

This is the submitted version of the article:

Đ. Marković, S. Ilić, D. Pavlović, J. Plavšić, and N. Ilich, ‘Multivariate and multi-scale generator based on non-parametric stochastic algorithms’, *Journal of Hydroinformatics*, vol. 21, no. 6, pp. 1102–1117, Nov. 2019, doi: [10.2166/hydro.2019.071](https://doi.org/10.2166/hydro.2019.071).



This work is licensed under the [Attribution-NonCommercial-NoDerivatives 4.0 International \(CC BY-NC-ND 4.0\)](https://creativecommons.org/licenses/by-nc-nd/4.0/)

1 **COMBINED HYDROLOGIC AND WEATHER GENERATOR BASED ON NON-PARAMETRIC**
2 **STOCHASTIC ALGORITHMS**

3

4 Short title: Non-parametric stochastic hydrologic and weather generator

5

6 Đurica Marković^{1*}, Siniša Ilić¹, Dragutin Pavlović², Jasna Plavšić², Nesa Ilich³

7

8 ¹University of Priština at Kosovska Mitrovica, Faculty of Technical Sciences, Kneza Miloša 7, 38220

9 Kosovska Mitrovica, Serbia, djurica.markovic@pr.ac.rs, sinisa.ilic@pr.ac.rs

10 ²University of Belgrade, Faculty of Civil Engineering, P.O. Box 42, 11120 Belgrade, Serbia,

11 dpavlovic@grf.bg.ac.rs, jplavsic@grf.bg.ac.rs

12 ³Optimal Solutions Ltd, 7128-5 Street NW, Calgary, AB T2K 1C8, Canada,

13 nilich@optimal-solutions-ltd.com

14 *corresponding author

15

16 **ABSTRACT**

17

18 A method for generating combined hydrologic and weather time series at multiple locations is presented. The
19 procedure is based on three steps: first, the Monte Carlo method generation of data with statistical properties
20 as close as possible to the observed series; second, the rearrangement of the order of simulated data in the
21 series to achieve target correlations; and third, the permutation of series for correlation adjustment between
22 consecutive years. The method is non-parametric and retains, to a satisfactory degree the properties of the
23 observed time series at the selected simulation time scale and at coarser time scales. The new approach is
24 tested on two case studies, where it is applied to the log-transformed streamflows and precipitation, using
25 weekly and monthly data. Special attention is given to the extrapolation of nonparametric cumulative
26 frequency distributions in their tail zones. The results show a good agreement of stochastic properties
27 between the simulated and the observed data. For example, for one of the case studies the average relative
28 errors of the observed and simulated weekly precipitation and streamflow statistics (up to skewness
29 coefficient) are in the range of 0.1–9.2%, and 0–5.4%, respectively.

30

31 Keywords: stochastic data generation, hydrologic time series, non-parametric methods, serial correlation,
32 cross-correlation

33

34 INTRODUCTION

35

36 Long hydrologic time series are required for effective water resources system planning, design and
37 operation. However, those are often too short, unreliable or non-existent. In these situations, various methods
38 can be used for generating synthetic time series of sufficient length with richer regimes (e.g., containing
39 more extreme values compared to those found in short observed series), while keeping the existing statistical
40 properties of the original series intact. Majority of these methods have been used to generate a single type
41 time series, e.g., streamflow, precipitation, temperature. More recently, there are examples of novel
42 stochastic simulation methods capable of dealing with multivariate stationary or cyclo-stationary processes
43 of any time scale with any marginal distribution and correlation structure. Some of these methods are based
44 on non-parametric approach (Ilich 2014; Srivastav et al. 2015) and some are based on the parametric
45 approach (Tsoukalas et al. 2018a, 2018b; Kossieris et al. 2019).

46

47 Hazen (1914) was probably the first one to use the notion of synthetic time series in hydrology. He generated
48 a 300-year long synthetic hydrologic series combining data from 14 watercourses. Since then, many other
49 approaches emerged for a hydrologic series generation.

50

51 Generating time-dependent hydrologic series, is much more complex than generating independent series
52 since its goal is to preserve not only the statistical distribution function of the original sample but also its
53 auto-correlation function for all significant lags. The well-known model by Thomas and Fiering (1962) with
54 serially correlated flows was the first model of this kind used for monthly flow generation at a single site.
55 This type of model reproduces the essential statistical characteristics of the series, but may lead to unrealistic
56 dependence patterns when combined with non-Gaussian white noise (Tsoukalas et al. 2018c). The problem
57 becomes more difficult for the multi-variate and/or multi-site generation (e.g., streamflows at multiple
58 gauging stations or streamflows and precipitation) where the interstation dependence (i.e., cross-correlation)
59 also has to be preserved in the generated series. The first multi-site stochastic flow generation model was
60 developed by Fiering (1964).

61

62 A number of models for stochastic hydrological time series generation are based on a stochastic processes
63 approach, such as the ARMA models (Box and Jenkins 1970). Despite advantages of the autoregressive (AR)
64 and the moving average (MA) group of models (including ARMA and ARIMA), they suffer from the “short
65 memory” problems, meaning that the serial correlation function quickly diminishes with the time lag
66 (Koutsoyiannis 2000). This approach involves the simultaneous fitting of a large number of parameters related
67 to the joint marginal probability distribution functions in order to comply with the spatial and temporal

68 covariance structure of the shorter historic time series. A detailed review of application of the Box-Jenkins
69 approach in hydrology is presented by Salas et al. (1980). However, recent papers (Tsoukalas et al. 2018a,
70 2018b) introduce methods that can be used for simulation of non-Gaussian univariate and multivariate
71 stationary processes, able to preserve any correlation structure and marginal distribution at any scale. The
72 method was also applied for simulating non-physical process (water demand) at fine scales from 1 hour up to 1
73 minute (Kossieris et al. 2019).

74
75 In the last twenty years, many authors have developed non-parametric methods for simulating hydrologic
76 processes. This became possible by the emergence of new mathematical procedures and methods, and the
77 advances in computational power and software tools. The methods mostly used are the moving block
78 bootstrap (Srinivas & Srinivasan 2005), K-nearest neighbour (K-NN; Sharif & Burn 2006, 2007), or kernel-
79 based methods (Sharma et al. 1997). The main advantage of these methods is that they do not rely on the
80 parameter estimates, while they suffer from the inability to extrapolate the probability distribution beyond
81 the observed data.

82
83 Multi-site streamflow series generation requires a stochastic model capable of reproducing the relevant
84 statistical characteristics of the observed data series. Ideally, the model should be capable of working with
85 selected time discretization (e.g., day, week or month), and also preserve the key statistical characteristics at
86 coarser time scales (e.g., annual). Furthermore, it should be able to extrapolate sensibly the distribution tails
87 for a particular time discretization. Finally, the model also needs to preserve the serial and cross-correlation
88 structure for each time scale, as well as the intra-annual cycle. All these requirements were discussed in
89 detail by Moran (1970), Salas et al. (1980), Koutsoyiannis (2005), Srinivas & Srinivasan (2005).

90
91 Stochastic methods are also used for generating precipitation time series. As precipitation is generally
92 modelled as an intermittent stochastic process, the models need to simulate both precipitation occurrence and
93 intensities/depths in time. Compared to streamflow generation methods, they have to reproduce additional
94 observed data characteristics, such as precipitation occurrence, duration, or distribution of consecutive wet
95 and dry days. Modelling intensities/depths in stochastic precipitation models is identical to modelling
96 streamflow distributions. For the occurrence of dry and wet spells, two types of models are commonly used:
97 Markov chain or renewal process based (Wilks & Wilby 1999). Those based on the Markov chains are often
98 used to specify the state of each spell as wet or dry. These models have been applied to data from various
99 climatic regions and series lengths; however, the structure of the model has to be adjusted to the local
100 conditions for each case study. In addition to the above- mentioned two-part models, resampling models,
101 transition probability matrix models and modifications of ARMA type models (e.g., using normalization
102 transformations or non-Gaussian white noise) are also used for generating rainfall (Srikanthan & McMahon,
103 2001). A good review of the topic is given by Srikanthan & McMahon (2001), Haberlandt et al. (2011), and
104 Serinaldi & Kilsby (2014).

105

106 Harrold et al. (2003a, 2003b) used non-parametric approach for modelling single site daily rainfall
107 occurrences and rainfall amounts, for a 140-year long rainfall record at Sydney, Australia. Their rainfall
108 simulation model is based on the K-NN resampling method where the Markov model was used to generate
109 sequences of dry and wet states (Harrold et al. 2003a). The model preserves short and long-term time series
110 characteristics, i.e., seasonal, annual and multi-annual properties of observed data series. Mehrotra et al.
111 (2006) also applied multi-site K-NN model for precipitation generation at 30 stations in Australia, along with
112 other two parametric generators, while Basinger et al. (2010) used non-parametric procedure based on
113 bootstrapped Markov chains for precipitation occurrence and resampling from observed data for
114 precipitation amounts.

115

116 In addition to the methods for generating single-variate hydro-meteorological series, there is a need to
117 develop approaches for a multi-variate series generation. Such a stochastic model is developed by Srivastav
118 & Simonovic (2014, 2015); this model uses the maximum entropy principle and the bootstrap method to
119 generate multiple variables at multiple sites. It reproduces data statistics, keeping the spatial and temporal
120 structure of data interdependence. The bootstrap method is implemented through the K-NN approach for data
121 generation. The model is tested on daily data (precipitation, maximum and minimum air temperature) from
122 22 gauging stations in the Thames River catchment (Ontario, Canada). However, the method does not
123 preserve the serial correlation between two consecutive years.

124

125 Unlike in some water resources areas related to modelling, such as for example river hydraulics, where there
126 are universally accepted modelling tools such as HEC-RAS, there is no similar tool in stochastic hydrology,
127 i.e. there is no universally accepted time series multi-variate generation model for simultaneous modelling of
128 flows and precipitation that is widely used by hydrologists around the world. In the recent work by Ilich and
129 Despotovic (2008), Ilich (2014) and Marković et al. (2015), a different approach to the generation of stochastic
130 streamflow series is developed that presents an essential departure from the previously established methods.
131 The proposed method consists of three steps: (1) independent data sets for the given time step are generated
132 using the Monte Carlo method, in which the statistical distribution functions of the observed series are fully
133 maintained; (2) data from the individual data sets are then rearranged to induce serial and cross-correlation
134 coefficients of the observed series, and (3) annual streamflows are rearranged to adjust their serial correlation
135 for time intervals that cross-connect two consecutive years. Such an approach has not been proposed by other
136 studies. Moreover, up to our best knowledge, other approaches do not deal explicitly with correlation between
137 data in the transition from one year to another, which is in our methodology done in step 3 by re-ordering
138 whole years in the generated weekly/monthly series. Ilich and Despotovic (2008) have applied this
139 methodology to weekly streamflows. Ilich (2014) has introduced the intermittent precipitation series along with
140 the continuous weekly streamflow series in the simulation procedure. Marković et al. (2015) made further
141 improvements in order to enhance the method's performance by employing the logarithmic transformation to

142 data in order to reduce skewness coefficient and the effect of outliers, and by including additional control to
143 simulate persistence of extremely low summer and autumn flows in dry years.

144

145 This paper builds on the previous work of Ilich and Despotovic (2008), Ilich (2014) and Marković et al. (2015)
146 by expanding the methodology for combined hydrologic and weather generation of time series. The
147 improvement of the methodology lies in introducing a new method for extrapolation of distribution tails, which
148 is different to the use of parametric distributions in Ilich and Despotovic (2008) and Ilich (2014). The main
149 advantages of the proposed methodology are: (1) starting from the shortest time step considered, the
150 methodology ensures that statistics are preserved for all larger steps, (2) the method preserves the serial
151 correlation between two consecutive years, (3) both continuous and intermittent hydrologic time series can be
152 generated, and (4) the procedure is completely automated with a set of default agreement criteria. The
153 application of the methodology in this paper includes streamflow and precipitation data in Canada and in
154 Serbia. While Markovic et al. (2015) generated streamflow data for both Canada and Serbia, in this paper
155 streamflow and precipitation data are jointly generated. By comparing these two sets of results, the efficiency
156 of the generating algorithm is evaluated in terms of multi-variate applications.

157

158 The next section gives an overview of the proposed methodology. It is followed by its application to two
159 datasets of weekly flows, one from Serbia (3 hydrologic stations and 1 meteorological station) and one from
160 Canada (7 hydrologic stations and 4 meteorological stations). The last section provides discussion and
161 conclusions with recommendations for further improvements.

162

163 METHODOLOGY

164

165 Hydrologic time series represent continuous natural processes and are defined in practical applications in a
166 discrete form of average flows or total precipitation for a selected time step, such as day, week, or month. They
167 are modelled as stochastic processes characterized by probability distributions and low-order summary statistics
168 (i.e., mean, variance, and skewness coefficient), and correlation structures.

169

170 The non-parametric stochastic generation method used in this study is formulated so to respect the principle
171 that the generated synthetic series should have distribution functions and the correlation structure very similar
172 to those of the observed series. In order to achieve this, statistics such as the mean, standard deviation, and skew
173 at each time step should be preserved in the generated series, and the serial and cross-correlations should match
174 the observed for any significant lag. Annual statistics of the simulated series, such as the annual mean, standard
175 deviation, and serial and cross-correlations should also match the annual statistics of the observed series.

176

177 The procedure of stochastic streamflow generation relies on the assumptions that observed data represent the
 178 natural hydrologic regime. This means it should be free from any effects of regulation, such as an upstream
 179 reservoir operation or diversion structures, and that the observed process at each time step has a unique
 180 statistical distribution that should be matched in the simulated series. This distribution function can be
 181 represented either by a theoretical parametric distribution that fits the data well or by using an empirical
 182 distribution, such as the non-parametric kernel-based distributions. The reason for using the non-parametric
 183 probability distributions is to avoid specifying any particular parametric distribution in the data generation
 184 process. A possible probability distribution model can be based on combining a non-parametric approach
 185 within the range of the observed data with a parametric distribution at tails, with smoothed inter-range
 186 transitions (Ilich 2014).

187
 188 The observed data, that represent the input for the generation procedure, are organized in a matrix \mathbf{X} , as
 189 shown in Figure SM1 in the supplementary material. The number of rows in the matrix is equal to the
 190 number of years n in the record. This matrix consists of K blocks of columns for each of the K stations
 191 considered. If, for example, weekly data are considered, each column in a block contains streamflow series
 192 for one week. For K stations, the total number of variables, i.e., columns in matrix \mathbf{X} , is $M = 52K$ for weekly
 193 data or $M = 12K$ for monthly data. Thus, the matrix \mathbf{X} is given with:

$$194 \mathbf{X} = [x_{ij}], \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, M \quad (1)$$

195 The columns X_j ($j = 1, 2, \dots, M$) of matrix \mathbf{X} represent the series for each selected time step:

$$196 X_j = [x_{ij}], \quad i = 1, 2, \dots, n \quad (2)$$

197
 198 For the given input matrix \mathbf{X} , the correlation matrix \mathbf{C} of size $M \times M$ contains correlation coefficients ρ_{ij}
 199 between two columns X_i and X_j (see Figure SM2 in the supplementary material):

$$200 \rho_{ij} = \text{Corr}(X_i, X_j), \quad i, j = 1, 2, \dots, M \quad (3)$$

201 Diagonal elements ρ_{ij} , when $i = j$, are equal to 1. Non-diagonal elements of matrix \mathbf{C} represent either serial
 202 correlation coefficients (for a single station) or cross-correlation coefficients (inter-station dependence). For
 203 example, for weekly data $\rho_{1,20}$ is the serial correlation between flows in 1st and 20th week at station 1, while
 204 $\rho_{2,72}$ is the cross-correlation between the 2nd week flow at station 1 and the 20th week flow at station 2.

205
 206 The three steps of the proposed procedure for data generation are described in the sequel providing basic
 207 theoretical background for each step (the full details are presented in Marković et al. (2015)) and using the
 208 pseudo codes with the aim of clarifying the method. The procedure is described for weekly flows, but it is
 209 equally valid for other temporal discretization

210

211 **Step 1 – Generation of independent data sets**

212 The first step is to generate N years of random weekly data having the statistical distributions for each time
213 step as close as possible to the target statistics of the historical series represented in the matrix \mathbf{X} of size $n \times$
214 M , where n is number of years of the observed data and M is the total number of columns (i.e. $M = 52K$ for K
215 stations). This step includes compiling the observed series and their log-transformation (to mitigate the
216 skewness intrinsic in the data), defining the target statistics (observed mean value, standard deviation and
217 coefficient of skewness) in the log-space for each week at every station, and then running the Monte Carlo
218 procedure for generating data from the observed distributions. In order to avoid the logarithm of zero, log-
219 transformation of zero precipitation is increased by a constant of 1 mm. For basins exhibiting zero flows, the
220 same would be applied. Generated data from this step are stored in the resulting matrix \mathbf{G} of the generated
221 independent data sets, which has M columns and N rows (in our study, $N = 1000$), but in general N can be as
222 large as necessary.

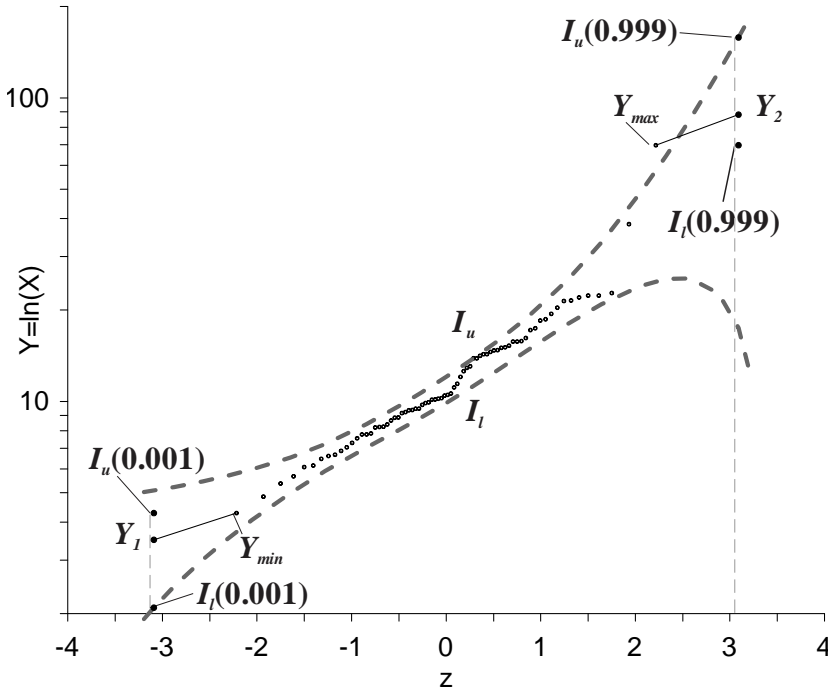
223
224 The probability distributions of the observed data for each week are defined using the non-parametric kernel
225 approach combined with an extrapolation algorithm for the distribution tails. The advantage of the non-
226 parametric approach is that it lends itself to a completely automatic procedure, which is a desirable feature.
227 However, the nonparametric kernel distributions perform poorly outside the range of the observed data. The
228 idea for distribution function extrapolation in the tail sections in this article originates from the work of
229 Scholz (1995). This extrapolation method linearizes distribution tails by utilizing linear dependence of a
230 variate (e.g., streamflow) on the standard variate of a theoretical distribution when plotted on a probability
231 paper. However, depending on the sample data and existence of outliers, extrapolating the lower tail could
232 produce negative values while extrapolating the upper tail could yield generated values much greater than the
233 maximum observed value. Both cases are undesirable.

234
235 To overcome the drawbacks of Scholz's approach, a different heuristic algorithm is applied here for
236 extrapolating distribution tails. The linear extrapolation is applied to the log-transformed variate $Y = \ln X$
237 plotted against the standard normal variate z (Figure 1). The developed algorithm assumes that the upper and
238 lower tail extrapolating lines must lie within the confidence interval of the observed distribution. However,
239 the confidence interval of a non-parametric distribution function cannot be constructed outside of the
240 observed data range. For this reason, the confidence interval limits outside of the observed range are
241 estimated by assuming the General Extreme Value (GEV) distribution. The GEV parameters are estimated
242 by the method of L-moments for each observed weekly series according to formulae given by Rao & Hamed
243 (2000). Each extrapolation line is determined by two points (Figure 1). At the lower tail, the first point is
244 defined by the log-transformed minimum observed value Y_{\min} , and the second point is a randomly selected
245 value Y_1 from the 90% confidence interval of the 0.1% GEV quantile (i.e., for cumulative distribution
246 function or CDF value of 0.1% or standard normal variate $z = -3.09$). Similarly, at the upper tail, the first

247 point is the log-transformed maximum observed value Y_{max} and the second point is a randomly selected value
 248 Y_2 from the 90% confidence interval of the 99.9% GEV quantile (with $z = 3.09$). The “randomness” of the
 249 choice of the points Y_1 and Y_2 from the 90% confidence interval (I_l, I_u) is restricted by the logical conditions:
 250 Y_1 cannot be greater than Y_{min} and Y_2 cannot be smaller than Y_{max} . These constraints can be formalized as:

251
 252
$$I_l(0.001) < Y_1 < \min\{Y_{min}, I_u(0.001)\} \tag{4}$$

253
 254
$$\max\{Y_{max}, I_l(0.999)\} < Y_2 < I_u(0.999) \tag{5}$$



255
 256 Figure 1 – Extrapolation of distribution tails applied in the model; the observed Y_{min} and Y_{max} are connected
 257 to randomly selected points (crosses) from the 90% confidence interval of 0.1% and 99.9% GEV quantiles
 258

259 The random values Y_1 and Y_2 are obtained from the restricted ranges given in Eqs. 4 and 5 by multiplying the
 260 range span by the uniformly distributed random number from the $[0,1]$ interval and by adding the product to
 261 the lower range limit $I_l(0.999)$ at the upper tail, or subtracting it from the upper limit $I_u(0.001)$ at the lower
 262 tail. The outermost points (i.e. selected GEV quantiles) are linearly connected to the log-transformed
 263 minimum and maximum observed values Y_{min} and Y_{max} . These linear dependencies on the log-normal
 264 probability plot are used for random sampling outside the observed data range (as shown in Figure 1).
 265

266 The Algorithm for Step 1 is presented by the pseudo-code for Step 1, which generates M data vectors by
 267 random sampling from the non-parametric distributions of the observed vectors using the pre-set criteria for
 268 agreement of the observed and simulated data statistics (mean, variance and skewness). The generation
 269 process ends when the generated statistics are close enough to the observed ones, as defined by specified
 270 criteria for each statistic. We have chosen to restrain the error in mean logarithmic flows to 0.001

271 (corresponding to an error of 0.1% in the original data space) for generating the first 10,000 data. In the case
 272 that desired mean value is not obtained from the first 10,000 data, the tolerance limit is relaxed to 0.003.
 273 Similarly, the tolerance limit in the skew of the logarithmic flows is 0.03 or 0.05 after 10,000 simulations.
 274

Algorithm 1: Pseudo-code for Step 1 – generating random data vectors

```

1: Load  $n$  years of historical data in matrix  $\mathbf{X}$  (size  $n \times M$ )
2: Pre-process  $\mathbf{X}$  (perform logarithmic transformation and sort each column independently) and
   store in matrix  $\mathbf{XLS}$ 
3: Initialize output matrix  $\mathbf{G}$  (size  $N \times M$ ) for  $N$  years of generated uncorrelated data ( $g_{ij} = 0$ )
4: for  $j = 1$  to  $M$  // for each column vector  $j$  in  $\mathbf{XLS}$ 
5:   calculate the observed ratio of zero values  $p_0$  in column  $j$ 
6:   calculate target statistics for non-zero values
7:   define target CDF by calculating observed non-parametric CDF and extrapolating the tails
   from data in column  $j$ 
8:   for  $i = 1$  to  $N$  // for each year  $i$  to be generated
9:     generate random number  $u$  from the range (0, 1)
10:    if ( $u < p_0$ ) then  $g_{ij} = 0$ 
11:    else  $g_{ij} =$  inverse of target CDF for  $u$ 
12:  end for (next  $i$ ) // generated data vector  $G_j$  created
13:  calculate statistics for non-zero values in the generated vector  $G_j$ 
14:  if statistics for  $G_j$  match target statistics then continue to line 23 (next  $j$ );
15:  else
16:    set new count  $k = 1$ ; set maximum number of iterations
17:    while statistics do not match target statistics or maximum number of iterations is reached
18:      generate new data value  $gg$  (like in lines 9-11)
19:      create trial data vector  $G_j$  by replacing  $g_{kj}$  by  $gg$ 
20:      calculate new statistics for non-zero values in trial data vector  $G_j$ 
21:      if new statistics are better than old statistics then approve a replacement on  $k^{\text{th}}$  position
22:      and set the new statistics is target statistics
23:    else continue generation process with  $k = k + 1$  (go to line 16)
24:  end while
25: end for (next  $j$ )
26: post-process data in  $\mathbf{G}$  from log-transformed to original data space

```

275
 276 If the generated series does not fulfill specified criteria after the first N simulations, the algorithm would
 277 continue to generate $(N+1)^{\text{st}}$ data value and to evaluate statistics of the series in the range $[2, N+1]$ by
 278 comparing it to those of the series in the range $[1, N]$. The process of generating one additional data value
 279 and sequential comparison of updated generated statistics with the observed ones is continued for each data
 280 vector until the specified criterion is met. At the end of the process, N years of log-transformed data are
 281 generated for each data vector, having the marginal distribution that corresponds to that of the observed
 282 vector. The generated series are then transformed back from the log space to the original data space and stored
 283 in matrix \mathbf{G} .

284
 285 Generating precipitation data takes into account that precipitation is an intermittent process and that the CDF
 286 $F(x)$ of a precipitation depth vector consists of two parts: probability p_0 of zero precipitation in one time

287 interval (i.e., dry interval occurrence), and the conditional CDF of precipitation depth during wet interval
 288 $F_1(x)$ weighted by the wet interval probability $(1 - p_0)$:

289
 290
$$F(x) = p_0 + (1 - p_0) \cdot F_1(x) \tag{6}$$

291
 292 Generating precipitation data, therefore, has two stages): (1) assessing the dry interval probability p_0 and the
 293 distribution of precipitation depths in wet intervals $F_1(x)$ from the observed data, and (2) random sampling of
 294 precipitation depths by sampling a random number u from the uniform $[0,1]$ distribution, evaluating F_1 that
 295 satisfies eq. (6) for $F(x) = u$, and finally estimating the corresponding precipitation depth quantile as $x_u =$
 296 $F_1^{-1}[(u - p_0)/(1 - p_0)]$. The remaining procedure is identical to generating streamflow data.

297
 298 **Step 2 – Adjusting the correlation structure of the generated series**

299
 300 The data vectors generated in Step 1 for each week or month represent uncorrelated streamflow or
 301 precipitation series, but they should also have the appropriate correlation structure of the observed series in
 302 order to describe realistically the natural hydrologic or precipitation regime at given locations. The
 303 correlation structure includes serial correlation between weekly or monthly data at each site and cross-
 304 correlation between the sites. In the case of streamflows, it is also important that the persistence of low flows
 305 within an extremely dry year is maintained in the generated time series, leading to the occurrence of
 306 extremely low annual flow.

307
 308 The algorithm for Step 2 is divided into two parts. The first part deals with data rearrangement to match the
 309 correlation of the observed weekly data (Algorithm 2.1), while the second part serves two purposes: it
 310 improves the fit between the distributions of the observed and generated annual minima, and allows user to
 311 control the fraction of extremely dry years in the generated data set (Algorithm 2.2).

312

 Algorithm 2.1: Pseudo-code for the first part of Step 2 – adjusting correlation structure

25: load matrices \mathbf{X} and \mathbf{G} from Step 1
 26: **for** M column vectors in matrix \mathbf{X} , calculate correlation matrix \mathbf{C} (of size $M \times M$) **end for**
 27: **if** \mathbf{C} is not positive definite matrix, **then** calculate the closest positive definite matrix and
 store it in \mathbf{C}
 28: apply the Iman-Conover method to rearrange elements in \mathbf{G} with correlation matrix closest
 to \mathbf{C}

Algorithm 2.2: Pseudo-code for the second part of Step 2 – adjusting extremely dry years

29: set the number nd of extremely dry years in generated data
 30: **for** each station
 31: create vector AO of annual sums of observed weekly data in \mathbf{X} for n years
 32: create vector AG of annual sums of generated weekly data in \mathbf{G} for N years
 33: find the smallest value AO_{\min} in AO
 34: find the smallest nd values in AG and their positions IAG
 35: **for** $i = 1$ to nd // perform a loop with respect to indices IAG in \mathbf{G}

```

36:   while  $AG(IAG_i) > AO_{\min}$  // while generated annual sum in row  $IAG_i$  is greater than the
      minimum observed annual sum
37:     find column  $j$  with the maximum value in row  $IAG_i$  of  $\mathbf{G}$  and store data cell position
       $pos1$ 
38:     find row  $k$  with the minimum value in column  $j$  of  $\mathbf{G}$  and store data cell position  $pos2$ 
39:     swap the values between positions  $pos1$  and  $pos2$ 
40:     recalculate  $AG(IAG_i)$ 
41:   end while
42: end for (next  $i$ )
43: end for (next station)

```

313

314 In the first part of Step 2, the algorithm of Iman and Conover (ICA) (Iman and Conover 1982) is used for data
315 permutations within the generated vectors to achieve target correlation structure. The matrix \mathbf{G} resulting from
316 Step 1 is the input for the algorithm, and its columns are the series to be rearranged. The observed data matrix
317 \mathbf{X} is here used to calculate the observed correlation matrix \mathbf{C} , which is set as a target correlation matrix for
318 ICA. The ICA application was presented in detail in Marković et al. (2015).

319

320 Considering that the purpose of the proposed stochastic method is to provide an input for the optimal design of
321 reservoir storage and/or optimal reservoir operation, it is important that the generated series cover a wide range
322 of input data and include events, such as long droughts, that could be critical for reservoir operation. These
323 events from the lower or the upper tail of flow distributions are not present in the observed series but are
324 expected to emerge within N years, which is usually much greater than the number of years with observations.
325 The critical events are very wet or dry years. The dry years with the total annual runoff below the observed
326 minimum are more critical for water allocation. Although the methodology generally yields the minimum
327 generated streamflow lower than the minimum weekly observed ones, the previously described rearrangement
328 for achieving the target correlation structure may not produce an extremely dry years in which low flows persist
329 over longer durations.

330

331 For this reason, the algorithm of Ilich (2014) is upgraded for additional rearrangement of the simulated data set
332 so that it contains a number of extremely dry years. This is achieved by additional swaps of the smallest weekly
333 flows while keeping previously achieved correlation structure, as explained by Markovic et al. (2015) and
334 shown in Algorithm 2.2 (code lines 35 to 42). One additional rearrangement yields one extremely dry year, but
335 the procedure can be repeated for an arbitrary number n_d of extremely dry years with n_d smallest annual flows.
336 The same procedure of additional rearrangement can be applied for the extreme wet years if they are of interest
337 for the reservoir operation management.

338

339 **Step 3 – Adjusting the correlation of weekly flows from one year to another and of mean annual flows**

340

341 Serial correlation of weekly or monthly hydrologic time series for different lags should not only be preserved
342 within one year but also from one year to another. For example, flows in weeks 1, 2, etc. in a year are

343 dependent on flows in weeks 50, 51 and 52 from the previous year. Such correlations in the observed data
 344 should, therefore, be reflected in the generated data. Also, annual streamflows also exhibit correlations that
 345 should be maintained in the generated series. These two requirements can be achieved by rearranging
 346 complete years (i.e., rows in matrix \mathbf{G}) with already arranged weekly streamflows (Ilich & Despotovic,
 347 2008). By doing so in Step 3 of the methodology, the generated random variates are effectively converted
 348 into time series with the required correlation structure.

349
 350 Algorithm 3 shows the pseudo-code for rearranging generated data to adjust the serial correlation of weekly
 351 data in the transition from one year to another and to adjust the serial correlation of the aggregated annual
 352 data. If s represents the index of the last time interval in a year ($s = 52$ for weekly data), then for any station
 353 from the given data set $\rho_{s,1}$ is the observed serial correlation coefficient between the 52nd week of the current
 354 year and the 1st week of the next year. Similarly, $\rho_{s-1,1}$ describes the correlation between week 51 in the
 355 current year and week 1 in the next year, etc. Performing additional rearrangement to adjust serial correlation
 356 over the time index range $[s - 1, 2]$ accounts for 2 time lags. The rearrangement criteria for station k is to
 357 minimise the statistic D_k representing the sum of squared differences between observed and simulated
 358 transitional correlations up to lag 2 (Ilich & Despotovic 2008):

$$359$$

$$360 D_k = (\rho_{s-1,1}^G - \rho_{s-1,1})^2 + (\rho_{s,1}^G - \rho_{s,1})^2 + (\rho_{s,2}^G - \rho_{s,2})^2 \quad (7)$$

361
 362 where superscript G denotes correlation coefficients in the generated data. The above statistic can be
 363 expanded to include correlations for any number L of weeks at the end and the beginning of year.

364
 365 The correlation structure of the observed annual flows or precipitation also has to be preserved in the
 366 simulated series. Similarly, if RAO_l and RAG_l denote annual serial correlation coefficients for lag l for the
 367 observed and generated data sets, respectively, the criteria D_k can be expanded by the term which measures
 368 the goodness of fit of the annual serial correlations up to lag m :

$$369$$

$$370 D_k = \sum_{p=1}^L \sum_{q=s-L+p}^s (\rho_{q,p}^{kG} - \rho_{q,p}^k)^2 + \sum_{l=1}^m (RAO_l^k - RAG_l^k)^2 \quad (8)$$

371
 372 where q and p are indices of weeks in the transition from one year to another and s is the number of weeks in
 373 a year. The serial correlation of weekly data can generally be adjusted up to an arbitrary lag L , while the
 374 annual serial correlation is adjusted up to the lag $m = N / 4$, where N is the number of data years in the
 375 observed series, as recommended by Box and Jenkins (1970). For all gauging stations, composite criteria
 376 statistic can be introduced as the sum of all D_k values, where K is the number of stations:

$$377$$

$$378 D = \sum_{k=1}^K D_k \quad (9)$$

379
380
381
382
383
384
385

The rearrangement of rows in matrix \mathbf{G} is performed until D is sufficiently small, i.e., smaller than a pre-set value D_0 . To find an appropriate order of years (i.e., rows in matrix \mathbf{G}) that satisfies the transitional weekly and annual correlations, the algorithm in this step combines forward and backward searches for substitute rows, starting from the first and the last row of \mathbf{G} simultaneously. The algorithm stops at the first encounter of satisfied criteria for statistic D .

Algorithm 3: Pseudo-code for Step 3 – adjusting transitional weekly correlation and annual correlation

```
44: load matrices  $\mathbf{X}$  and  $\mathbf{G}$  from Step 2
45: set the value for the number  $L$  of ending/starting weeks in a year to be included in the adjustment
46: set the value for the number  $m$  of lags in annual serial correlation function to be included in the adjustment
47: set the value for the tolerance limit  $D_0$  for the criteria statistic  $D$ 
48: for each station  $k$  of  $K$ 
49:   find transitional weekly correlations:
50:     from  $\mathbf{X}$  extract  $L$  last columns with rows from 1 to  $n - 1$  and  $L$  first columns with rows from 2 to  $n$ 
51:     from  $\mathbf{G}$  extract  $L$  last columns with rows from 1 to  $N - 1$  and  $L$  last columns with rows from 2 to  $N$ 
52:     calculate  $L(L+1)/2$  correlation coefficients  $\rho_{qp}$  between the extracted columns in each:  $\mathbf{X}$  and  $\mathbf{G}$ 
53:     calculate  $D_1(k)$  as the sum of differences between observed and generated  $\rho_{qp}$  for all lags
54:     find annual correlations:
55:       create vector  $AO$  of annual sums of observed weekly data in  $\mathbf{X}$  for  $n$  years
56:       create vector  $AG$  of annual sums of generated weekly data in  $\mathbf{G}$  for  $N$  years
57:       calculate autocorrelation functions  $RAO$  and  $RAG$  of observed/generated annual data up to lag  $m$ 
58:       calculate  $D_2(k)$  as the square sum of differences between  $RAO$  and  $RAG$  for all lags
59:   end for (next station)
60:   calculate statistic  $D = D_1 + D_2$ 
61:   start rearrangement algorithm on matrix  $\mathbf{G}$ : set initial best statistic  $DB = D$ 
62:   for  $i_{asc} = first\ year$  to  $last\ year$  with increment +1
63:     for  $i_{desc} = last\ year$  to  $first\ year$  with increment -1
64:       if  $i_{asc} \diamond i_{desc}$ 
65:         trial swap of data values between rows  $i_{asc}$  and  $i_{desc}$ 
66:         re-calculate correlation coefficients and statistic  $D$ 
67:         if  $D < DB$  then accept trial swap and set  $DB = D$ 
68:         if  $DB < D_0$  then break
69:       end if //  $i_{asc} \diamond i_{desc}$ 
70:     end for (next  $i_{desc}$ )
71:   end for (next  $i_{asc}$ )
```

386
387
388

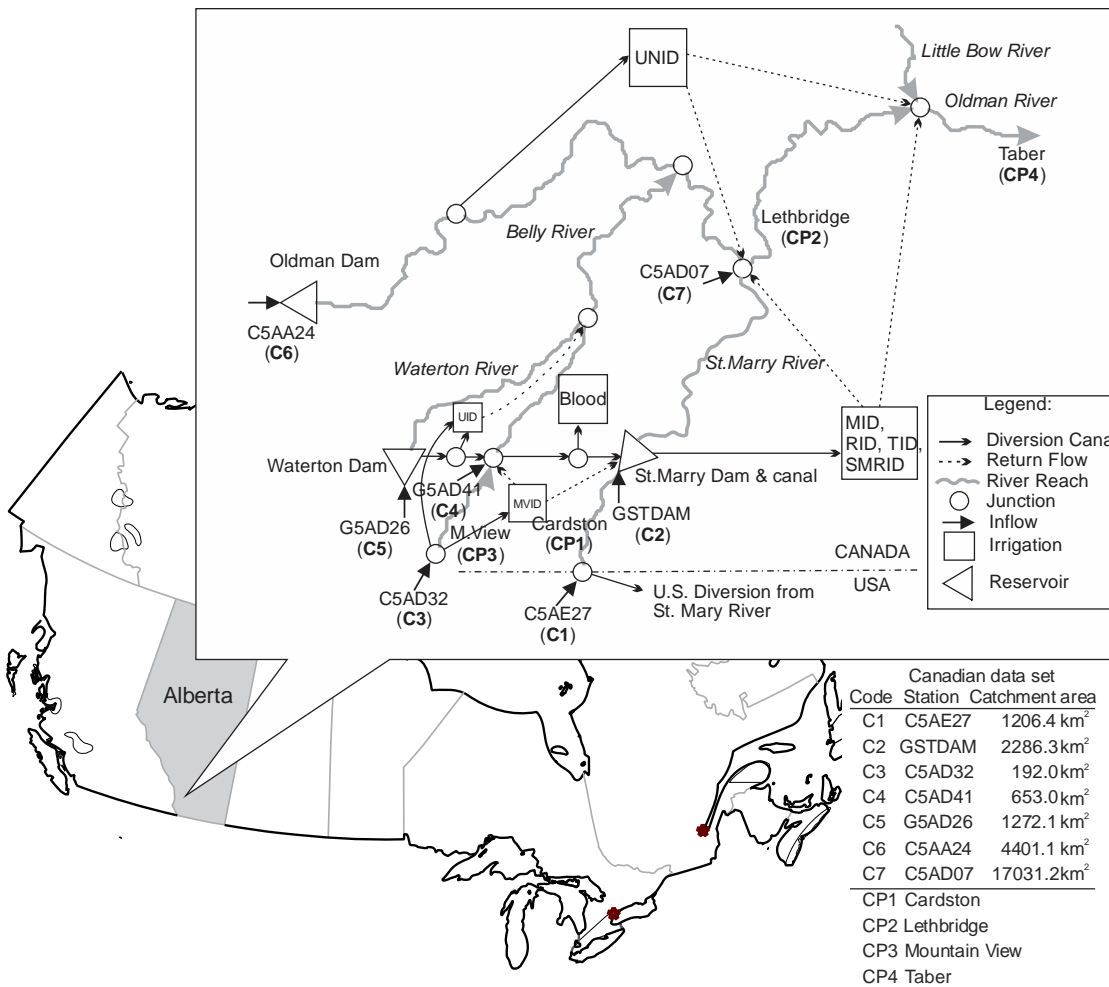
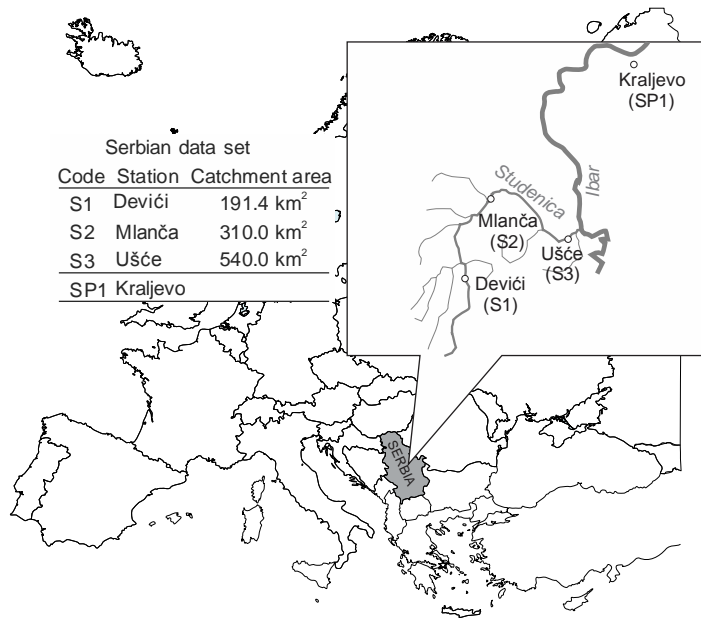
APPLICATION

389

Models and Data Sets

390 The presented method for multi-variate, multi-site and multi-temporal stochastic hydrologic generation is
391 applied to two data sets, one from Serbia and one from Canada, consisting of streamflow and precipitation
392 data series from a different number of stations. For both data sets, two models are applied: (1) model for
393 generation of streamflow series, denoted here MG-Q, and (2) model for generation of streamflow and
394 precipitation series, denoted MG-QP. Both models are applied for two time discretization: weekly and

395 monthly (symbolized by letters *w* and *m* respectively; e.g., MG-Q(*w*) is the model for generating streamflows
 396 on a weekly scale.
 397



398

399 Figure 2 – The map of the study area in Serbia (top) and Canada (bottom) – short codes and full names for
400 the stations used in the model application

401

402 The Serbian data set comprises daily data from three hydrologic stations (Devići, Mlanča, and Ušće) on the
403 Studenica River and meteorological station Kraljevo with the precipitation data (upper part of Figure 2). The
404 streamflow data represent natural flows because there are no water control facilities on the Studenica River.
405 Prior to the application, data were subjected to quality control procedures. Minor gaps were filled using the
406 regression analysis with other stations. The record is 49 years long, from 1964 to 2012.

407

408 The Canadian study region is the Oldman River in Southern Alberta with two of its tributaries, Waterton
409 River and St. Mary River (lower part of Figure 2). Naturalized weekly flows were obtained from Alberta
410 Environment’s natural flow database, with an available record from 1912 to 2001, and for precipitation from
411 1928 to 2001. Ilich (2014) used this data set as an example for his original procedure. Table SM1 in the
412 supplementary material summarizes information for all stations.

413

414 The presented method for stochastic streamflow generation is coded in the MATLAB environment according
415 to the Algorithms 1, 2 and 3 and executed on various computing machines from laptop to desktop PCs. Our
416 experience is that the execution is substantially dependant on the number of variables that are of interest
417 (streamflows, precipitation, temperatures, etc.), the number of gauging stations and the length of the
418 simulation time step. In the case of application of MG-QP model to the Canadian data set (7 streamflow and
419 4 precipitation stations) at weekly time scale, the computational time is as follows: 1h for Step 1, 3 min for
420 Step 2, and 20 h for Step 3. Faster execution would be possible if the code was implemented in computer
421 languages that can be compiled.

422

423 The paper of Markovic et al. (2015) presented the results of simulations involving only the MG-Q model
424 (streamflow data generation only) for Canadian and Serbian data sets. This paper presents the results for
425 MG-QP model that includes both streamflow and precipitation data from Canada and Serbia. These two sets
426 of results enable comparing the efficiency of the algorithm in generating streamflows by taking into account
427 either streamflow dependence structure only, or streamflow-precipitation dependence structure.

428

429 **Results**

430

431 *Results for Step 1 – generation of random series.* The distributions of the generated weekly vectors obtained
432 by the MG-QP model are almost identical to the observed ones. Figure 3(a) shows the empirical distributions
433 of the observed and simulated 10th-week precipitation for the Serbian precipitation station SP1. For
434 comparison reasons, some of the most commonly used parametric distribution functions (Gumbel, Pearson 3,
435 log-Pearson 3, two-parameter gamma) are also applied to the data in Figure 3(a). It can be seen that in this

436 example the employed parametric distributions do not have necessary flexibility to describe the data at
437 distribution tails, while the non-parametric distributions provide that the generated data have almost the same
438 empirical distribution as the observed data. Also, the non-parametric distributions are more appropriate at the
439 lower tail, where some parametric distributions would yield negative values. The same results for stations
440 CP1 and S2 are given in Figures SM3 and SM4, leading to the same conclusions.

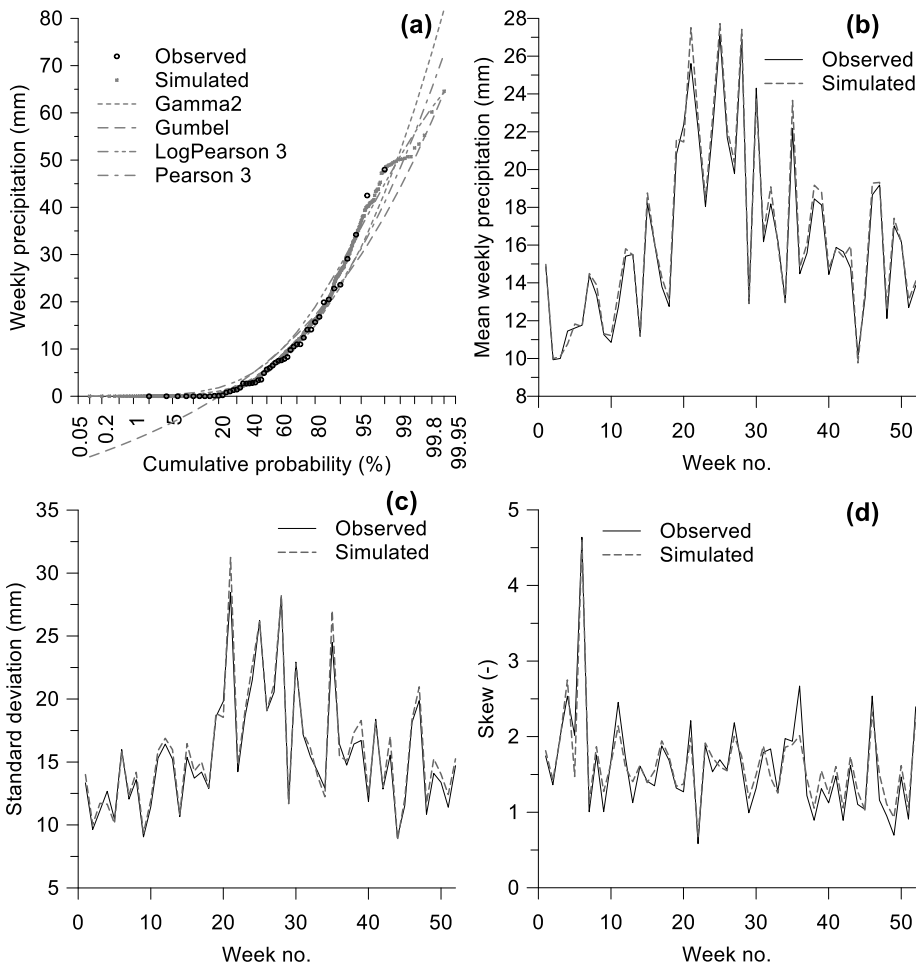
441

442 The good fit of the distributions of the observed and generated vectors also lead to a good fit in the vector
443 statistics. The means, standard deviations and skew coefficients of weekly precipitation are almost identical
444 for the observed and simulated series, as shown in plots (b), (c) and (d) of Figure 3. For example, the relative
445 errors in mean weekly flows/precipitation data are in the range of 0.2–6.4% for station S1 (mean 2.1%), 0.1–
446 6.2% for station S2 (mean 2.2%), 0–5.9% for station S3 (mean 2.3%) and 0.1–7.9% for station SP1 (mean
447 2.7%). Complete results on errors in means are given in Table SM2, showing that the errors for the shorter
448 Serbian data set are comparable with those for the longer Canadian data set.

449

450 The generated data sets have greater maxima than the observed ones, as expected in the longer series (Figure
451 SM6). Similarly, simulated minimum flows are smaller than the observed, as shown by Markovic et al.
452 (2015). With zero being the most frequent minimum value in the observed precipitation series, the same is
453 the case in the simulated series. Also, the percentages of zero values in the observed and the generated
454 precipitation series are very similar (panel (c) in Figure SM6).

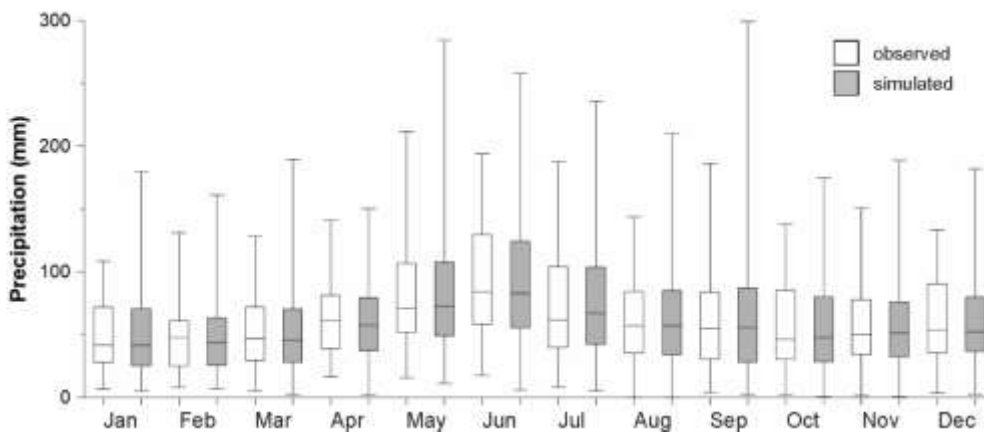
455



456

457 Figure 3 – Model MG-QP(w), precipitation station SP1: (a) empirical distributions of the observed and
 458 simulated precipitation for week 10 compared to four commonly used distributions fitted to the observed
 459 data; (b)-(d) observed and simulated means, standard deviations and skew coefficients of data vectors for
 460 each week.

461



462

463 Figure 4 – Model MG-QP(m), Box-and-whiskers plot of the observed (white) and simulated (grey) monthly
 464 precipitation at meteorological station SP1

465

466 The same conclusions can be made about good reproduction of the distributions of observed monthly
467 vectors. Figure 4 compares these distributions using the box plots. The errors in mean monthly data are in the
468 range of 0.1–4.8% for Serbian stations and 0.0–4.4% for Canadian stations (Table SM3).

469

470 *Results for Step 2 – Serial and cross-correlation.* The data rearrangement resulting from the application of
471 the ICA results in a good fit between the observed and generated correlation structure. Lag 1 and lag 2 serial
472 correlations for stations SP1 and S2 are compared in Figure SM7. Equally good results are obtained for
473 higher lags and for all stations. It is important to notice that the algorithm reproduces not only high
474 correlations but also the small ones, which are below the significance level. The average and maximum
475 differences of the observed and simulated correlation coefficients for weekly data (derived from the
476 correlation matrices for corresponding data) are 0.035 and 0.273 for Serbian stations, respectively, and 0.033
477 and 0.308 for Canadian stations, respectively.

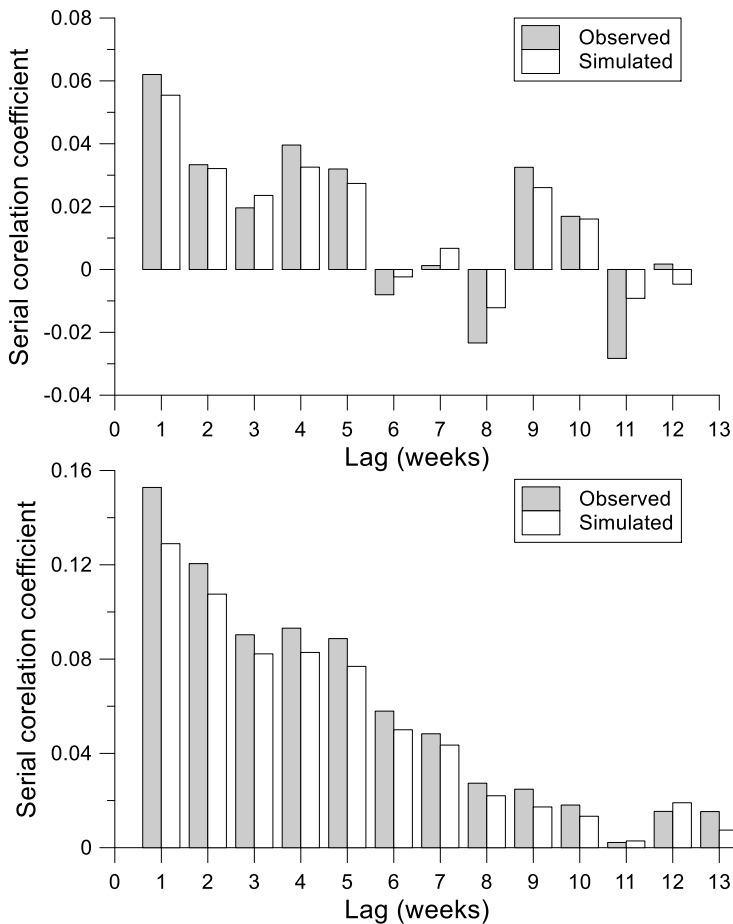
478

479 For the monthly data reproduction of autocorrelation is also good (Figures SM8). The average and maximum
480 differences of the observed and simulated correlation coefficients for monthly data are 0.021 and 0.118 for
481 Serbian data, and 0.022 and 0.169 for Canadian data.

482

483 *Results for Step 3 – transitional weekly correlation and annual correlation.* The data rearrangements in step
484 3 lead to adjustment of the correlation coefficients in the year-to-year transition and therefore at the end of
485 this step the generated data represent the time series with the completely reproduced autocorrelation function
486 (ACF) of the observed time series. The simulation results show that the transitional year-to-year correlations
487 for weekly data are well simulated (Table SM4). The differences between the observed and generated
488 transitional correlations are generally very small (in average 0.036), but the greatest differences (up to 0.383)
489 are attributed to Serbian hydrologic stations. As a result, the ACFs of the observed and generated data are in
490 good agreement. The examples of the ACFs for weekly precipitation are given in Figure 5, showing that the
491 correlation structure is preserved even for small correlations close to zero. Similarly, comparison of the
492 cross-correlation functions for weekly data at selected stations (Figures SM9 and SM10) also shows good
493 agreement.

494



495

496 Figure 5 – Model MG-QP(w), comparison of serial correlation functions of the observed and simulated
 497 weekly precipitation at station SP1 (top) and CP1 (bottom)

498

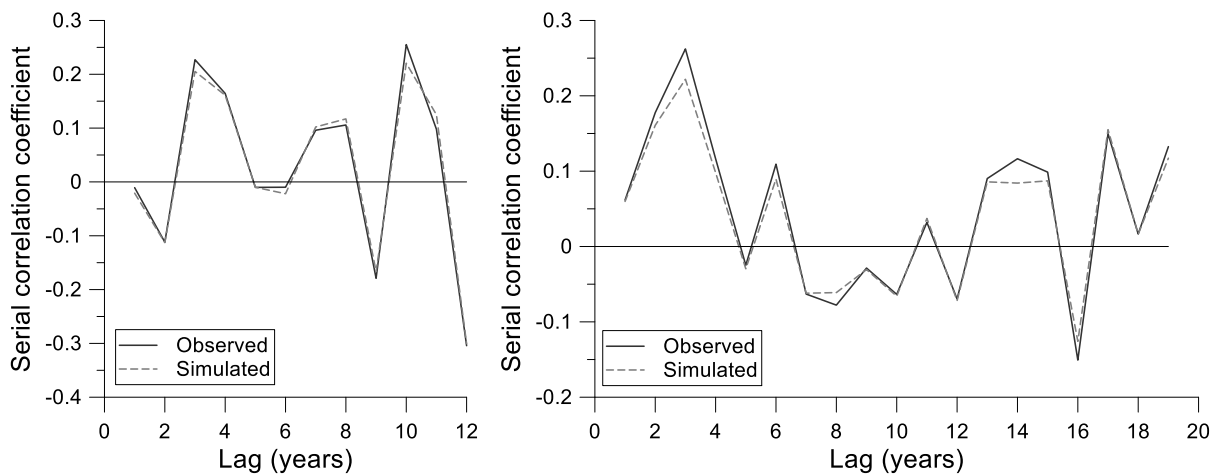
499 When aggregated on a coarser temporal scale, the generated data are comparable to the aggregated observed
 500 data in terms of the annual statistics, distributions and correlation structure. This is shown by aggregating
 501 generated weekly data to 4-weeks scale and to annual scale. An example of the observed and simulated
 502 annual precipitation distribution functions is shown in Figure SM11. This figure also illustrates the effect of
 503 additional treatment at the end of Step 2 over the years with low annual precipitation, which results in a
 504 better agreement of the lower distribution tail.

505

506 The main statistics for the weekly streamflow data aggregated to the 4-weeks scale for one station are
 507 presented in Figure SM12, also showing good agreement. The comparison of the statistics of the annual
 508 streamflows and precipitation aggregated from weekly data is given in Table SM5, showing remarkable
 509 agreement. Differences in the means do not exceed 2.3% and 2.9% for Serbian and Canadian stations
 510 respectively, while the differences in standard deviations are almost negligible for flows and somewhat
 511 greater for precipitation due to its more random nature.

512

513 Serial correlation is also preserved in the aggregated series. Annual ACF of weekly precipitation aggregated
 514 to annual scale for two stations are shown in Figure 6. Statistically insignificant correlations are here well
 515 reproduced in the simulated series for up to 12 lags. The cross-correlation of the annually aggregated weekly
 516 data is also preserved (Table SM6). The average and maximum deviations of the observed and simulated
 517 cross-correlations are 0.008 and 0.048 for Serbian data respectively, while the corresponding values for the
 518 Canadian data are 0.016 and 0.08 (these are slightly greater because more precipitation stations were
 519 included in imposing the correlation structure). Figure SM13 presents ACFs of weekly data aggregated to 4-
 520 weeks scale.
 521



522
 523 Figure 6 – Model MG-QP(w), ACFs of annually aggregated weekly data for stations SP1 (left) and CP1
 524 (right)
 525

526 In the 1000-year long generated series, the annual extreme values should exceed those found in the observed
 527 series. The generated and the observed annual minima or maxima should generally be evaluated in terms of
 528 their distributions. To avoid deciding on the goodness-of-fit of the theoretical distributions to the observed
 529 and generated data, we compare the ranges of theoretical quantiles obtained by fitting some of the commonly
 530 used theoretical distributions to both observed and simulated annual data (we used log-normal, Gumbel,
 531 Pearson 3, log-Pearson 3, two-parameter gamma distributions). The ranges of theoretical quantiles of the
 532 minimum and maximum annual weekly streamflows for station C1 are compared in Figure SM14. The
 533 ranges of theoretical quantiles of generated maxima mostly overlap with those for the observed maxima,
 534 although are somewhat wider. The ranges of theoretical quantiles of generated minima also mostly overlap
 535 with those for the observed minima and can be lower than their observed counterparts for greater
 536 probabilities. This indicates the direction for future improvement of the model.
 537

538 The effects of the rearrangement algorithms in Steps 2 and 3 can also be seen through marginal
 539 improvements in achieving dependence structure of the generated data after Step 1, Step 2 and Step 3. This is

540 illustrated for weekly ACFs, transitional weekly correlations and serial correlation of weekly streamflows
541 aggregated to annual scale in Figures SM15, SM16 and SM17.

542

543 The results for monthly data show equally good agreement of transitional year-to-year correlations (Table
544 SM7) and of complete ACFs (Figure SM18). Comparison of the statistics of the annually aggregated
545 monthly data is shown in Tables SM8 and SM9.

546

547 *Comparison of MG-QP and MG-Q models.* By comparing the results for the streamflows simulated by the
548 MG-QP model presented in this paper with the results of the simulations with the MG-Q model presented in
549 Markovic et al. (2015), no significant differences in the model performance can be seen. For example,
550 empirical distributions of observed and simulated series, observed and simulated weekly flow means,
551 standard deviations and skew coefficients are almost the same for both models (Figures SM4 and SM5).
552 Also, the relative errors in the means of weekly streamflows by the MG-Q model range from 0.0% to 3.65%,
553 which is virtually the same as with the MG-QP model. Additional comparisons of the results of two models
554 are given in the supplementary material (Figures SM19, SM20, and SM21) showing that the model
555 performance is not deteriorated with the introduction of a greater number of variables and more complicated
556 dependence structure of the multi-variate setup.

557

558 CONCLUSIONS

559

560 This paper presents the development and application of the stochastic model for generating simultaneous
561 multi-variate hydrological time series for weekly or monthly temporal scale. The following are the main
562 characteristics of the proposed methodology:

- 563 - It uses non-parametric distributions coupled with the extrapolation algorithm for data generation and
564 non-parametric rearrangement algorithms to achieve the target correlation structure.
- 565 - The heuristic extrapolation algorithm provides robust solution for extrapolating tails and allows fully
566 automated execution of the algorithm.
- 567 - The methodology ensures that the empirical statistic properties of the processes are preserved to a
568 satisfactory degree at the simulation time scale as well as at coarser time scales (e.g. by aggregating
569 from weekly to monthly or annual scale).
- 570 - The method preserves the serial correlation on the transition from one year to another.
- 571 - Both continuous and intermittent hydrologic time series can be generated.
- 572 - The generation process is based on the log-transformed data in order to reduce the effect of outliers
573 and avoid negative generated values.
- 574 - The procedure is completely automated with a set of default agreement criteria.

575

576 The results derived from the two independent data sets (from Serbia and Canada) show that the model can
577 satisfactorily reproduce probability distributions of multi-variate observed series. This is evident from the
578 good match between the main statistics (mean, variance and skewness coefficient) of the generated and the
579 observed data series. For example, the average relative errors of the observed and simulated weekly
580 precipitation and streamflow series are in the range of 0.1–9.2%, and 0–5.4%, respectively (Table SM2), for
581 the Canadian case study. The agreement is achieved by a careful application of nonparametric probability
582 distributions on log-transformed observed data and by using the developed algorithm for the extrapolation of
583 the nonparametric probability distribution.

584

585 The logarithmic transformation of the observed data mitigates the influence of outliers and/or skew in data
586 on the resulting long synthetic data series. The algorithm for the extrapolation of the nonparametric
587 probability distribution uses the linear extrapolation of the cumulative distribution functions using the log-
588 normal probability plot. The extrapolation is performed in the range of the 90% confidence interval of the
589 GEV probability distribution for the 1000-year quantiles. This algorithm enables equally successful
590 simultaneous generation of long streamflow and precipitation series in a hydrologically homogeneous region.

591

592 Two model setups that are considered, one based solely on streamflow data (presented in Marković et al.,
593 2015) and another based on streamflow and precipitation data (presented in this paper), generate series of
594 almost identical stochastic and marginal characteristics to those observed.

595

596 Further research should go in the direction of algorithm refinement regarding computational efficiency for a
597 large number of gauging sites with long records and short time steps (e.g., daily time step). Another
598 improvement can be found in development of a more efficient method for the optimization algorithm in step
599 3.

600

601 ACKNOWLEDGEMENT

602

603 The authors would like to thank Republic Hydro meteorological Service of Serbia for providing data on
604 streamflows and precipitation for this study. Code availability: The authors are open to considering source
605 code sharing request subject to internal procedures among the authors. Data availability: Data are available at
606 request from the first author djurica.markovic@pr.ac.rs under condition that they are not distributed to the
607 third parties.

608

609 REFERENCES

- 610 Basinger, M., Montalto, F., Lall, U. 2010 A rainwater harvesting system reliability model based on
611 nonparametric stochastic rainfall generator. *Journal of Hydrology*, 392, 105–118
- 612 Box G. & Jenkins G. 1970 *Time Series Analysis Forecasting and Control*, Oakland, California, Holden-Day.
- 613 Fiering, M. B. 1964 Multivariate technique for synthetic hydrology, *Journal of the Hydraulics Division*,
614 90(5), 43–60.
- 615 Haberlandt U., Hundecha Y., Pahlow M. & Schumann A. H. 2011 Rainfall generators for application in
616 flood studies. In *Flood Risk Assessment and Management* (pp. 117–147), Springer.
- 617 Harrold T., Sharma A. & Sheather S. 2003a A nonparametric model for stochastic generation of daily
618 rainfall occurrence. *Water Resources Research*, **39**(10), 1300.
- 619 Harrold T., Sharma A. & Sheather S. 2003b A nonparametric model for stochastic generation of daily
620 rainfall amounts. *Water Resources Research*, **39**(12), 1343.
- 621 Hazen A. 1914 Storage to be provided in impounding reservoirs for municipal water supply, *Transactions of*
622 *the American Association of Civil Engineers*, **77**, 1539-1669.
- 623 Ilich N. 2014 An effective three-step algorithm for multi-site generation of stochastic weekly hydrological
624 time series. *Hydrologic Sciences Journal*, **59**(1), 1-14.
- 625 Ilich N. & Despotovic J. 2008 A simple method for effective multi-site generation of stochastic hydrologic
626 time series. *Stochastic Environmental Research and Risk Assessment*, **22**(2), 265-279.
- 627 Iman R. & Conover W. 1982 A distribution-free approach to inducing rank correlation among input
628 variables. *Communications in Statistics - Simulation and Computation*, **11**(3), 311-334.
- 629 Koutsoyiannis D. 2000 A generalized mathematical framework for stochastic simulation and forecast of
630 hydrologic time series. *Water Resources Research*, **36**(6), 1519–1533.
- 631 Koutsoyiannis, D. 2005 Stochastic Simulation of Hydrosystems. In *Water Encyclopedia*. John Wiley &
632 Sons, Inc. <https://doi.org/10.1002/047147844X.sw913>
- 633 Kossieris, P., Tsoukalas, I., Makropoulos, C. & Savic, D. 2019 Simulating Marginal and Dependence
634 Behaviour of Water Demand Processes at Any Fine Time Scale. *Water*, 11(5), 885.
635 <https://doi.org/10.3390/w11050885>

636 Marković Đ., Plavšić J., Ilich N. & Ilic S. 2015 Non-parametric Stochastic Generation of Streamflow Series
637 at Multiple Locations. *Water Resources Management*, **29**(13), 4787-4801.

638 Mehrotra, R., Srikanthan, R., Sharma, A. 2006 A comparison of three stochastic multi-site precipitation
639 occurrence generators. *Journal of Hydrology*, 331, 280-292

640 Moran, P. 1970 Simulation and Evaluation of Complex Water Systems Operations. *Water Resources*
641 *Research*, **6**(6), 1737–1742. <https://doi.org/10.1029/WR006i006p01737> Rao R. & Hamed K. 2000
642 *Flood Frequency Analysis*. CRC Press, Boca Raton.

643 Salas J.D., Delleur J.W., Yevjevich V. & Lane W. 1980 *Applied Modeling of Hydrologic Time Series*. Water
644 Resources Publications, Littleton, CO, USA.

645 Sharif M. & Burn D. 2006 Simulating climate change scenarios using an improved K-nearest neighbor
646 model. *Journal of Hydrology*, **325**, 179-196.

647 Sharif M. & Burn D. 2007 Improved K-Nearest Neighbor Weather Generating Model. *Journal of Hydrologic*
648 *Engineering*, **12**(1), 42-51.

649 Sharma A., Tarboton D. & Lall U. 1997 Streamflow simulation: A nonparametric approach. *Water*
650 *Resources Research*, **33**(2), 291-308.

651 Scholz F. 1995 *Nonparametric Tail Extrapolation*, Boeing Information & Support Services (available at:
652 <http://faculty.washington.edu/fscholz/Reports/ISSTECH-95-014.pdf> (accessed 20 February 2019).

653 Serinaldi, F. & Kilsby C. G. 2014 Simulating daily rainfall fields over large areas for collective risk
654 estimation. *Journal of Hydrology*, 512, 285–302. <https://doi.org/10.1016/j.jhydrol.2014.02.043>

655 Srikanthan R. & McMahon T. 2001 Stochastic generation of annual, monthly and daily climate data: A
656 review. *Hydrology and Earth System Sciences*, **5**(4), 653-670.

657 Srinivas V. & Srinivasan K. 2005 Hybrid moving block bootstrap for stochastic simulation of multi-site
658 multi-season streamflows. *Journal of Hydrology*, **302**, 307-330.

659 Srivastav R. & Simonovic S. 2014 An analytical procedure for multi-site, multi-season streamflow
660 generation using maximum entropy bootstrapping. *Environmental Modelling & Software*, **59**, 59–75.

661 Srivastav R. & Simonovic S. 2015 Multi-site, multivariate weather generator using maximum entropy
662 bootstrap. *Climate Dynamics*, **44**(11–12), 3431–3448

- 663 Thomas H.A. Jr. & Fiering M.B. 1962 Mathematical synthesis of streamflow sequences for the analyses of
664 river basins by simulation. In: *The design of water resources systems*. Harvard University Press,
665 Cambridge, Massachusetts, 459-493.
- 666 Tsoukalas, I., Efstratiadis, A. & Makropoulos, C. 2018a Stochastic Periodic Autoregressive to Anything
667 (SPARTA): Modeling and simulation of cyclostationary processes with arbitrary marginal
668 distributions. *Water Resources Research*, 54(1), 161–185. <https://doi.org/10.1002/2017WR021394>
- 669 Tsoukalas, I., Makropoulos, C. & Koutsoyiannis, D. 2018b Simulation of Stochastic Processes Exhibiting
670 Any-Range Dependence and Arbitrary Marginal Distributions. *Water Resources Research*, 54(11),
671 9484-9513.
- 672 Tsoukalas, I., Papalexiou, S., Efstratiadis, A., & Makropoulos, C. 2018c A Cautionary Note on the
673 Reproduction of Dependencies through Linear Stochastic Models with Non-Gaussian White Noise.
674 *Water*, 10(6), 771. <https://doi.org/10.3390/w10060771>
- 675 Wilks D. & Wilby R. 1999 The weather generation game: a review of stochastic weather models. *Progress*
676 *in Physical Geography*, 23(3), 329-357.