



19

Abstract

20 Stochastically simulated data have been employed for hydrological variables in critical water-
21 related risk management. The simulated data can be utilized to assess the existing flood protection
22 structure and future mitigation frameworks. Disaggregation of the simulated annual data to a lower
23 time scale is often required since water resource management and flood mitigation plans should
24 be done in a fine scale such as a monthly or quarter-monthly. In the current study, the randomized
25 random block length was proposed for the nonparametric disaggregation model since one of the
26 major weakness points for the nonparametric disaggregation model is repetition of similar patterns
27 in the disaggregated data. Furthermore, long-term dependence structure was also mainly focused
28 to preserve since consistent high-flow results devastating damages to inundated area. The proposed
29 model was compared with the existing parametric and nonparametric disaggregation models. The
30 annual net basin supplies (NBS) of the Lake Champlain–Richelieu River (LCRR) Basin was
31 employed to test the performance of the proposed model by reproducing the critical statistics of
32 the 2011 flood in the LCRR Basin. The 2011 flood occurred and was sustained for a few months.
33 The results show that the existing parametric and nonparametric models have limitations and
34 shortcoming and do not provide sufficient temporal dependence. In contrast, the proposed random
35 block-based nonparametric disaggregation (RB-NPD) model with further model enhancement by
36 the genetic algorithm mixture illustrates that the proposed RB-NPD model can be a comparable
37 alternative and that its enhancement is suitable for disaggregating the annual NBS data for the
38 LCRR Basin.

39

40



41 **1. Introduction**

42 Disaggregation models for hydrological variables have been developed in a number of
43 studies to downscale simulated data of a coarse time scale to a fine time scale since water
44 management should be performed in a fine time scale such as monthly or quarter-monthly.
45 Valencia and Schaake (1973) proposed a parametric disaggregation model. The major shortcoming
46 of the model by Valencia and Schaake (1973) is that there is no consideration for the previous year.
47 Mejia and Rousselle (1976) improved the model by including the additional term for the last month
48 of the previous year. However, the models require a significant number of parameters. To avoid
49 parsimoniousness, some models have been proposed in a number of studies (Stedinger and Vogel,
50 1984; Lane and Frevert, 1990; Santos and Salas, 1992) (Santos and Salas, 1992). Furthermore,
51 Koutsoyiannis and Manetas (1996) proposed the accurate adjusting procedure (AAP), which
52 integrates a model for the higher scale (e.g., yearly) and a model for the lower scale (e.g., monthly)
53 by matching the generated sequences at each time scale.

54 Alternative nonparametric disaggregation methods have been proposed in a number of
55 studies (Srikanthan and McMahon, 1982; Porter and Pink, 1991; Tarboton et al., 1998; Prairie et
56 al., 2007; Lee et al., 2010; Lee and Jeong, 2014; Lee and Park, 2017). Prairie et al. (2007) employed
57 the K-nearest neighbor resampling technique, and Lee et al. (2010) improved the model by
58 including the genetic algorithm (GA) mixture.

59 In disaggregation models of hydrologic variables, higher time scale data (here, annual) are
60 disaggregated into lower time scales (here, monthly) according to the relationship between the
61 annual and monthly data. It is relatively easy to preserve the inner-annual relationship (i.e., month-
62 to-month). However, the interannual relationship of the disaggregated monthly data cannot be



63 easily captured since a disaggregation procedure is commonly performed with the current annual
64 data value without the condition of the previous and future relationships. This might result in a
65 discontinuity between the current monthly data and the previous (also following) year for the
66 disaggregated monthly data. It is, in general, not very problematic for disaggregated data since the
67 interannual relationship contains data at a higher time scale (i.e., annual). However, the interannual
68 relationship is critical for the reproduction of a certain event, such as the 2011 flood in the Lake
69 Champlain–Richelieu River (LCRR) Basin, which lasted approximately 3 months. An extreme
70 event that consistently occurs over a long time cannot be generated unless the interannual
71 relationship is appropriately considered.

72 The U.S. and Canadian governments launched an initiative to identify how flood forecasting,
73 preparedness and mitigation can be improved in the LCRR Basin. Synthetic net basin supply (NBS)
74 series of the LCRR Basin are crucial to evaluate the adequacy of flood risk mitigation measures
75 and management strategies under a number of potential hydrological scenarios that might occur in
76 the future. Furthermore, the regulations and management plans for the LCRR Basin have been
77 performed on a monthly or quarter-monthly scale. Therefore, an appropriate temporal
78 disaggregation model to provide more specific information for the basin should be applied, if any,
79 or developed to meet the specific statistical characteristics of the 2011 flood in the LCRR Basin.

80 In the current study, comparable existing parametric and nonparametric models were tested
81 to disaggregate the annual NBS data for the LCRR Basin. The performance of the existing models
82 was carefully examined. Furthermore, a novel approach based on nonparametric techniques was
83 proposed to improve the performance of the existing models, especially for the reproduction of the
84 critical statistics related to the 2011 flood event in the LCRR Basin. Specifically, efforts were



85 made to find and devise a disaggregation model that appropriately captures the interannual
86 relationship in disaggregated data.

87 The present report is organized as follows. A mathematical description of the employed
88 models is described in Section 2. The model development is described in Section 3. The data
89 description and application methodology are presented in Section 4. The results of the compared
90 models and the proposed model are shown in Section 5, followed by a summary and conclusion in
91 Section 6.

92

93 **2. Mathematical Background**

94 **2.1. Parametric Disaggregation**

95 Valencia and Schaake (1973) introduced the basic temporal disaggregation model for annual
96 flows to seasonal flows, followed by its extended version by Mejia and Rousselle (1976), defined
97 as:

$$98 \quad \mathbf{Y}_t = \mathbf{A}X_t + \mathbf{B}\boldsymbol{\varepsilon}_t \quad (1)$$

99 where X_t is the annual time series at year t , \mathbf{Y}_t is the seasonal data for year t , $\tau = 1, \dots, n_m$ as $\mathbf{Y}_t =$
100 $[Y_{t,1}, Y_{t,2}, \dots, Y_{t,\tau}, \dots, Y_{t,n_m}]^T$, and n_m is the number of seasons. \mathbf{A} and \mathbf{B} are $n_m \times 1$ and $n_m \times n_m$
101 parameter matrices, respectively. $\boldsymbol{\varepsilon}_t$ is the $n_m \times 1$ column noise vector uncorrelated with each
102 element distributed as a standard normal. Mejia and Rousselle (1976) included an additional term
103 to preserve the lag-1 correlation between the current year and the past year as:

$$104 \quad \mathbf{Y}_t = \mathbf{A}X_t + \mathbf{B}\boldsymbol{\varepsilon}_t + \mathbf{C}\mathbf{Y}_{t,n_m} \quad (2)$$



105 where \mathbf{C} is the $n_m \times 1$ parameter vector. These disaggregation models suffer from the parsimonious
 106 problem since the models require too many parameters, sometimes more than observations.

107 To avoid this drawback, Lane and Frevert (1990) proposed the condensed version of the
 108 parametric disaggregation model on a one-season-at-a-time basis as:

$$109 \quad Y_{t,\tau} = A_\tau X_t + B_\tau \varepsilon_{t,\tau} + C_\tau Y_{t,\tau-1} \quad (3)$$

110 where, A_τ , B_τ , and C_τ are parameters at each season τ . These parameters can be estimated by
 111 applying the covariance matrices as:

$$112 \quad \hat{A}_\tau = [\mathbf{S}_{YX}(\tau, \tau) - \mathbf{S}_{YY}(\tau, \tau - 1)\mathbf{S}_{YY}^{-1}(\tau - 1, \tau - 1)\mathbf{S}_{YX}(\tau - 1, \tau)] \cdot [\mathbf{S}_{XX}(\tau, \tau) -$$

$$113 \quad \mathbf{S}_{XY}(\tau, \tau - 1)\mathbf{S}_{YY}^{-1}(\tau - 1, \tau - 1)\mathbf{S}_{YX}(\tau - 1, \tau)]^{-1} \quad (4)$$

$$114 \quad \hat{C}_\tau = [\mathbf{S}_{YY}(\tau, \tau - 1) - \hat{A}_\tau \mathbf{S}_{XY}(\tau, \tau - 1)]\mathbf{S}_{YY}^{-1}(\tau - 1, \tau - 1) \quad (5)$$

$$115 \quad \hat{B}_\tau \hat{B}_\tau^T = \mathbf{S}_{YY}(\tau, \tau) - \hat{A}_\tau \mathbf{S}_{XY}(\tau, \tau) - \hat{C}_\tau \mathbf{S}_{YY}^{-1}(\tau - 1, \tau) \quad (6)$$

116 where, $\mathbf{S}_{YX}(a, b)$ indicates a covariance between $Y_{t,a}$ and $X_{t,b}$. Note that \hat{B}_τ is estimated from the
 117 $\hat{B}_\tau \hat{B}_\tau^T$ with either using eigenvalue and eigenvectors or estimating a lower triangular form matrix
 118 (Bras and Rodriguez-Iturbe, 1994). To meet the additive condition, adjustment must be made as:

$$119 \quad Y_{t,\tau}^* = Y_{t,\tau} \times X_t / \sum_{\tau=1}^{n_m} Y_{t,\tau} \quad (7)$$

120 **2.2. Nonparametric Disaggregation (NPD)**

121 Lee et al. (2010) proposed a nonparametric disaggregation model based on k-nearest neighbor
 122 resampling and a genetic algorithm for streamflow applications and further developed it for daily



123 precipitation (Lee and Jeong, 2014; Lee and Park, 2017). The procedure is briefly described as
124 follows:

125 Let the annual, x_t , and seasonal observations $\mathbf{y}_t = [y_{t,1}, \dots, y_{t,n_m}]$ and $t=1, \dots, N$, where N is the
126 record length. In addition, X_t is the target annual variable. The objective is to disaggregate the
127 annual time series X_t to the seasonal time series $\mathbf{Y}_t = [Y_{t,1}, \dots, Y_{t,n_m}]$.

128 Assumed that the number of nearest neighbors, k , is already known, the temporal
129 disaggregation procedure is as follows:

130 (1) The distances between the target annual value X_t and the observed annual variable are
131 estimated as:

$$132 \quad D_i = \begin{bmatrix} X_t - x_i \\ Y_{t-1,n_m} - y_{i-1,n_m} \end{bmatrix}^T \Xi \begin{bmatrix} X_t - x_i \\ Y_{t-1,n_m} - y_{i-1,n_m} \end{bmatrix} \quad i = 2, \dots, N \quad (8)$$

133 where the distances are measured for $i=2, \dots, N$, and Ξ is the variance–covariance
134 matrix of $[x_i, y_{i-1,n_m}]$. Here, the target annual value is considered. In addition, the
135 simulated seasonal value of the last month of the previous year is also taken into
136 account to preserve the dependence of the previous set, as in Lane’s model in Eq.(3).

137 (2) The estimated distances from Step (1) are arranged in ascending order, the first k
138 distances (i.e., the smallest k values) are selected, and the time indices of the smallest
139 k distances are reserved.

140 (3) One of the stored k time indices is randomly chosen with the weighting probability
141 given by:

$$142 \quad w_m = \frac{1/m}{\sum_{j=1}^k 1/j}, \quad m = 1, \dots, k \quad (9)$$



143 (4) The seasonal values of the selected time index (denoted as p) are assigned from Step
144 (3) as $\mathbf{y}_p = [y_{p,1}, \dots, y_{p,n_m}]$.

145 (5) The following steps are executed for GA mixing:

146 (5-1) Reproduction: One additional time index is selected using Steps (1) through (4)
147 and this index is denoted as p^* . The corresponding seasonal values are
148 obtained, $\mathbf{y}_{p^*} = [y_{p^*,1}, \dots, y_{p^*,n_m}]$. The subsequent two GA operators use the two
149 selected vectors, \mathbf{y}_p and \mathbf{y}_{p^*} .

150 (5-2) Crossover: Each element $y_{p,\tau}$ is replaced with $y_{p^*,\tau}$ at the crossover probability P_c ,
151 as:

$$152 \quad Y_{t,\tau} = \begin{cases} y_{p^*,\tau} & \text{if } r < P_c \\ y_{p,\tau} & \text{otherwise} \end{cases} \quad (10)$$

153 where r is a uniform random number between 0 and 1.

154 (5-3) Mutation: Each element (i.e., each season, $\tau=1, \dots, n_m$) is replaced with the one
155 chosen from all observations of this season with the mutation probability P_m , i.e.,

$$156 \quad Y_{t,\tau} = \begin{cases} y_{a,\tau} & \text{if } r < P_m \\ y_{p,\tau} & \text{otherwise} \end{cases} \quad (11)$$

157 where $y_{a,\tau}$ is selected from $[y_{1,\tau}, \dots, y_{N,\tau}]$ with equal probability for $i=1, \dots, N$.

158 (6) The GA mixed values are adjusted as follows to preserve the additive condition
159 as in Eq. (7).

160 (7) Steps (1)-(5) are repeated until the target data are generated.



161 The characteristics of the mutation probability P_m and the crossover probability P_c were
162 studied well by Lee et al. (2010) and Lee (2008). In the current study, two probabilities were used
163 as tuning parameters to manipulate preservation of the historical statistics for the generated
164 monthly time series. The selection of the number of nearest neighbors (k) has been studied (Lall
165 and Sharma, 1996; Lee and Ouarda, 2011). The most common and simplest selection method was
166 applied in the current study by setting $k = \sqrt{N}$. This heuristic approach has commonly been
167 employed in simulation studies with KNNR (Lall and Sharma, 1996; Lee and Ouarda, 2011; Lee
168 et al., 2017; Lee and Ouarda, 2019).

169 2.3. Normal Copula Standardization

170 Seasonal variables, especially in hydroclimatological fields, are commonly skewed. A
171 number of transformation methods have been attempted, such as box-cox, log, and gamma
172 transformation. Among others, copula normal standardization can be a good alternative due to its
173 simplicity and preservation of the marginal statistics as follows:

$$174 \quad Z_t = F_\phi^{-1}[F_Y(Y; \boldsymbol{\theta})] \quad (12)$$

175 where F_Y is the selected distribution for the Y variable, such as gamma, and F_ϕ^{-1} is the inverse
176 standard normal distribution. With this normalization, the result variable (F_Z) definitely has a
177 standard normal distribution. For a marginal distribution, gamma was chosen since this distribution
178 has been commonly used in hydroclimatological variables and fitted well to positively nonnegative
179 skewed variables. Applications have been made in the statistical downscaling in climate change
180 studies (Lee and Singh, 2018). Its back-transformation can be performed by:

$$181 \quad Y_t = F_Y^{-1}[F_\phi(Z; \boldsymbol{\theta})] \quad (13)$$



182 **3. Model Development**

183 **3.1. Random Block-based Nonparametric Disaggregation (RB-NPD)**

184 To avoid discontinuation between the current year and following year, the random block length
185 was used instead of the fixed length (i.e., $l=12$). The proposed method is similar to the original
186 NPD model. However, both the current and following years for annual data must be considered in
187 selecting the candidate. It can be described as follows: (1) to generate the block length (l) from a
188 discrete distribution; (2) to estimate the distances between the observed and generated data, such
189 as the current year and following year values of the observed and target annual data; and (3) to
190 mix the selected block and one additional block. A detailed description is provided in Figure 1:

191 i. A block length, L_B , is generated randomly from a discrete distribution (e.g., geometric or
192 Poisson) for the length of the following seasonal values that follow $Y_{t,\tau}$. Among other
193 distributions, a Poisson distribution is used because the distribution shape is close to a
194 gaussian distribution centered on the mean. More information on the selection of this
195 discrete distribution in block bootstrapping can be found in previous studies (Lee and
196 Ouarda, 2012; Lee and Ouarda, 2019). The Poisson distribution with its parameter (α) is:

$$197 \quad L_B \sim \frac{e^{-\alpha} \alpha^{l-1}}{(\alpha-1)!} \quad l = 1, 2, \dots \quad (14)$$

198 Note that the parameter (α) is the mean of L_B . The parameter for this Poisson distribution
199 (τ) was set as the number of seasons used (e.g., $\alpha = 12$ for monthly) so that the average
200 block length was the same as the number of seasons. In Figure 1, $l=10$ is generated.

201 ii. Distances are estimated to collect close observations to the current status with KNNR as
202 follows:



$$203 \quad D_i = \begin{bmatrix} X_t - x_i \\ Y_{t,\tau-1} - y_{i,\tau-1} \end{bmatrix}^T \Xi_\tau^{-1} \begin{bmatrix} X_t - x_i \\ Y_{t,\tau-1} - y_{i,\tau-1} \end{bmatrix} \quad i = 1, \dots, N \quad (15)$$

$$204 \quad D_i = \begin{bmatrix} X_t - x_i \\ X_{t+1} - x_{i+1} \\ Y_{t,\tau-1} - y_{i,\tau-1} \end{bmatrix}^T \Xi'_\tau^{-1} \begin{bmatrix} X_t - x_i \\ X_{t+1} - x_{i+1} \\ Y_{t,\tau-1} - y_{i,\tau-1} \end{bmatrix} \quad i = 2, \dots, N - 1 \quad (16)$$

205 Eq. (15) or (16) should be used according to whether the generated block length overrides
 206 the next year. For example, in Figure 1, the block length $l=10$ overrides the following year
 207 since the starting season for the current simulation is $\tau=11$, and Eq. (16) must be employed.
 208 Ξ_τ and Ξ'_τ are the variance–covariance matrices for the considered elements as follows:

$$209 \quad \Xi_\tau = \begin{bmatrix} \text{var}(X_t) & \text{cov}(X_t, Y_{t,\tau-1}) \\ \text{cov}(X_t, Y_{t,\tau-1}) & \text{var}(Y_{t,\tau-1}) \end{bmatrix} \quad (17)$$

$$210 \quad \Xi'_\tau = \begin{bmatrix} \text{var}(X_t) & \text{cov}(X_t, X_{t+1}) & \text{cov}(X_t, Y_{t,\tau-1}) \\ \text{cov}(X_t, X_{t+1}) & \text{var}(X_{t+1}) & \text{cov}(X_{t+1}, Y_{t,\tau-1}) \\ \text{cov}(X_t, Y_{t,\tau-1}) & \text{cov}(X_{t+1}, Y_{t,\tau-1}) & \text{var}(Y_{t,\tau-1}) \end{bmatrix} \quad (18)$$

211 The first element ($j=1$) in both equations and the last element in Eq. (16) are omitted
 212 because the data are not available.

213 iii. From the k numbers of the smallest distances among $j = 2, \dots, N$ (or $N-1$), one of the time
 214 indices is chosen with the probability in Eq. (9). Assume that the selected points are p , and
 215 its following sequence is obtained. For example, $l=10$ in Figure 1 and the following
 216 sequence is selected as the generated sequence: $[Y_{t-1,11}, Y_{t-1,12}, Y_{t,1}, \dots, Y_{t,8}] = [y_{p-1,11}, y_{p-1,12},$
 217 $y_{p,1}, \dots, y_{p,8}]$.

218 iv. For the GA crossover, one more sequence is chosen with the steps above (ii)-(iii) and the
 219 additional time index is assumed as p^* in Figure 1 as $[y_{p^*-1,11}, y_{p^*-1,12}, y_{p^*,1}, \dots, y_{p^*,8}]$. Each



220 element of the first chosen sequence is replaced with the crossover probability as in Eq.
 221 (10). In Figure 1, the elements $[y_{p-1,12}, y_{p,5}, y_{p-1,7}]$ are replaced with $[y_{p^*-1,12}, y_{p^*,5}, y_{p^*-1,7}]$

222 v. For the GA mutation, each element is substituted into the gamma random number with the
 223 probability P_m (i.e., $r < P_m$, where r is a uniform random number between 0 and 1), i.e.,

$$224 \quad Y_{t,\tau}^{new} \sim \text{Gamma}(\alpha_\tau, \beta_\tau) \quad (19)$$

225 where α_τ and β_τ are parameters that can be estimated by fitting the gamma distribution to
 226 the seasonal data, i.e., $y_{t,\tau}$ and $t=1, \dots, N$. In Figure 1, $Y_{t,2}$ is simulated from Eq. (19).

227 vi. Simulated seasonal data at each year are adjusted to meet the additive condition as in Eq.(7).

228 vii. Steps (i)-(vi) are repeated until the required data are disaggregated.

229 **3.2. Model Enhancement**

230 Further consideration was tested to preserve the interconnection in the mutated values from
 231 Eq.(19) by replacing the value with the following condition:

$$232 \quad Z_{t,\tau} = \begin{cases} Z_{t,\tau}^{new} & \text{if } Z_{t,\tau-1}Z_{t,\tau}^{new} + Z_{t,\tau}^{new}Z_{t,\tau+1} > Z_{t,\tau-1}Z_{t,\tau} + Z_{t,\tau}Z_{t,\tau+1} \\ Z_{t,\tau} & \text{otherwise} \end{cases} \quad (20)$$

233 where $Z = \Phi^{-1}[F_Y(Y); \alpha_\tau, \beta_\tau]$. Note that the condition of $Z_{t,\tau-1}Z_{t,\tau}^{new} + Z_{t,\tau}^{new}Z_{t,\tau+1} >$
 234 $Z_{t,\tau-1}Z_{t,\tau} + Z_{t,\tau}Z_{t,\tau+1}$ indicates that the newly proposed value ($Z_{t,\tau}^{new}$) has a higher correlation than
 235 the original value ($Z_{t,\tau}$) since Z is the standard normal variable and its multiplication indicates the
 236 correlation. This modification is named after ‘the enhanced correlation algorithm in crossover’ and
 237 denoted as ‘ECAco’.

238 Furthermore, we tested an additional algorithm by simulating the mutation value as



239
$$Z_{t,\tau}^{new} = a_1 Z_{t,\tau-1} + a_2 \tilde{Z}_t + \varepsilon_t \quad (21)$$

240
$$Y_{t,\tau}^{new} = F_Y^{-1}[\Phi(Z_{t,\tau}^{new}); \alpha_\tau, \beta_\tau] \quad (22)$$

241 where $Z_{t,\tau-1}$ is the standard normal variable transformed from $Y_{t,\tau-1}$ with the gamma distribution
242 and \tilde{Z}_t is the standard normal variable of X_t . This enhancement of the parametric simulation
243 algorithm in mutation is denoted as ‘PSAm’.

244 4. Data Description and Application Methodology

245 4.1. Data Description

246 The annual and monthly data of the net basin supply (NBS) series for the Lake Champlain–
247 Richelieu River system (LCRR) were applied in the current study. The water supplies to a lake or
248 a river are referred to as NBS, and they are estimated with both component-based and residual-
249 based methods (Croley and Lee, 1993). The component-based NBS series is used due to its
250 accuracy and popularity in the literature (Fagherazzi et al., 2011; Ouarda and Charron, 2019). The
251 LCRR Basin has an area of 23,900 km² with approximately 84% of the basin in northeastern New
252 York and northwestern Vermont in the US and 16% in Quebec in Canada, as shown in Figure 2.
253 In the spring of 2011 in the LCRR Basin, the worst flooding ever recorded in the past 100 years
254 occurred, which damaged homes, businesses, and farms. Several annual stochastic series for the
255 NBS were simulated in different studies, and further disaggregation to monthly data is required
256 since the regulations and water management were performed on a monthly or quarter-monthly
257 time scale.



258 **4.2. Application Methodology**

259 A number of disaggregation models were considered in the current study, including the parametric
260 VS and MR models and the NPD model (the original nonparametric disaggregation model from
261 Lee et al. (2010)). Furthermore, the proposed RB-NPD model was fully tested with/without the
262 enhanced algorithms in the GA, ECaco and PSAm. In application, the RB-NPD model implies
263 the model without applying the enhanced algorithms and the RB-NPD model with ECaco (called
264 the selective model) is only the enhancement applied in the crossover while the RB-NPD model
265 with ECaco+ PSAm (called hybrid model) presents both the enhancements applied to the
266 crossover and mutation process. Note that the parametric VS and MR models were applied to the
267 copula transformed data with Eq. (12) and the disaggregated data for these models were back-
268 transformed with Eq. (13). Additionally, the Lane model was also tested. However, its
269 performance was not satisfactory and was much worse than those of the VS and MR models.
270 Therefore, its results were not included in this manuscript.

271 To check the performance of the disaggregation models considered in the current study, the
272 observed annual NBS data were disaggregated, and 200 series were produced. Note that the
273 proposed disaggregation model is based on the simulation technique, and an infinite number of
274 series can be produced from the simulation-based disaggregation model. The key statistics of the
275 disaggregated data were estimated and presented by boxplots following the comparison with the
276 observed data. In a boxplot, the boxes represent the interquartile range (IQR), and the whiskers
277 extend up to 1.5 IQR. The horizontal line inside the box shows the median of the data. Data beyond
278 the whiskers (1.5 IQR) are indicated by a plus sign (+).

279 Furthermore, the interesting feature of the disaggregated NBS is the reproduction of the
280 observed high values, especially in a consistent manner. In other words, the flood in the LCRR



281 Basin in 2011 continued for a few months. It is important to generate these continuous high values
282 for a few months in the disaggregated data. Therefore, the key statistics of the accumulated data
283 up to six months were also tested. Note that the n-month accumulation was performed by averaging
284 the monthly data of the previous months. For example, 3 months accumulation for Month 2 (i.e.,
285 $\tau=2$) is the average value of the $Y_{t-1,12}$, $Y_{t,1}$, and $Y_{t,2}$.

286 5. Results

287 5.1. Parametric Disaggregation

288 5.1.1. Valencia–Schaake (VS) model

289 In Figure 3, the observed annual NBS data (top panel) and the observed and disaggregated
290 monthly NBS data are presented (bottom panel), indicating that the disaggregated data reproduce
291 the variability of the observed monthly data with a higher maximum than the observed data. Figure
292 4 presents the basic statistics of the disaggregated monthly data with boxplots and the statistics of
293 the data observed by the dotted line with cross markers. The figure illustrates that the mean and
294 standard deviation are reproduced as well as the extrema (i.e., maximum and minimum), while
295 significant underestimation is found in the skewness.

296 This underestimation of the skewness results stems from the fact that the gamma marginal
297 distribution employed with the copula transformation was not good enough to reproduce these
298 statistics. Three gamma distributions including the location parameter were also tested to
299 reproduce this statistic. However, the location parameter induces another problem: no smaller
300 values than the location parameter were simulated if the location parameter is greater than zero,
301 and negative values were simulated when it is smaller than zero. Other distributions also have



302 similar problems. Furthermore, the original transformation method, such as box-cox, log, and
303 power transformation, produces a larger problem of generating exceptionally large values and
304 negative values after back-transformation (Jeong and Lee, 2015). Additionally, the extreme
305 statistics are underestimated in several months, as shown in the right middle panel in Figure 4.
306 This might be induced from the underestimation of the skewness by the normal copula
307 transformation with the gamma marginal distribution. High skewness cannot be always reproduced
308 in a gamma distribution, and it affects the magnitude and frequency of extreme events in the
309 disaggregated data.

310 The marginal cumulative distribution function and probability distribution function for the
311 disaggregated data with the VS model and observed data are presented in Figure 5 and Figure 6,
312 respectively. As shown in Figure 5, high disaggregated values are lower than those observed with
313 the same CDF in most months (i.e., the blue thick solid line is located on the left side of the dotted
314 red line with a cross marker, especially in months 1, 2, 6, 8, and 9). The skewness as well as the
315 maximum of these months are highly underestimated, as shown in the left and right middle panels
316 in Figure 4.

317 The PDFs in Figure 6 show that the disaggregated marginal distribution does not match the
318 observed one. The observed PDF is more positively skewed than the median of the 200
319 disaggregated series. This indicates that the disaggregated data have lower skewness than the
320 observed data. Furthermore, the observed PDF presents a thicker tail than the disaggregated
321 median, implying that the disaggregated maximum should be lower than the observed maximum.
322 As shown in Figure 4, the disaggregated data from the VS model often underestimate the maximum
323 of the observed data.



324 Another critical point of the performance of the VS model is in the lag-1 autocorrelation
325 (ACF1). The ACF1 of the first month is not preserved as in Figure 4. The ACF1 of the first month
326 indicates the relationship between the last month of the previous year and the first month of the
327 current year (i.e., $\text{corr}(Y_{t-1,n_m}, Y_{t,1})$). This is because there is no term considering the previous year
328 in the VS model in Eq. (1). This is one of the major reasons why the MR mode was devised, as in
329 Eq. (2). Therefore, the extremes of the accumulated data are expected to be underestimated.

330 As shown in Figure 7, the maximum values of the accumulated monthly data up to six
331 months disaggregated with the VS model are underestimated in a number of months. For example,
332 the maximum values of Months 5 and 6 are generally underestimated in most accumulated datasets,
333 as are Months 1 and 2. The reproduction of this statistic for these months is important since the
334 current study started from the 2011 flood that occurred on April 13 and lasted 67 days until June
335 19. The major cause of this underestimation might be the discontinuity between the previous year
336 and the marginal gamma distribution.

337 5.1.2. Meija–Rousselle (MR) model

338 To avoid the discontinuity of the VS model with the previous year, the MR model includes
339 an additional term by taking the last month of the previous year into account, as in Eq. (2). Its
340 performance improvement can be observed in the lag-1 correlation of the basic statistics in the
341 right bottom panel in Figure 8. The first month of ACF1 was improved. Even if it was just one
342 simple improvement, it implies that the disaggregated data are connected to the previous year. The
343 same behavior as the VS model can be observed for the other statistics since the same normal
344 copula transformation with the gamma marginal distribution was applied to this MR model. In
345 Figure 9, the maximum of the accumulated data presents little improvement compared to that of



346 the VS model shown in Figure 7. The underestimation for Months 5 and 6 still often occurred in
347 the MR model.

348 **5.2. Nonparametric Disaggregation**

349 *5.2.1. Original NPD Model*

350 The original NPD model by Lee et al. (2010) was tested. For the NPD model, the results
351 with $P_c=0.1$ and $P_m=0.01$ for the GA mixture in Eqs. (10) and (11) were presented. Rather high
352 values of these probabilities produced more diverse scenarios. The results of the key statistics are
353 shown in Figure 10. All statistics were reproduced well from the NPD model except the slight
354 underestimation of ACF1, which was not significant. The slight underestimation of ACF1 was
355 induced from the GA mixture. Lowering the probabilities could lead to less diverse disaggregated
356 scenarios since these probabilities control the magnitude of mixing and mutating the scenarios
357 from the observed sequences. Combined with the additive adjustment, totally new values and
358 patterns could be produced from the NPD model.

359 The maximum of the accumulated data presents better performance than the parametric
360 model, as shown in Figure 11. The statistics for Month 5 improved in this model, while those for
361 Month 6 were still underestimated. Additionally, the other months, such as Months 2 and 1, also
362 improved. This improvement might be induced from the marginal distribution. As shown in Figure
363 12 and Figure 13 for the CDF and PDF, respectively, the marginal distribution of the disaggregated
364 monthly data reproduced the observed marginal distribution well, especially comparable to the
365 results of the parametric model in Figure 5 and Figure 6. Since the NPD model does not use
366 parameters, especially for a marginal distribution, the characteristics of the observed marginal



367 distribution were preserved well. This feature remains in the other NPD model (i.e., RB-NPD) as
368 well.

369 There were some tangible improvements in the original NPD model compared to the
370 parametric model. However, the NPD model always disaggregated the annual data with a 12-
371 month length basis. The applied combinations (called blocks) for a year were always from one of
372 the observed monthly combinations for a year. In other words, one of the 12 month observed
373 combinations must always be taken for the disaggregation (as $y_{t,1}, y_{t,2}, \dots, y_{t,12}$). Even if the GA
374 mixture is applied to overcome this feature, it still suffers from this limitation. This might result in
375 a weak lagged correlation, especially in the early few months, such as Months 1 and 2. Figure 14
376 presents the lagged correlation starting at each month, i.e., lag-5 correlation at Month 2 is
377 $\text{corr}(Y_{t-1, n_m-4}, Y_{t,2})$. The lagged correlations at Months 1-4 were underestimated in a number of
378 lags, while those of the other months were preserved well.

379 5.2.2. Proposed RB-NPD Model

380 To avoid the fixed length and the effects on the lagged correlation, the random block-based
381 NPD model (i.e., RB-NPD) was devised in the current study. Instead of the fixed 12-month block
382 by Lee et al. (2010), the block length was randomly selected from a Poisson distribution as in
383 Eq.(14). The parameter $\alpha=12$ was used to make the average of the block length the same as the
384 original NPD model. This random block allows the change point of the block to be different from
385 the first month of a year in the original NPD model. By changing the block length into a random
386 block, the distance and its related covariance matrix must be changed at each month. Furthermore,
387 the adjustment to meet the additive condition must be made whenever the random block length
388 overrides the following year.



389 The key statistics of the disaggregated data with the RB-NPD model are presented in Figure
390 15. Note that relatively high probabilities of crossover and mutation for the GA mixture were
391 employed as $P_m=0.1$ and $P_c=0.3$ to produce diverse disaggregated scenarios. Further testing was
392 performed to check whether the employed probabilities were feasible with root mean square error
393 (RMSE), especially the maximum for accumulated data. One of the major objectives was to
394 produce extreme events such as the 2011 flood that had consistently high NBS values for a long
395 period.

396 As shown in Figure 16, increasing the crossover probability (P_c) did not lower the RMSE.
397 Since a high P_c can be beneficial for producing diverse scenarios, $P_c =0.3$ can be acceptable. In
398 addition, a lower P_m value might be better for all accumulated data except Acc-2 shown in the
399 panel in Figure 16(b). With $P_c =0.3$, $P_m =0.1$ is the best choice for Acc-1 (see the blue dotted line
400 with circle in Figure 16 (a)). Additionally, the choice of $P_m =0.1$ is an acceptable choice, showing
401 the second lowest RMSE for all accumulated data with $P_c =0.3$. Note that $P_m =0.01$ can be a good
402 alternative, but this might lower the role of the ECACO and PSAM algorithms in the model
403 enhancement. Therefore, $P_m=0.1$ and $P_c=0.3$ were applied to the following results. Diverse values
404 of P_m and P_c were also tested in the simulation, and no significantly better results were found.

405 All observed statistics were preserved well except ACF1. The underestimation of ACF1 was
406 induced by the high probabilities of mutation and crossover. The enhancement was made by
407 choosing high correlation in the crossover as in Eq. (20), denoted as the enhanced correlation
408 algorithm in the crossover (ECACO, see the model development section). Note that $Z_{t,\tau-1}Z_{t,\tau} +$
409 $Z_{t,\tau}Z_{t,\tau+1}$ indicates the correlation since for variables a and b, $\text{corr}(a,b)=\text{cov}(a,b)/\text{std}(a)\text{std}(b)$ and
410 $\text{std}(a)$ and $\text{std}(b)$ are one and $\text{mean}(a)$ and $\text{mean}(b)$ are zero for standard normal variables (i.e.,



411 $\text{corr}(a, b) = E[a \cdot b]$). Further enhancement was made in the mutation by replacing the monthly
412 sequence with the sequence generated by the parametric simulation algorithm, called PSAm.

413 The basic statistics of the enhanced RB-NPD model with ECaco and PSAm are presented
414 in Figure 17. The figure shows that all statistics were reproduced well, including ACF1. Even
415 slight overestimation of ACF1 is shown. The maximum of the accumulated data in Figure 18 was
416 better preserved in this model compared to the original NPD model in Figure 11, especially Months
417 5 and 6, which were the most critical months when floods occurred in the LCRR Basin.
418 Furthermore, the lagged correlation at each month in Figure 19 was better preserved than that of
419 the NPD model in Figure 14. In particular, the lagged correlations of Months 1, 2, and 3 improved
420 in the enhanced RB-NPD model.

421 Further model testing was performed with wavelet analysis (Foufoula-Georgiou and Kumar,
422 1994; Grinsted et al., 2004) for the RB-NPD models (1) without enhancement (Basic), with ECaco
423 only (Selective), and with ECaco and PSAm (Hybrid). The magnitude-squared coherences (C^2)
424 between the observed data and the disaggregated data from the RB-NPD models are presented in
425 Figure 20. At lower frequencies, strong coherences can be observed, and this is rational since the
426 aggregated data (annual NBS) of the disaggregated data (monthly NBS) are the same as those
427 observed from the additive condition. Figure 21 presents the magnitude-squared coherence for the
428 selected frequencies with high magnitudes. This illustrates that the selective and hybrid models
429 show higher coherence than the basic model and shows that the selective and hybrid algorithms
430 mimic the spectral frequency behavior of the observed data. The results indicate that the proposed
431 algorithm can be a reasonable alternative to disaggregate the annual NBS data to monthly or
432 quarter-monthly scale data.



433 **6. Summary and Conclusions**

434 Based on the 2011 flood in the LCRR Basin, the assessment of the existing flood protection
435 structure and future mitigation frameworks requires simulation scenarios. Furthermore, the
436 scenarios must be on a monthly or quarter-monthly scale. Here, the disaggregation model
437 development was made in the current study focusing on preserving the consistent extreme event
438 of the 2011 flood that was sustained for more than three months.

439 The existing parametric models, VS and MR, and the nonparametric model of the NPD were
440 tested. The VS and MR models employed the normal copula transformation with the gamma
441 marginal distribution instead of the traditional log, box-cox, or power transformation to avoid
442 producing exceptionally large values and negative values. The results were reasonable, but the
443 skewness and maximum statistics were underestimated. Furthermore, the maximum of the
444 accumulated data was not reproduced well, especially for Months 5 and 6, which were critically
445 related to flood events in the LCRR Basin. In contrast, the NPD model reproduced all basic
446 statistics, including skewness and maximum statistics. However, the lagged correlation at each
447 month was underestimated, especially in the first few months (i.e., Months $\tau=1, 2,$ and 3) due to
448 the fixed number of blocks, 12 months, in the disaggregation procedure. Additionally, the
449 maximum of the accumulated data was not appropriately reproduced, especially in Month 6.

450 To overcome the shortages induced by the fixed number of blocks, the random block was
451 suggested in the NPD model as the RB-NPD model. The proposed RB-NPD model varies the
452 starting month of the block, while the starting month of the original NPD model for the block is
453 always Month 1 ($\tau=1$). Further model enhancement was made in the GA mixture to improve the
454 cross-correlation by adding the ECaco and PSAm algorithms. The enhanced RB-NPD model



455 reproduced the lagged correlation as well as the maximum of the accumulated data, especially in
456 Months 5 and 6, which are critical for the purpose of the current disaggregation model.
457 Furthermore, the disaggregated data with the enhanced RB-NPD model present better preservation
458 of the lagged correlation at each month than the NPD model, as well as spectral coherence.

459 Therefore, the results indicate that the proposed RB-NPD model can be a comparable
460 alternative and that its enhancement is suitable for disaggregating the annual NBS data for the
461 LCRR Basin.

462



463 **Abbreviations**

464 Abbreviations

465 LCRR Lake Champlain–Richelieu River

466 NBS Net Basin Supply

467 VS Valencia and Schaake

468 MR Mejia and Rousselle

469 NPD Nonparametric Disaggregation

470 GA Genetic Algorithm

471 KNNR K-Nearest Neighbor Resampling

472 CDF Cumulative Distribution Function

473 PDF Probability Density Function

474 RB-NPD Random Block-based Nonparametric Disaggregation

475 PSAm Parametric Simulation Algorithm in mutation

476 ECACO Enhanced Correlation Algorithm in crossover

477

478

479



480 **7. Code and data availability**

481 The model code and example data are available at Mendeley Data in
482 <<https://data.mendeley.com/datasets/jrrwbc4cx6/1>>. The net basin supply (NBS) observed data
483 are available from the International Joint Commission (IJC, <https://www.ijc.org/en>) upon request

484 **Competing interests**

485 The author declares that they have no conflict of interest.

486 **Author Contribution**

487 L.T. performed writing and collection data as well as programming while O.T. carried out
488 the research plan and supervising.

489

490 **Acknowledgement**

491 The first author acknowledges that this research was partially supported by a grant (2022-
492 MOIS63-001) of Cooperative Research Method and Safety Management Technology in National
493 Disaster funded by Ministry of Interior and Safety(MOIS, Korea).

494



495 **References**

- 496 Bras, R.L., Rodriguez-Iturbe, I., 1994. Random Functions and Hydrology. Dover Books on
497 Advanced Mathematics. Dover, New York, 559 pp.
- 498 Croley, T.E., II, Lee, D.H., 1993. EVALUATION OF GREAT LAKES NET BASIN SUPPLY
499 FORECASTS. JAWRA Journal of the American Water Resources Association, 29(2): 267-
500 282. DOI:10.1111/j.1752-1688.1993.tb03207.x
- 501 Fagherazzi, L., Salas, J.D., Sveinsson, O., 2011. Stochastic Modeling and Simulation of the Great
502 Lakes System, International Upper Great Lakes Quebec.
- 503 Fofoula-Georgiou, E., Kumar, P. (Eds.), 1994. Wavelets in Geophysics. Wavelet analysis and its
504 applications. Academic Press, San Diego, CA, 373 pp.
- 505 Grinsted, A., Moore, J.C., Jevrejeva, S., 2004. Application of the cross wavelet transform and
506 wavelet coherence to geophysical times series. Nonlinear Processes in Geophysics, 11(5-
507 6): 561-566.
- 508 Jeong, C., Lee, T., 2015. Copula-based modeling and stochastic simulation of seasonal intermittent
509 streamflows for arid regions. Journal of Hydro-Environment Research, 9(4): 604-613.
510 DOI:10.1016/j.jher.2014.06.001
- 511 Koutsoyiannis, D., Manetas, A., 1996. Simple disaggregation by accurate adjusting procedures.
512 Water Resources Research, 32(7): 2105-2117.
- 513 Lall, U., Sharma, A., 1996. A nearest neighbor bootstrap for resampling hydrologic time series.
514 Water Resources Research, 32(3): 679-693.
- 515 Lane, W.L., Frevert, D.K., 1990. Applied Stochastic Techniques, Personal Computer, Version 5.2,
516 User's Manual, U.S. Bureau of Reclamation, Denver, Colorado.
- 517 Lee, T., Jeong, C., 2014. Nonparametric statistical temporal downscaling of daily precipitation to
518 hourly precipitation and implications for climate change scenarios. Journal of Hydrology,
519 510: 182-196. DOI:10.1016/j.jhydrol.2013.12.027
- 520 Lee, T., Ouarda, T.B.M.J., 2011. Identification of model order and number of neighbors for k-
521 nearest neighbor resampling. Journal of Hydrology, 404(3-4): 136-145.
- 522 Lee, T., Ouarda, T.B.M.J., 2012. Stochastic simulation of nonstationary oscillation hydroclimatic
523 processes using empirical mode decomposition. Water Resources Research, 48(2).
524 DOI:doi:10.1029/2011WR010660
- 525 Lee, T., Ouarda, T.B.M.J., 2019. Multivariate Nonstationary Oscillation Simulation of Climate
526 Indices With Empirical Mode Decomposition. Water Resources Research, 55(6): 5033-
527 5052. DOI:10.1029/2018WR023892



- 528 Lee, T., Ouarda, T.B.M.J., Yoon, S., 2017. KNN-based local linear regression for the analysis and
529 simulation of low flow extremes under climatic influence. *Climate Dynamics*, 49(9-10):
530 3493-3511. DOI:10.1007/s00382-017-3525-0
- 531 Lee, T., Park, T., 2017. Nonparametric temporal downscaling with event-based population
532 generating algorithm for RCM daily precipitation to hourly: Model development and
533 performance evaluation. *Journal of Hydrology*, 547: 498-516.
534 DOI:10.1016/j.jhydrol.2017.01.049
- 535 Lee, T., Salas, J.D., Prairie, J., 2010. An enhanced nonparametric streamflow disaggregation
536 model with genetic algorithm. *Water Resources Research*, 46(8).
537 DOI:10.1029/2009WR007761
- 538 Lee, T., Singh, V.P., 2018. *Statistical Downscaling for Hydrological and Environmental*
539 *Applications*, 1. CRC press, Boca Raton, FL, 181 pp.
- 540 Lee, T.S., 2008. *Stochastic simulation of hydrologic data based on nonparametric approaches*, Ph.
541 D. Dissertation, Colorado State University, Fort Collins, CO., USA, 346 pp.
- 542 Mejia, J.M., Rousselle, J., 1976. *Disaggregation Models in Hydrology Revisited*. *Water Resources*
543 *Research*, 12(2): 185-186.
- 544 Ouarda, T.B.M.J., Charron, C., 2019. *Frequency analysis of Richelieu River flood flows, Lake*
545 *Champlain flood level and NBS to the Richelieu River basin*, Final project report, INRS-
546 ETE, Quebec, Canada.
- 547 Porter, J.W., Pink, B.J., 1991. A method of synthetic fragments for disaggregation in stochastic
548 data generation, *Int. Hydrol. and Water Resour. Symp.* The Institution of Engineers,
549 Australia, Canberra, pp. 187-191.
- 550 Prairie, J., Rajagopalan, B., Lall, U., Fulp, T., 2007. A stochastic nonparametric technique for
551 space-time disaggregation of streamflows. *Water Resources Research*, 43(3): W03432.
552 DOI:W03432
- 553 10.1029/2005wr004721
- 554 Santos, E.G., Salas, J.D., 1992. Stepwise Disaggregation Scheme for Synthetic Hydrology. *Journal*
555 *of Hydraulic Engineering-Asce*, 118(5): 765-784.
- 556 Srikanthan, R., McMahon, T.A., 1982. Simulation of Annual and Monthly Rainfalls - a
557 Preliminary-Study at 5 Australian Stations. *Journal of Applied Meteorology*, 21(10): 1472-
558 1479.
- 559 Stedinger, J.R., Vogel, R.M., 1984. Disaggregation Procedures for Generating Serially Correlated
560 Flow Vectors. *Water Resources Research*, 20(1): 47-56.
- 561 Tarboton, D.G., Sharma, A., Lall, U., 1998. Disaggregation procedures for stochastic hydrology
562 based on nonparametric density estimation. *Water Resources Research*, 34(1): 107-119.



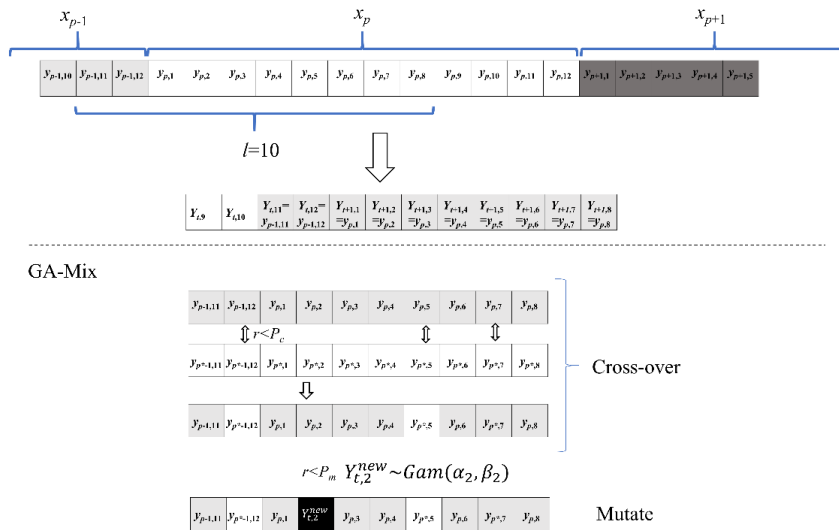
563 Valencia, D., Schaake, J.C., 1973. Disaggregation Processes in Stochastic Hydrology. Water
564 Resources Research, 9(3): 580-585.

565



566 **Figures**

567

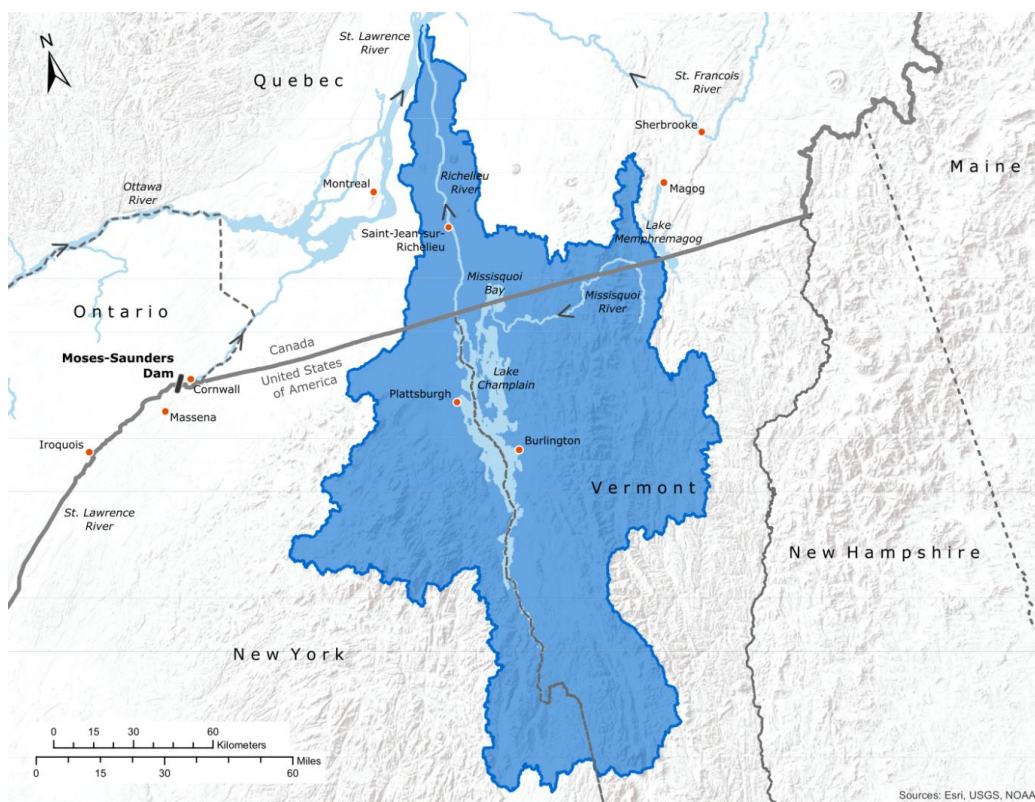


568

569 Figure 1. Diagram of the proposed random block-based nonparametric disaggregation (RB-NPD)
 570 from the current study.

571

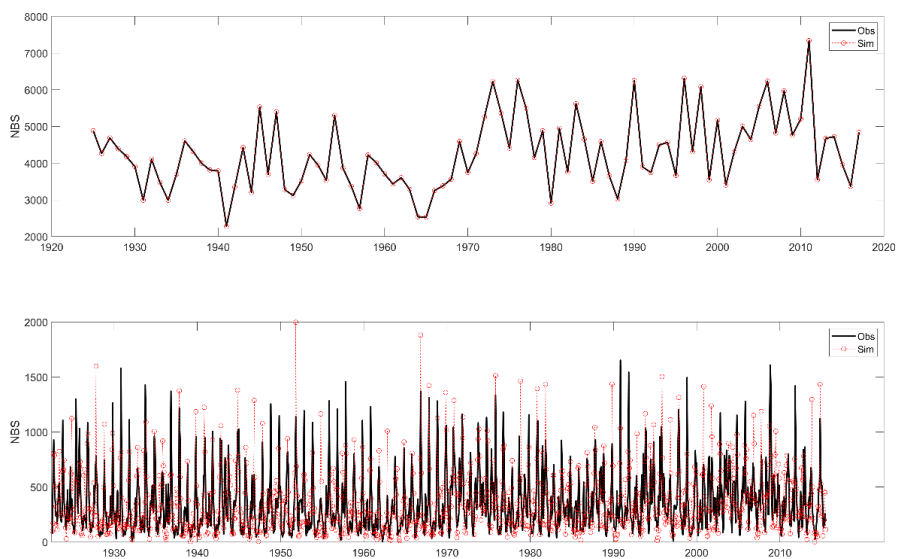
572



573

574 Figure 2. Map of the LCRR Basin. Note that dark blue represents the whole area of the LCRR
575 Basin, light blue inside dark blue represents Lake Champlain, and the northward line from Lake
576 Champlain represents the Richelieu River. Note that the map was provided by the International
577 Joint Commission.

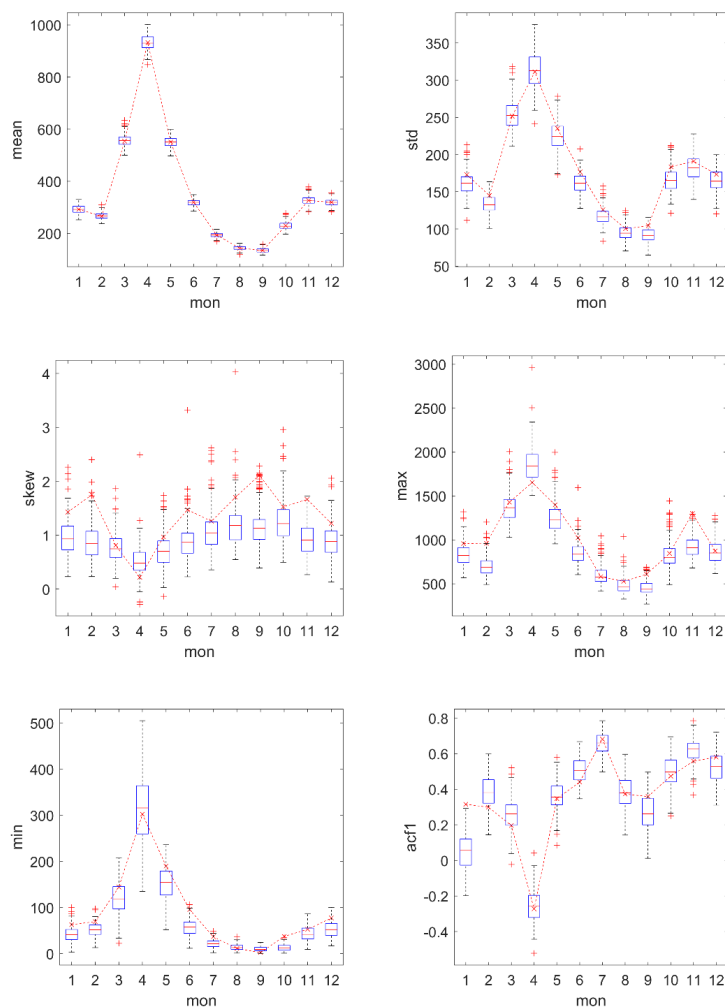
578



579

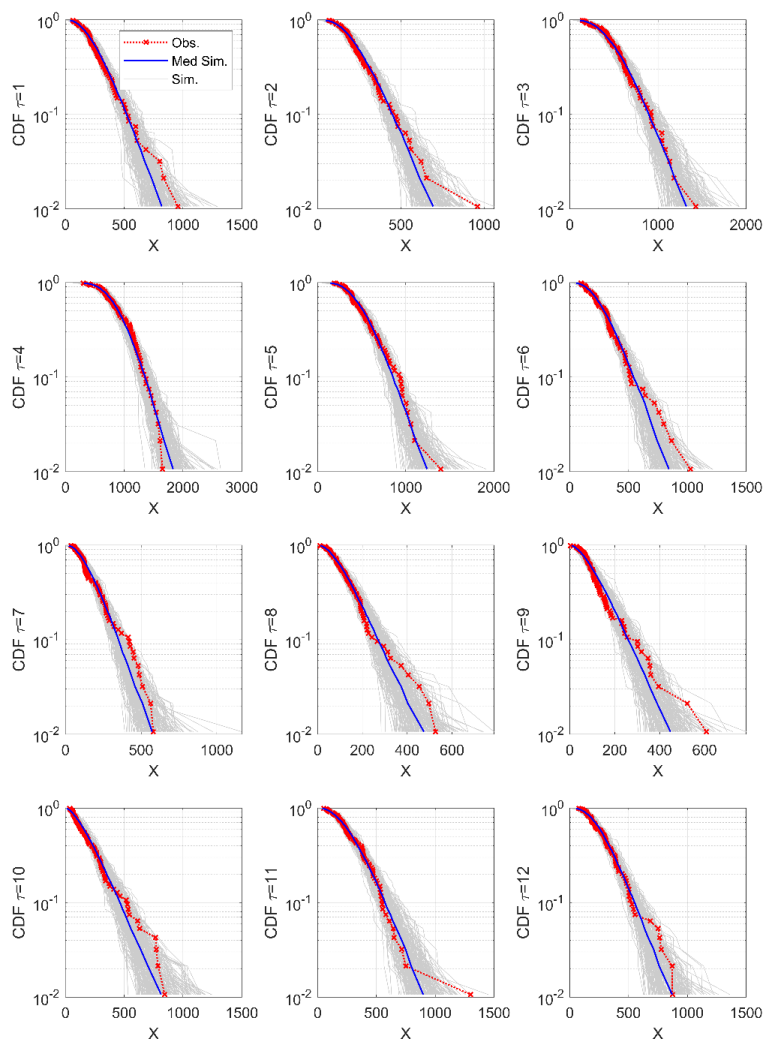
580 Figure 3. Time series of the annual (top panel) and monthly (bottom panel) data for the observed
581 (thick solid line) and disaggregated (dotted line with circle) simulation data with the VS model.

582



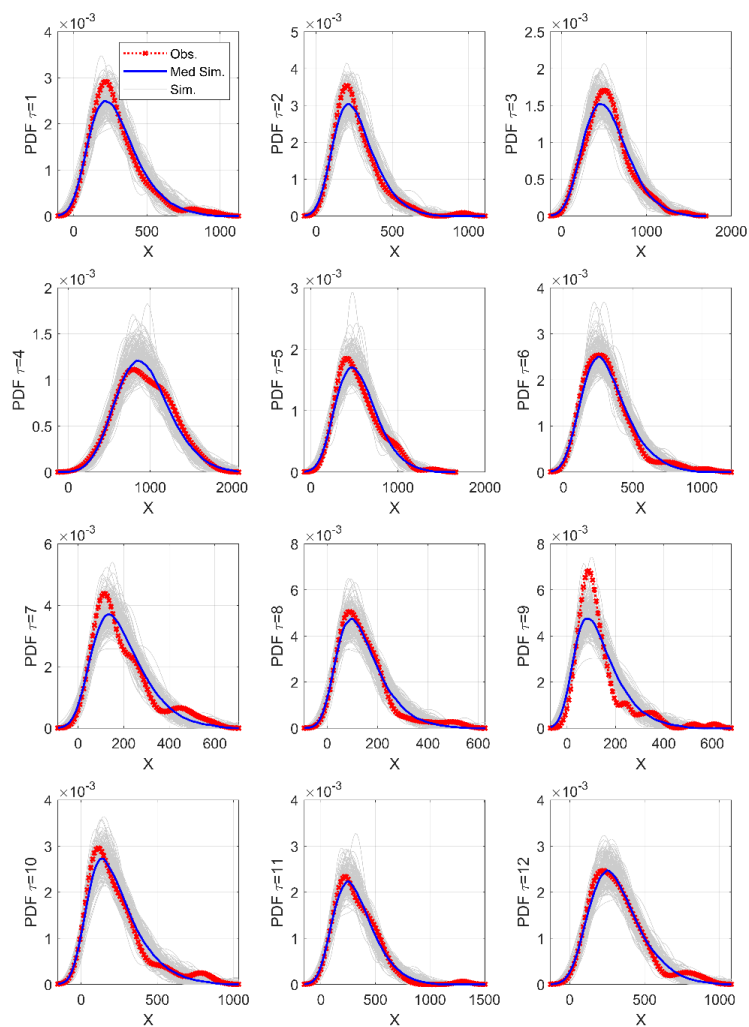
583

584 Figure 4. Boxplots of the basic statistics of the disaggregated monthly data from the annual data
585 with the VS model for the NBS of the LCRB basin. Note that the statistics of the observed data
586 are also represented with the dotted line and cross marker (.x.). The boxes represent the
587 interquartile range (IQR), and the whiskers extend to 1.5 IQR. The horizontal lines inside the
588 boxes depict the data median. Data beyond the whiskers (1.5 IQR) are shown by a plus sign (+).



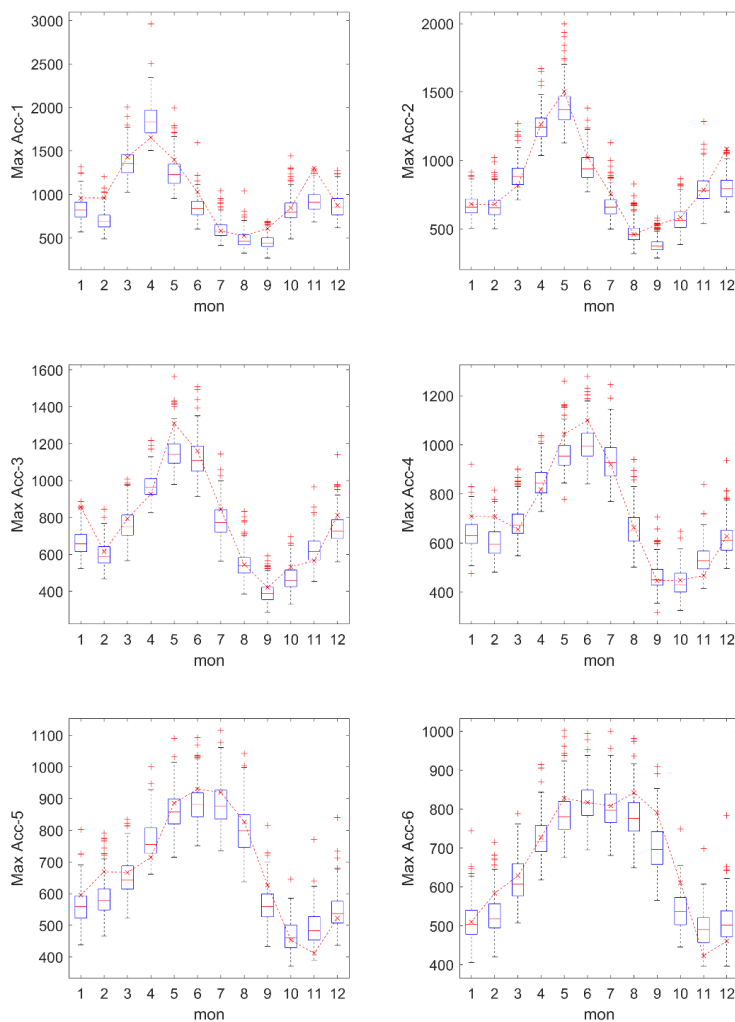
589

590 Figure 5. Cumulative distribution function (CDF) of the disaggregated data in each month with
591 the VS model from the annual NBS of the LCRB Basin. Note that the observed data are
592 represented with the dotted line and cross marker (.x.); (2) the 200 disaggregated simulation
593 series are shown with thin gray lines, while their median is represented with the thin blue line;
594 and (3) τ indicates the month from 1 to 12.



595

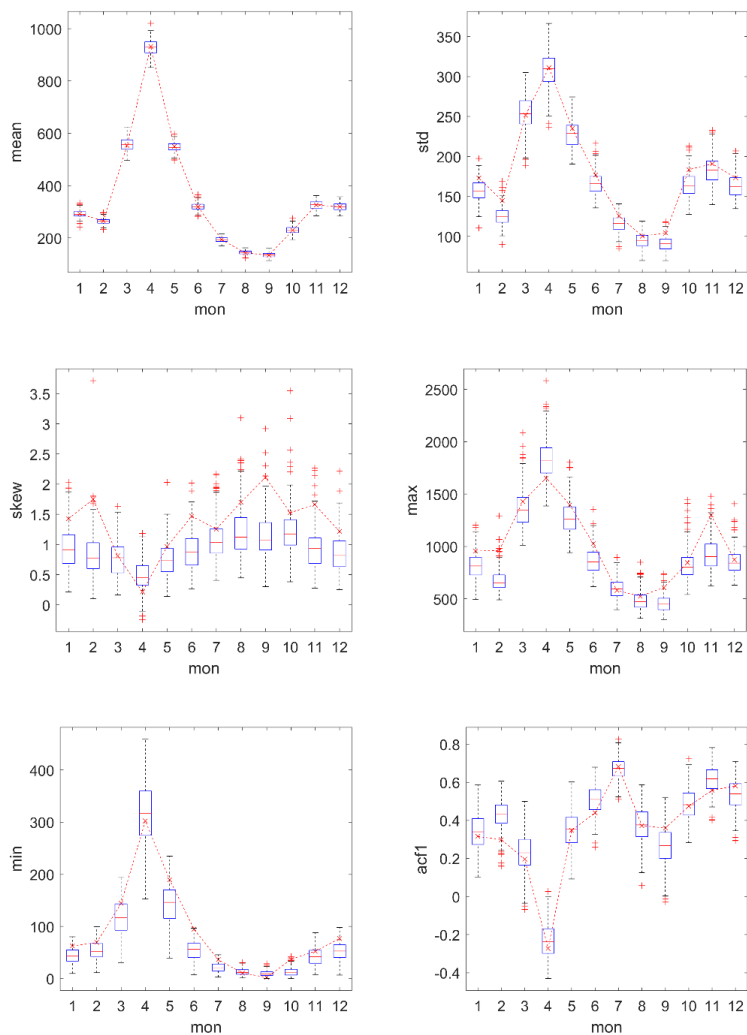
596 Figure 6. Probability density function (PDF) of the disaggregated data at each month with the VS
597 model with $P_m=0.01$ and $P_c=0.1$ from the annual NBS of the LCRR Basin. Note that the
598 observed data are represented with the dotted line and cross marker (.x.); and (2) the 200
599 disaggregated simulation series are shown with thin gray lines, while their median is represented
600 with the thick blue line.



601

602 Figure 7. Boxplots of the maximum of the accumulated data for 1-6 months at each month of the
603 disaggregated data with the VS model from the annual to the monthly NBS of the LCRR Basin.
604 Note that the statistics of the observed data are also represented with the dotted line and cross
605 marker (.x.); (2) the accumulation was performed for the previous months. For example, the acc-
606 4 data at Month 6 are obtained by summing the monthly data of 6, 5, 4, and 3 months.

607
608

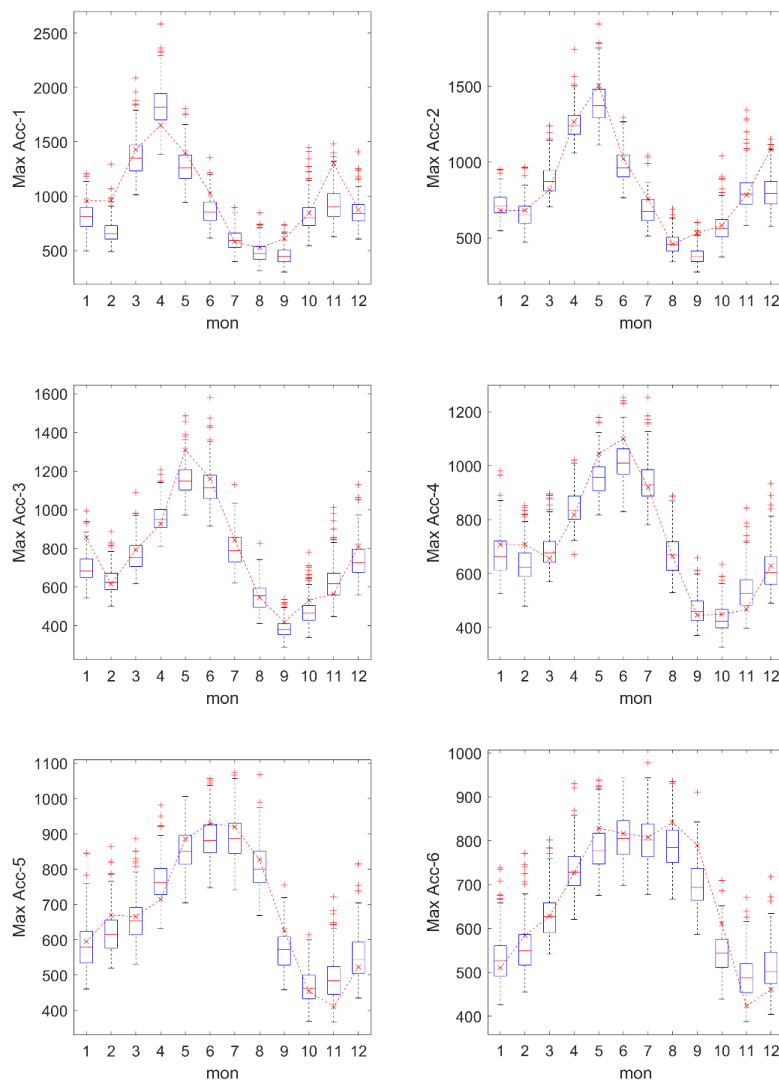


609

610 Figure 8. Same as Figure 4 but for downscaled data from the MR model.

611

612

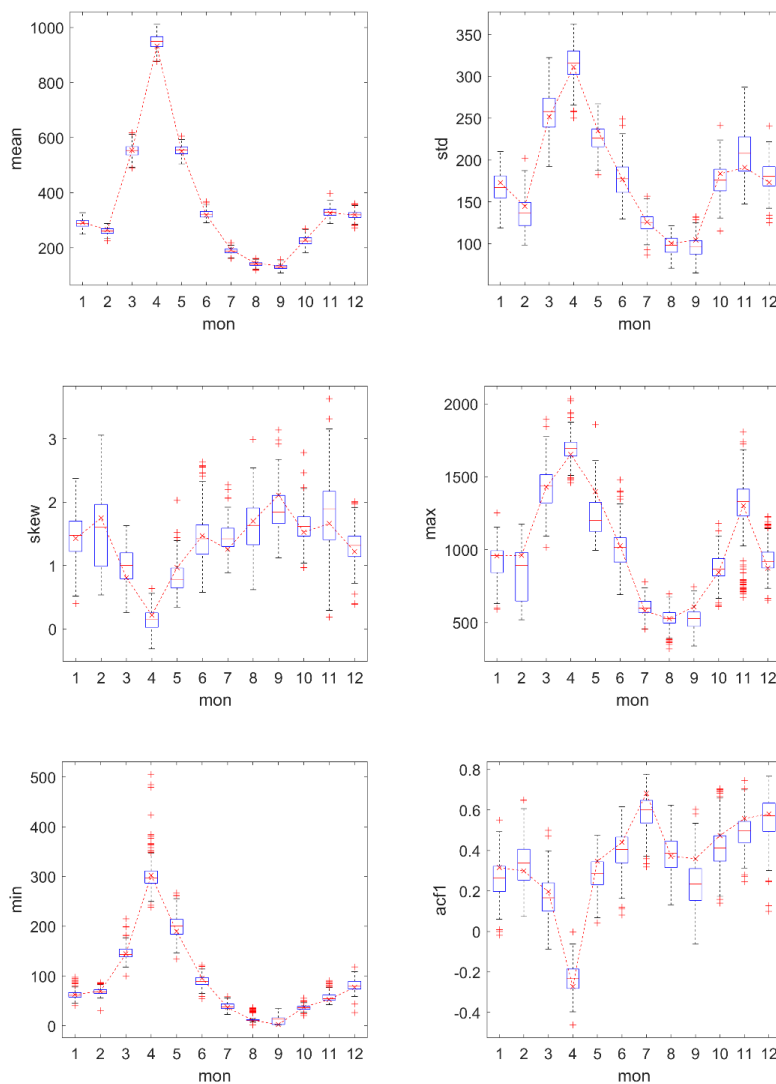


613

614 Figure 9. Same as Figure 7 but for downscaled data from the MR model.

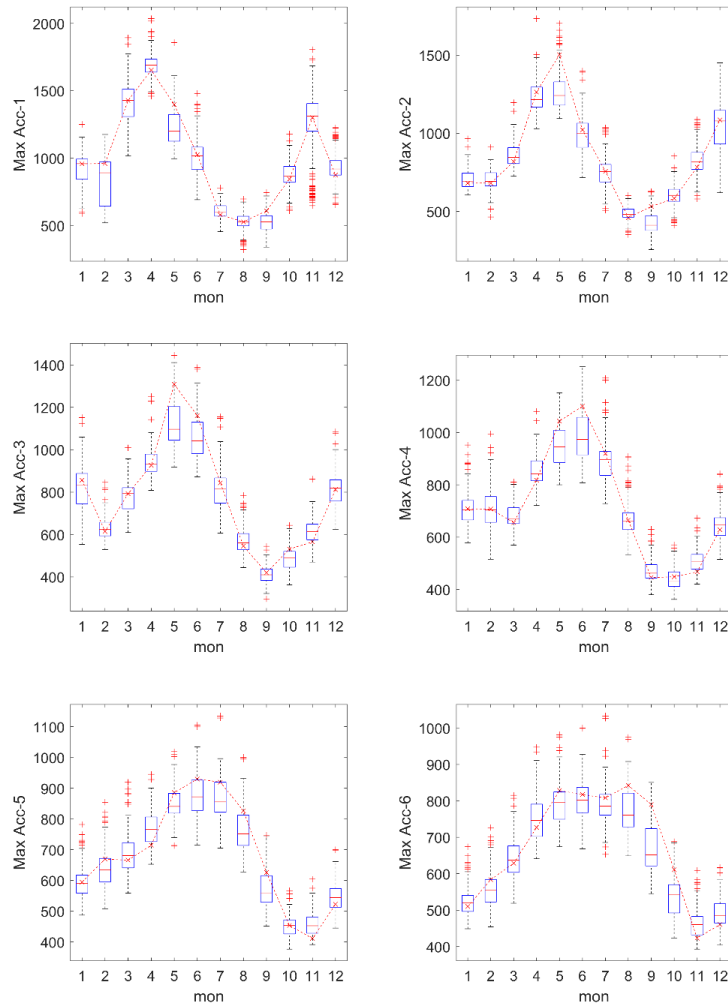
615

616



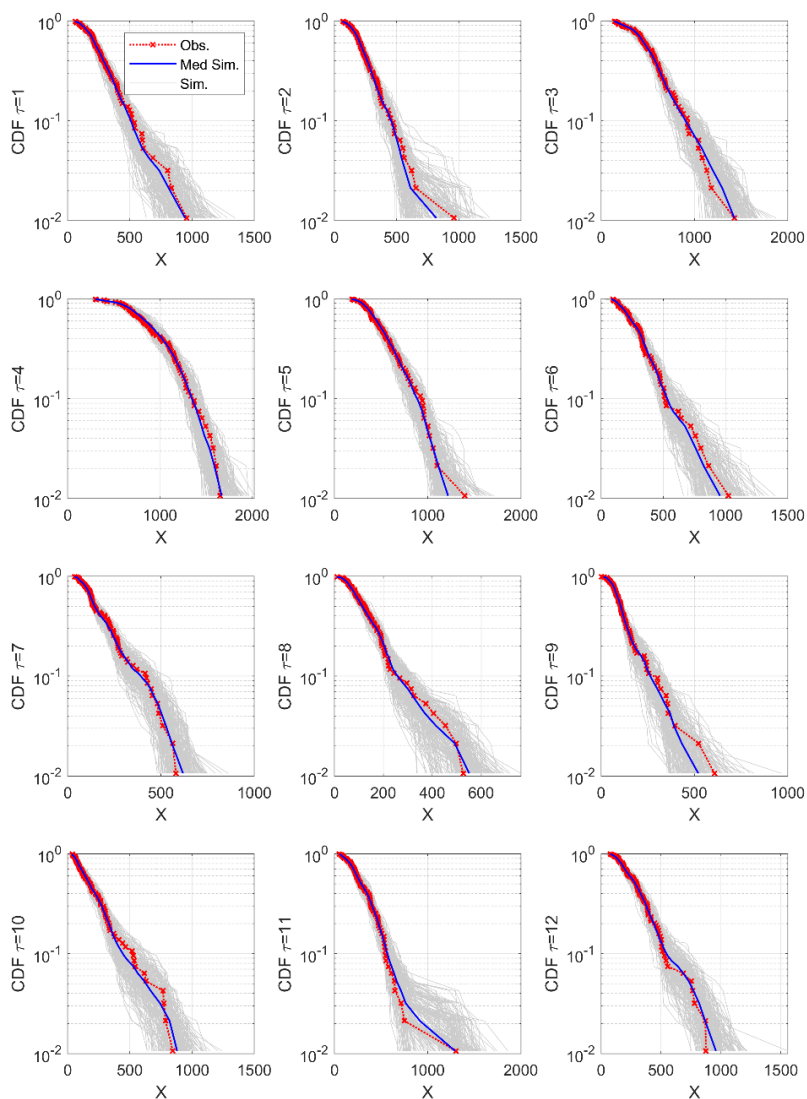
617
618 Figure 10. Boxplots of the basic statistics of the disaggregated monthly data from the annual data
619 with the NPD model with $P_c=0.1$ and $P_m=0.01$ for the NBS of the LCRR Basin. Note that the
620 statistics of the observed data are also represented with the dotted line and cross marker (.x.).

621



622
623 Figure 11. Boxplots of the maximum of the accumulated data for 1-6 months at each month of
624 the disaggregated data by the NPD model with $P_c=0.1$ and $P_m=0.01$ from the annual to the
625 monthly NBS of the LCRR Basin. Note that the statistics of the observed data are also
626 represented with the dotted line and cross marker (.x.); and (2) the accumulation was made for
627 the previous months. For example, the acc-4 data at Month 6 are obtained by summing the
628 monthly data of 6, 5, 4, and 3 months.

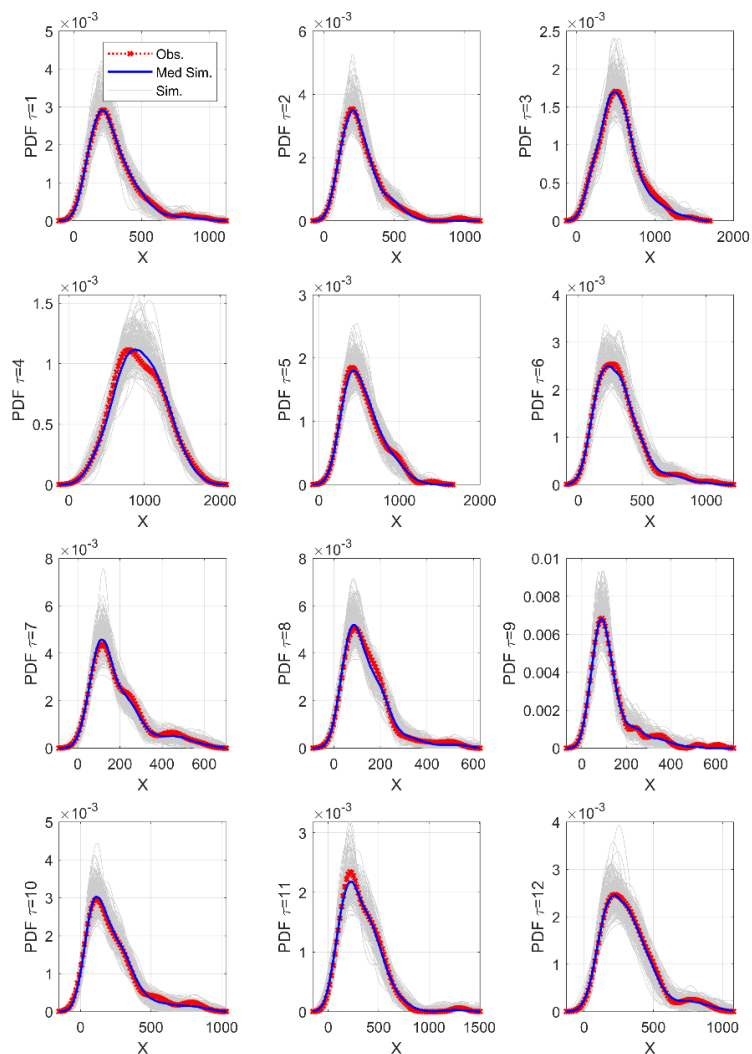
629



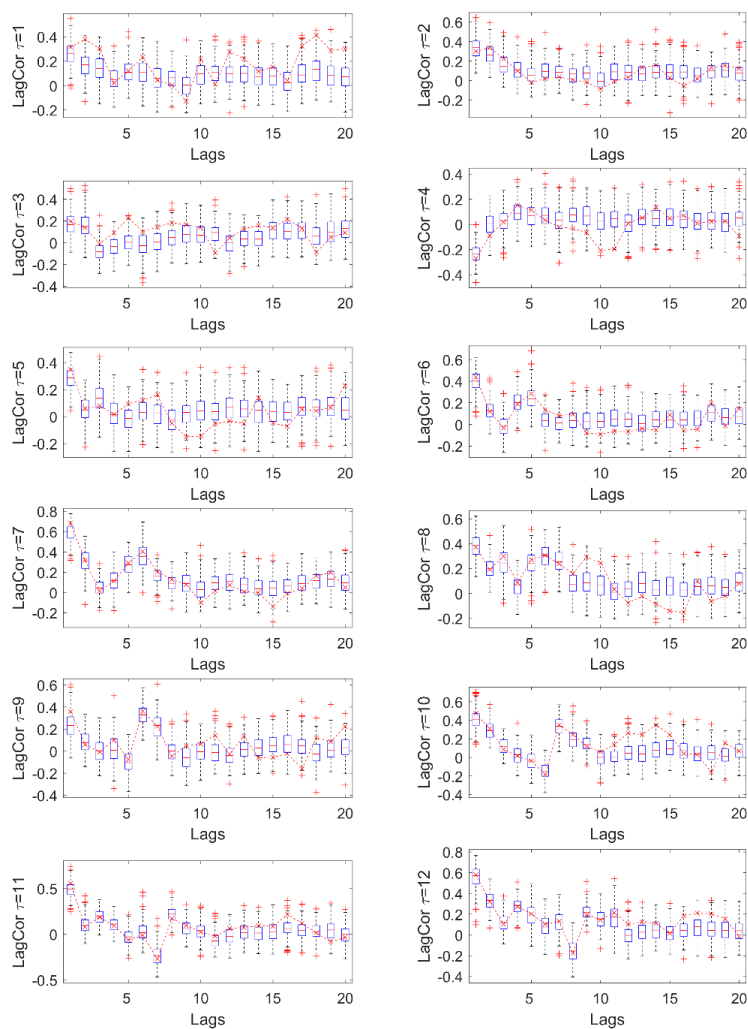
630
631 Figure 12. Cumulative distribution function (CDF) of the disaggregated data at each month with
632 the NPD model with $P_m=0.01$ and $P_c=0.1$ from the annual NBS of the LCRB Basin. Note that the
633 observed data are represented with the dotted line and cross marker (.x.); and (2) the 200
634 disaggregated simulation series are shown with the thin gray lines, while their median is
635 represented with the thick blue line.



636

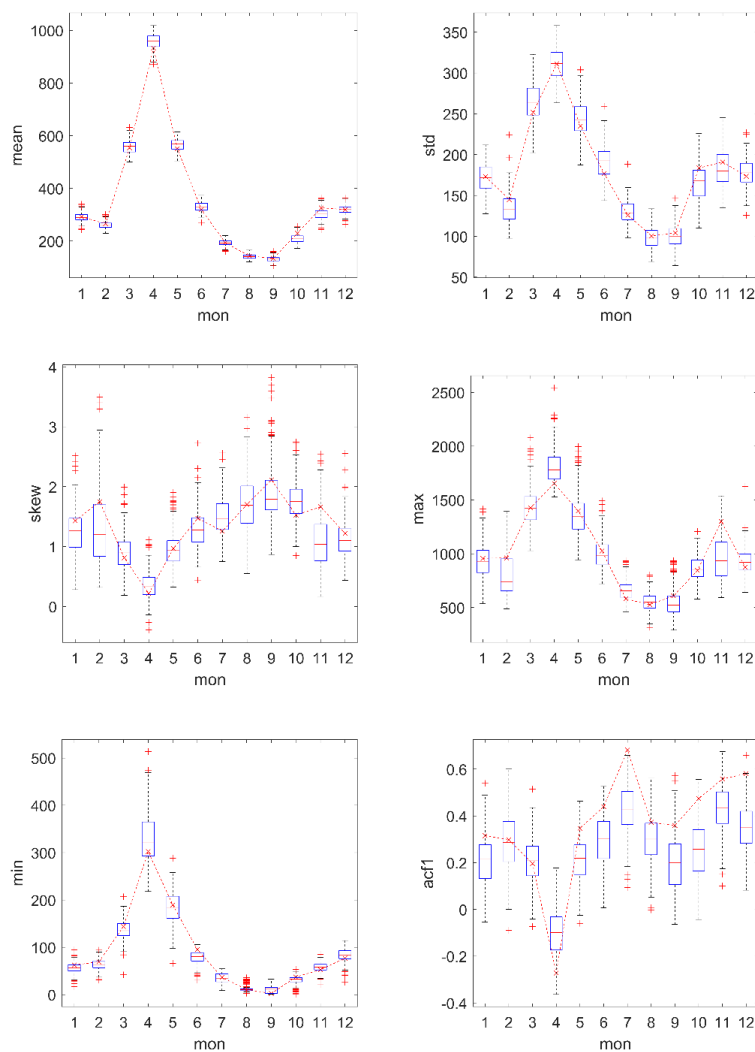


637
638 Figure 13. Probability density function (PDF) of the disaggregated data at each month with the
639 NPD model with $P_m=0.01$ and $P_c=0.1$ from the annual NBS of the LCRR Basin. Note that the
640 observed data are represented with the dotted line and cross marker (.x.); and (2) the 200
641 disaggregated simulation series are shown with the thin gray lines, while their median is
642 represented with the thick blue line.



643
644 Figure 14. Boxplots of the lagged correlation for the disaggregated monthly data with the NPD
645 model with $P_m=0.01$ and $P_c=0.1$ from the annual NBS of the LCRR Basin. Note that the statistics
646 of the observed data are also represented with the dotted line and cross marker (.x.); and (2) the
647 lagged correlation was estimated at each month. For example, the lag-2 correlation at $\tau=1$ was
648 estimated with the Month-1 data of the current year and the Month-11 data of the previous year.

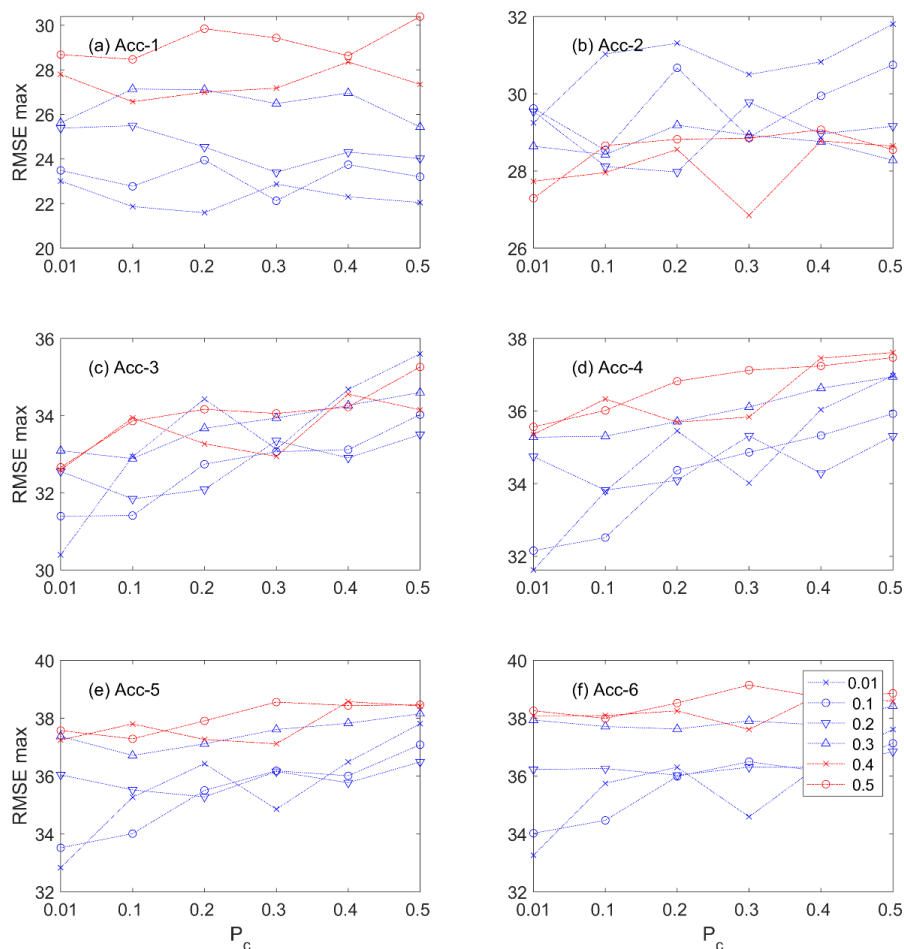
649



650
651
652
653
654

Figure 15. Boxplots of the basic statistics of the disaggregated monthly data from the annual data with the RB-NPD model with $P_c=0.3$ and $P_m=0.1$ for the NBS of the LCRB Basin. Note that the statistics of the observed data are also represented with the dotted line and cross marker (x).

655

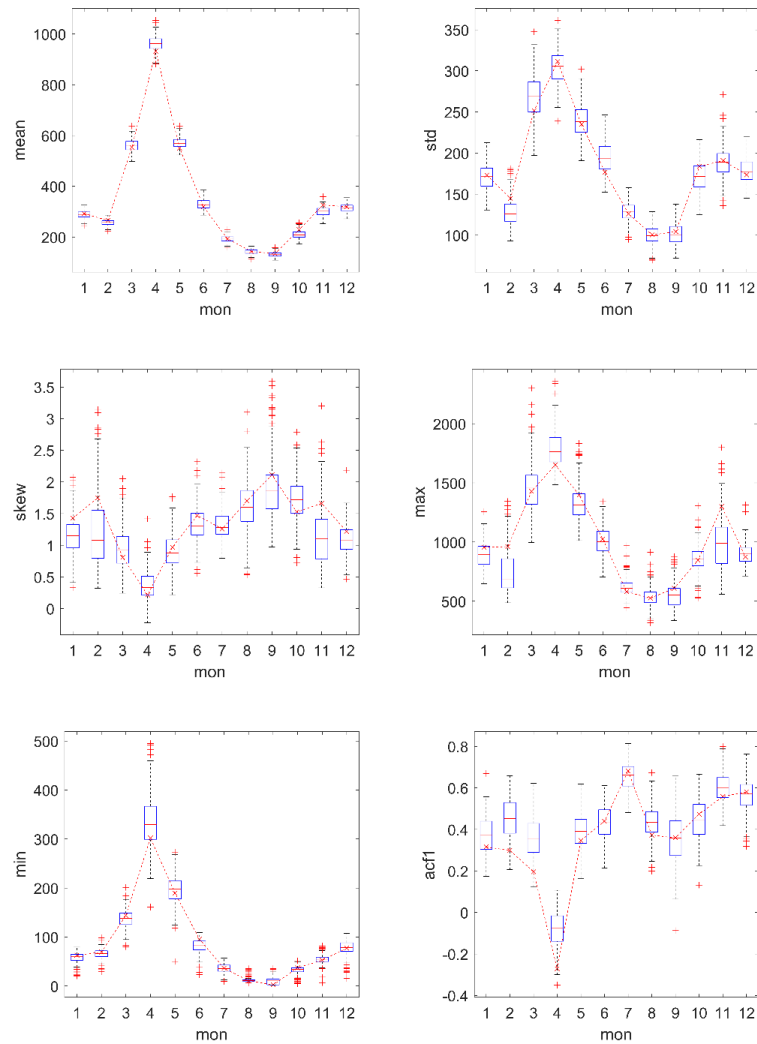


656

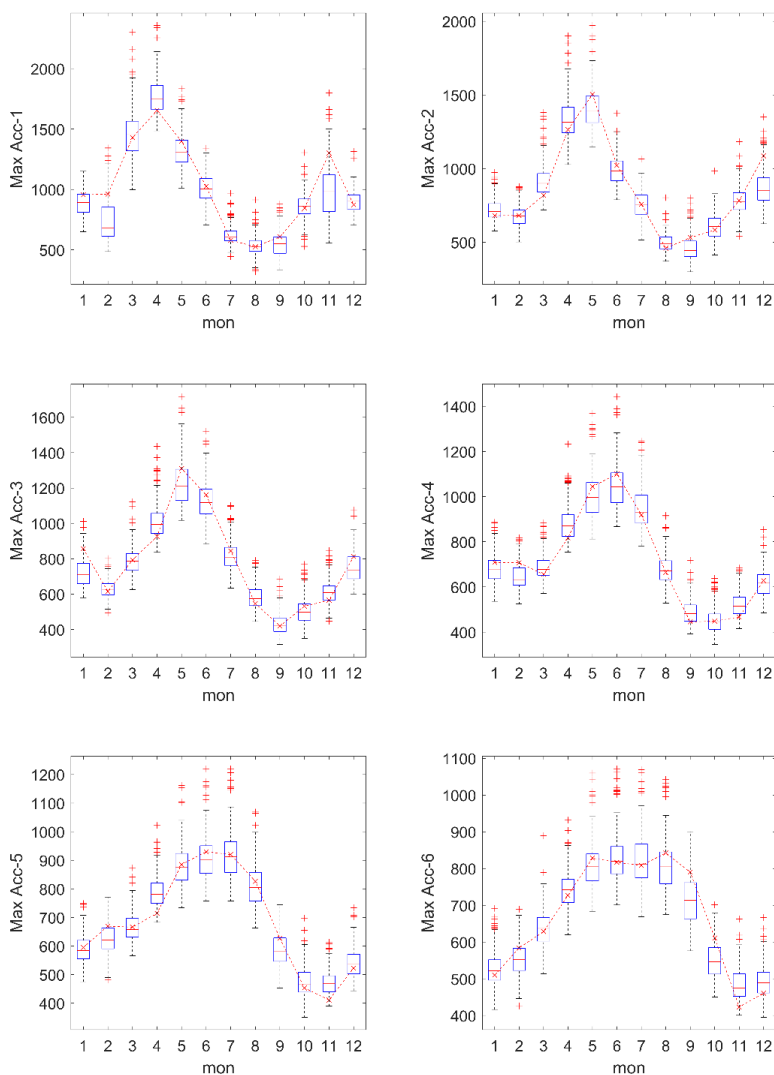
657 Figure 16. Root mean square error between the observed maximum values for the average of the
658 accumulated monthly NBS and the simulated maximum values with different crossover and
659 mutation probabilities of 0.01, 0.1, 0.2, 0.3, 0.4, and 0.5.

660

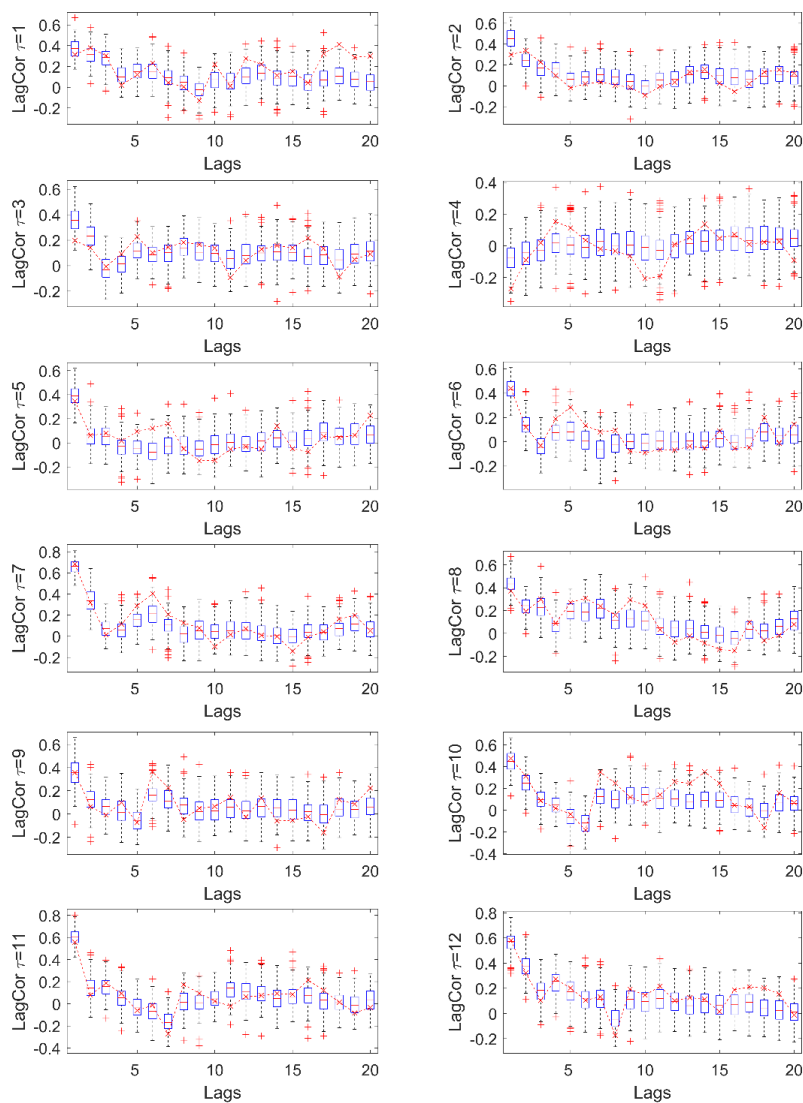
661



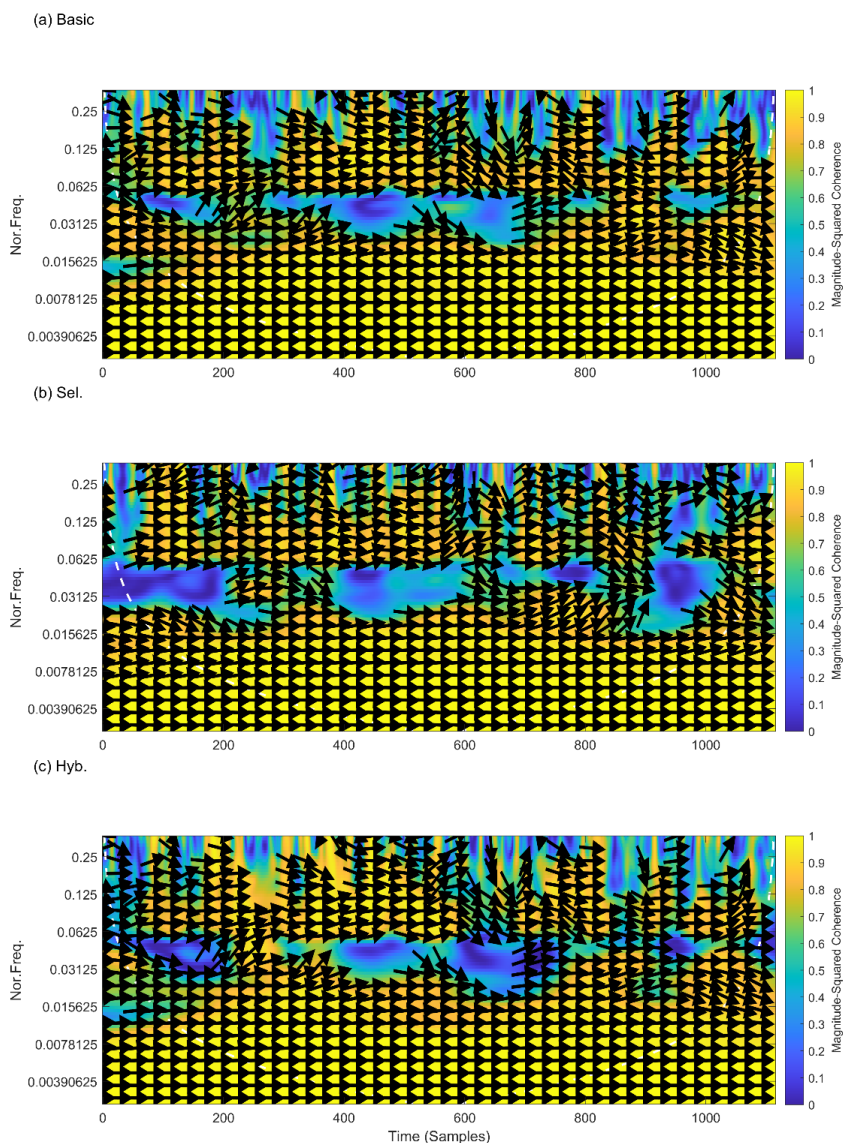
662
663 Figure 17. Boxplots of the basic statistics of the disaggregated monthly data from the annual data
664 for the RB-NPD model with ECACo+PSAm (hybrid model) as well as $P_c=0.3$ and $P_m=0.1$ for the
665 NBS of the LCRR Basin. Note that the statistics of the observed data are also represented with
666 the dotted line and cross marker (.x.).



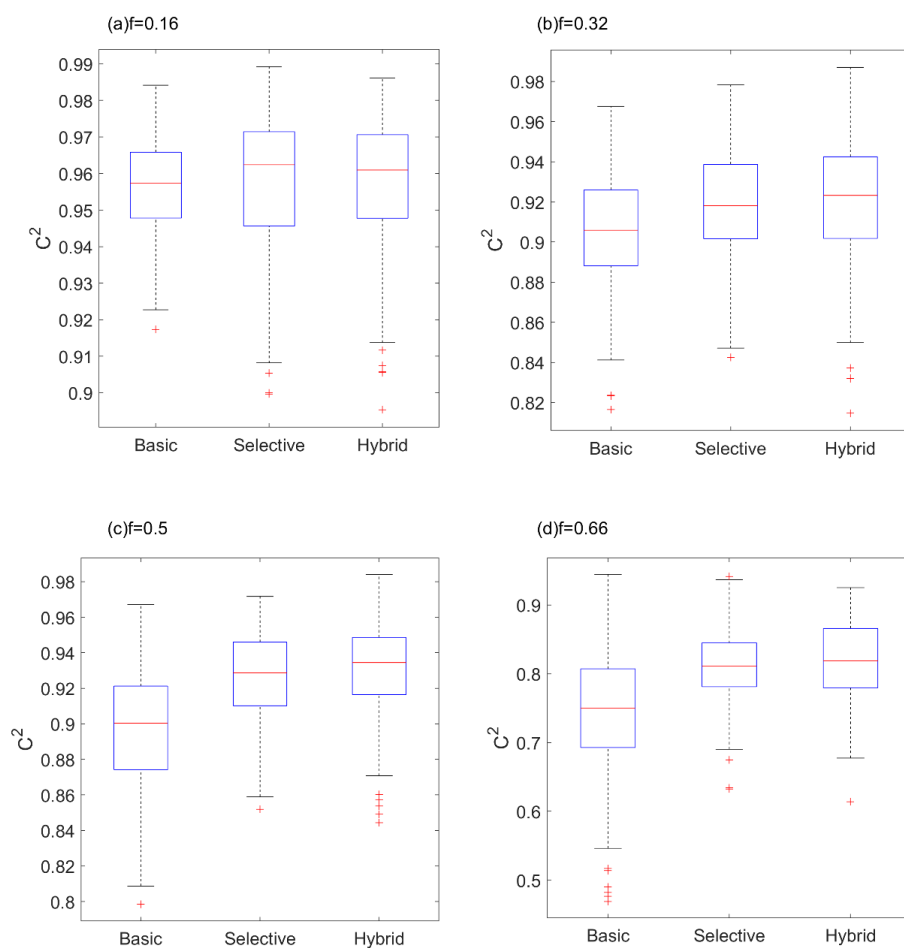
667
 668 Figure 18. Boxplots of the maximum of the accumulated data for 1-6 months at each month of
 669 the disaggregated data from the RB-NPD model with the ECAco+PSAm as well as $P_c=0.3$ and
 670 $P_m=0.1$ from the annual to the monthly NBS of the LCRR Basin. Note that the statistics of the
 671 observed data are also represented with the dotted line and cross marker (.x.); and (2) the
 672 accumulation was made for the previous months. For example, the acc-4 data at Month 6 are
 673 obtained by summing the monthly data of 6, 5, 4, and 3 months.



674
675 Figure 19. Boxplots of the lagged correlation for the disaggregated monthly data for the RB-
676 NPD model with the ECaco+PSAm as well as $P_c=0.3$ and $P_m=0.1$ from the annual NBS of the
677 LCRR Basin. Note that the statistics of the observed data are also represented with the dotted
678 line and cross marker (.x.); and (2) the lagged correlation was estimated at each month. For
679 example, the lag-2 correlation at $\tau=1$ was estimated with the Month-1 data of the current year
680 and the Month-11 data of the previous year.



681 Figure 20. Magnitude-squared coherence (C^2) of all frequencies between the observed monthly
682 NBS and the example of the disaggregated data from the RB models of the (a) basic, (b)
683 selective, and (c) hybrid algorithms. Note that (1) $C^2(f) = |S_{xy}(f)| / S_{xx}(f) \cdot S_{yy}(f)$, where
684 $S_{ab}(f)$ is the cross power spectrum of two signals, a and b , at frequency f ; (2) very strong
685 coherence can be seen in lower normalized frequencies; (3) lower frequency indicates long-term
686 variability; and (4) this is sound since the disaggregated data have the same annual values from
687 the additive condition.



688
689 Figure 21. Magnitude-squared coherence (C^2) of selected high frequencies ($f=0.16, 0.32, 0.5,$
690 0.66) between the observed monthly NBS and the disaggregated from the RB models of the
691 basic, selective, and hybrid models. Note that high coherence indicates that the model mimics the
692 spectral frequencies of the observed data well.
693
694
695
696