 (7–8). <https://doi.org/10.1080/02626667.2014.935780>

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

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

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

465 Klemeš, V. (1986). Operational testing of hydrologic simulation models. *Hydrological Sciences Journal*,
466 31(1), 13–24. <https://doi.org/10.1080/02626668609491024>

467 Kling, H., Stanzel, P., Fuchs, M., & Nachtnebel, H.-P. (2015). Performance of the COSERO precipitation–
468 runoff model under non-stationary conditions in basins with different climates. *Hydrological*
469 *Sciences Journal*, 60 (7–8). <https://doi.org/10.1080/02626667.2014.959956>

470 Krysanova, V., Donnelly, C., Gelfan, A., Gerten, D., Arheimer, B., Hattermann, F., & Kundzewicz, Z. W.
471 (2018). How the performance of hydrological models relates to credibility of projections under
472 climate change. *Hydrological Sciences Journal*, 63(5), 696–720.
473 <https://doi.org/10.1080/02626667.2018.1446214>

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

476 Leleu, I., et al. (2014). La refonte du système d'information national pour la gestion et la mise à
477 disposition des données hydrométriques. *La Houille Blanche*, (1), 25–32.
478 <https://doi.org/10.1051/lhb/2014004>

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

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

485 Magand, C., Ducharne, A., Le Moine, N., & Brigode, P. (2015). Parameter transferability under changing
486 climate: case study with a land surface model in the Durance watershed, France. *Hydrological*
487 *Science Journal*, 60 (7–8). <https://doi.org/10.1080/02626667.2014.993643>

488 Montanari, A., et al. (2013). “Panta Rhei—Everything Flows”: Change in hydrology and society—The
489 IAHS Scientific Decade 2013–2022. *Hydrological Sciences Journal*, 58(6), 1256–1275.
490 <https://doi.org/10.1080/02626667.2013.809088>

491 Mosteller, F., & Tukey, J. W. (1968). *Data Analysis, Including Statistics. The Collected Works of John W.*
492 *Tukey (1988) Graphics pp. 1965-1985, vol. 5 (123)*

493 Motavita, D.F., R. Chow, A. Guthke & W. Nowak. (2019). The comprehensive differential split-sample
494 test: A stress-test for hydrological model robustness under climate variability, *Journal of*
495 *Hydrology*, 573: 501-515, <https://doi.org/10.1016/j.jhydrol.2019.03.054>

496 Nicolle, P., et al. (2014). Benchmarking hydrological models for low-flow simulation and forecasting on
497 French catchments. *Hydrol. Earth Syst. Sci.*, 18(8), 2829–2857. [https://doi.org/10.5194/hess-18-](https://doi.org/10.5194/hess-18-2829-2014)
498 [2829-2014](https://doi.org/10.5194/hess-18-2829-2014)

499 Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andréassian, V., Anctil, F., & Loumagne, C. (2005). Which
500 potential evapotranspiration input for a lumped rainfall–runoff model?: Part 2-Towards a simple
501 and efficient potential evapotranspiration model for rainfall–runoff modelling. *Journal of*
502 *Hydrology*, 303(1–4), 290–306. <https://doi.org/10.1016/j.jhydrol.2004.08.026>

503 Perrin, C., Michel, C., & Andréassian, V. (2003). Improvement of a parsimonious model for streamflow
504 simulation. *Journal of Hydrology*, 279(1–4), 275–289. [https://doi.org/10.1016/S0022-](https://doi.org/10.1016/S0022-1694(03)00225-7)
505 [1694\(03\)00225-7](https://doi.org/10.1016/S0022-1694(03)00225-7)

506 Popper, K. (1959). *The logic of scientific discovery*. London: Routledge.

507 Rau, P., et al. (2019). Assessing multidecadal runoff (1970–2010) using regional hydrological modelling
508 under data and water scarcity conditions in Peruvian Pacific catchments. *Hydrological Processes*,
509 33(1), 20–35. <https://doi.org/10.1002/hyp.13318>

510 Refsgaard, J. C., & Henriksen, H. J. (2004). Modelling guidelines—terminology and guiding principles.
511 *Advances in Water Resources*, 27, 71–82. <https://doi.org/10.1016/j.advwatres.2003.08.006>

512 Refsgaard, J. C., & Knudsen, J. (1996). Operational validation and intercomparison of different types of
513 hydrological models. *Water Resources Research*, 32(7), 2189–2202.
514 <https://doi.org/10.1029/96WR00896>

515 Refsgaard, J. C., et al. (2013). A framework for testing the ability of models to project climate change
516 and its impacts. *Climatic Change*, 122(1–2), 271–282. [https://doi.org/10.1007/s10584-013-](https://doi.org/10.1007/s10584-013-0990-2)
517 [0990-2](https://doi.org/10.1007/s10584-013-0990-2)

518 Roberts, D. R., et al. (2017). Cross-validation strategies for data with temporal, spatial, hierarchical, or
519 phylogenetic structure. *Ecography*, 40(8), 913–929. <https://doi.org/10.1111/ecog.02881>

520 Royer-Gaspard, P., Andréassian, V., and Thirel, G. (2021) Technical note: PMR – a proxy metric to assess
521 hydrological model robustness in a changing climate. In review for *Hydrol. Earth Syst. Sci.*
522 *Discuss.*, <https://doi.org/10.5194/hess-2021-58>

523 Seibert, J. (2003). Reliability of model predictions outside calibration conditions. *Nordic Hydrology*,
524 34(5), 477–492. <https://doi.org/10.2166/nh.2003.0019>

525 Seifert, D., Sonnenborg, T. O., Refsgaard, J. C., Højberg, A. L., and Trolborg, L. (2012). Assessment of
526 hydrological model predictive ability given multiple conceptual geological models, *Water*
527 *Resour. Res.*, 48, W06503, doi:10.1029/2011WR011149.

528 Seiller, G., Anctil, F., & Perrin, C. (2012). Multimodel evaluation of twenty lumped hydrological models
529 under contrasted climate conditions. *Hydrology and Earth System Sciences*, 16(4), 1171–1189.
530 <https://doi.org/10.5194/hess-16-1171-2012>

531 Tanaka, T. and Tachikawa, Y. (2015). Testing the applicability of a kinematic wave-based distributed
532 hydrological model in two climatically contrasting catchments. *Hydrological Sciences Journal*, 60
533 (7–8). <https://doi.org/10.1080/02626667.2014.967693>

534 Taver, V., et al. (2015). Feed-forward vs recurrent neural network models for non-stationarity
535 modelling using data assimilation and adaptivity. *Hydrological Sciences Journal*, 60 (7–8).
536 <https://doi.org/10.1080/02626667.2014.967696>

537 Teutschbein, C. & Seibert, J. 2013. Is bias correction of regional climate model (RCM) simulations
538 possible for non-stationary conditions? *Hydrology and Earth System Sciences*, 17, 5061–5077.,
539 <https://doi.org/10.5194/hess-17-5061-2013>

540 Thirel, G., Andréassian, V., & Perrin, C. (2015b). On the need to test hydrological models under
541 changing conditions. *Hydrological Sciences Journal*, 60(7–8), 1165–1173.
542 <https://doi.org/10.1080/02626667.2015.1050027>

543 Thirel, G., et al. (2015a). Hydrology under change: an evaluation protocol to investigate how
544 hydrological models deal with changing catchments. *Hydrological Sciences Journal*, 60(7–8),
545 1184–1199. <https://doi.org/10.1080/02626667.2014.967248>

546 Vaze, J., Post, D. A., Chiew, F. H. S., Perraud, J. M., Viney, N. R., & Teng, J. (2010). Climate non-
547 stationarity - Validity of calibrated rainfall-runoff models for use in climate change studies.
548 *Journal of Hydrology*, 394(3–4), 447–457. <https://doi.org/10.1016/j.jhydrol.2010.09.018>

549 Vidal, J.-P., Martin, E., Franchistéguy, L., Baillon, M., & Soubeyroux, J.-M. (2010). A 50-year high-
550 resolution atmospheric reanalysis over France with the Safran system. *International Journal of*
551 *Climatology*, 30(11), P. 1627–1644. DOI: 10.1002/joc.2003. <https://doi.org/10.1002/joc.2003>

552 Vormoor, K., Heistermann, M., Bronstert, A., & Lawrence, D. (2018). Hydrological model parameter
553 (in)stability – “crash testing” the HBV model under contrasting flood seasonality conditions.
554 *Hydrological Sciences Journal*, 63(7), 991–1007.
555 <https://doi.org/10.1080/02626667.2018.1466056>

556 Wilby, R. L. (2019). A global hydrology research agenda fit for the 2030s. *Hydrology Research*, 50(6):
557 1464–1480. <https://doi.org/10.2166/nh.2019.100>

558 Xu, C. (1999). Climate Change and Hydrologic Models: A Review of Existing Gaps and Recent Research
559 Developments. *Water Resources Management*, 13(5), 369–382.
560 <https://doi.org/10.1023/A:1008190900459>

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

563 **11 Appendix – Checking the impact of the partial overlap between**
564 **calibration and validation periods in the RAT**

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

569

570 • **Evaluation method**

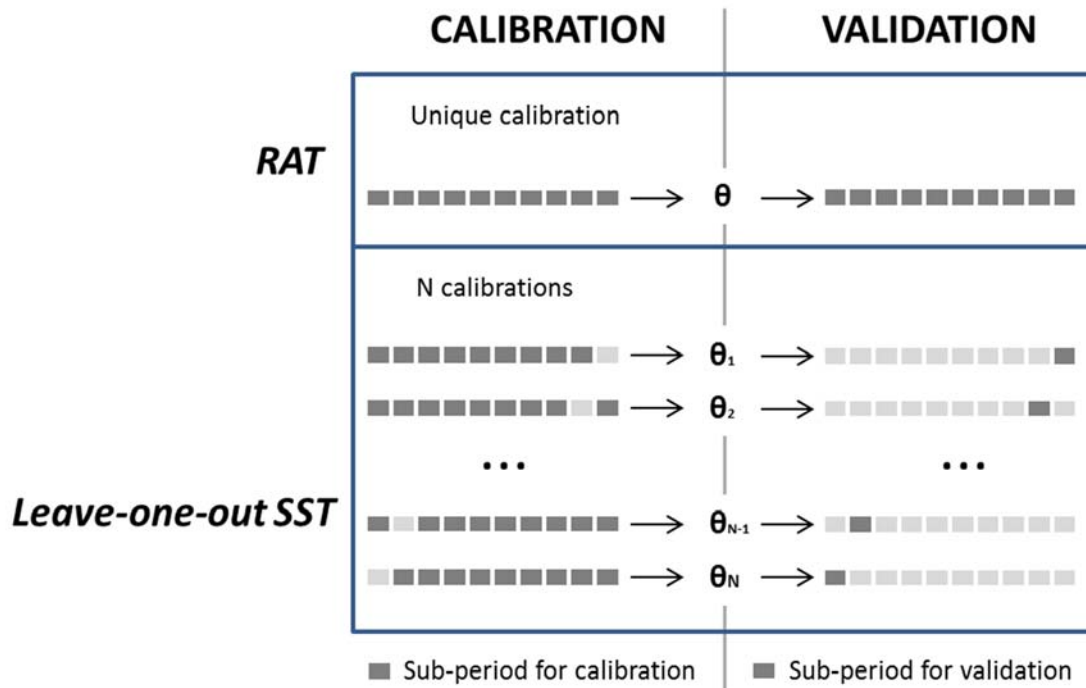
571 In order to check the impact of the partial overlap between calibration and validation periods in the
572 RAT, it is possible (provided one works with a calibrated model) to compare the RAT with a “leave-one
573 out” version of it, which is a classical variant of the Split Sample Test (SST): instead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

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

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

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

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



588

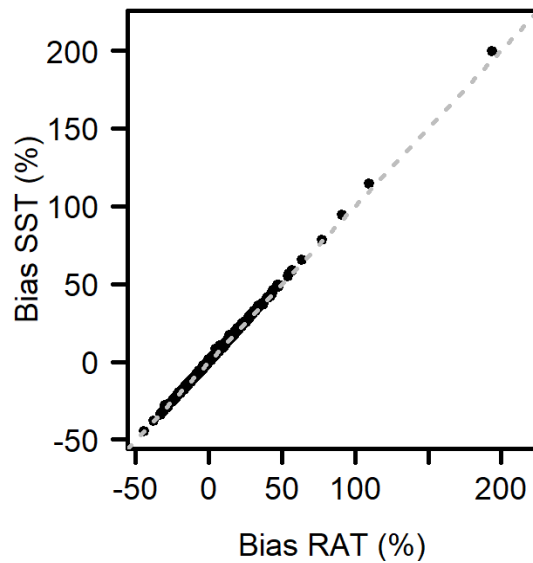
589 **Figure 11. Comparison of the RAT procedure with a leave-one-out split-sample test (SST). Both methods have**
 590 **N validation periods (one per year). The RAT needs only one calibration, whereas the SST requires N**
 591 **calibrations. Dark grey squares represent the years used for calibration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

592

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

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

598

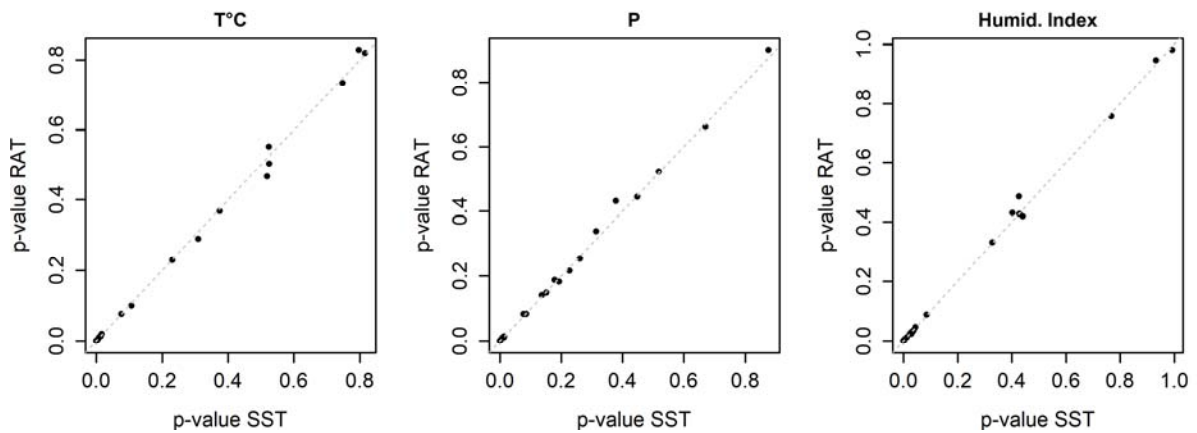


599

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

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

607



608

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

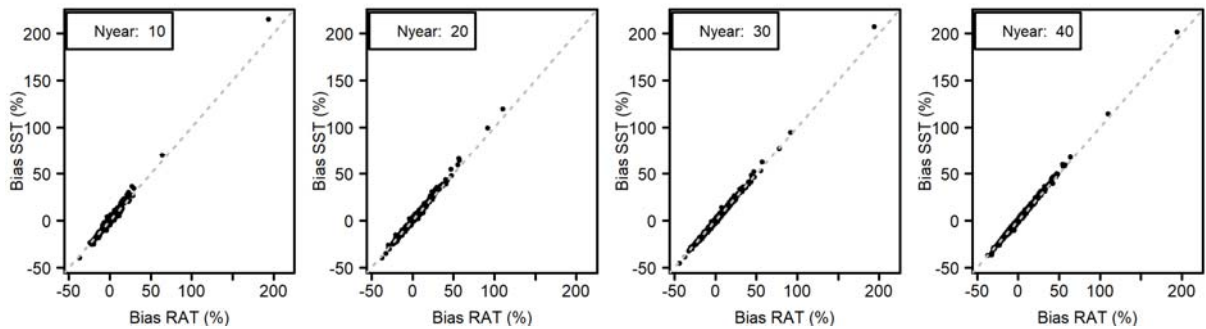
611

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

613 It is also interesting to investigate the limit of our hypothesis (i.e. that the relative weight of one year
 614 within a long time series is very small) by progressively reducing the period length: indeed, the shorter

615 the data series available to calibrate the model, the more important the relative weight of each
 616 individual year. Figure 14 compares the annual bias obtained with the RAT procedure with the annual
 617 bias obtained with the leave-one-out SST, for 10-, 20-, 30-, and 40-year period lengths (selection of the
 618 shorter periods was realized by sampling 10, 20, 30, and 40 years regularly among the complete time
 619 series). The shorter the calibration period, the larger the differences between both approaches (wider
 620 points scatter): there, we reach the limit of the single calibration procedure. We would not advise to
 621 use RAT with time series of less than 20 years.

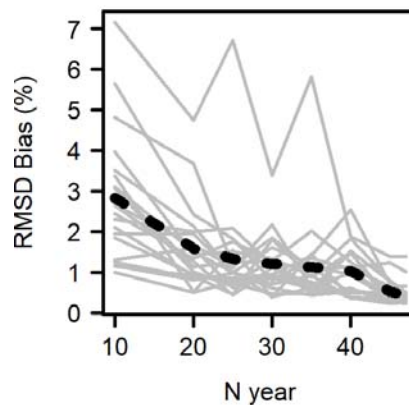
622



623

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

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



631

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