



**HAL**  
open science

# Stochastic downscaling of precipitation: From dry events to heavy rainfalls

M. Vrac, P. Naveau

► **To cite this version:**

M. Vrac, P. Naveau. Stochastic downscaling of precipitation: From dry events to heavy rainfalls. Water Resources Research, 2007, 43 (7), 10.1029/2006WR005308 . hal-03191925

**HAL Id: hal-03191925**

**<https://hal.science/hal-03191925>**

Submitted on 7 Apr 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Stochastic downscaling of precipitation: From dry events to heavy rainfalls

M. Vrac<sup>1,2</sup> and P. Naveau<sup>3</sup>

Received 30 June 2006; revised 22 February 2007; accepted 27 March 2007; published 4 July 2007.

[1] Downscaling precipitation is a difficult challenge for the climate community. We propose and study a new stochastic weather typing approach to perform such a task. In addition to providing accurate small and medium precipitation, our procedure possesses built-in features that allow us to model adequately extreme precipitation distributions. First, we propose a new distribution for local precipitation via a probability mixture model of Gamma and Generalized Pareto (GP) distributions. The latter one stems from Extreme Value Theory (EVT). The performance of this mixture is tested on real and simulated data, and also compared to classical rainfall densities. Then our downscaling method, extending the recently developed nonhomogeneous stochastic weather typing approach, is presented. It can be summarized as a three-step program. First, regional weather precipitation patterns are constructed through a hierarchical ascending clustering method. Second, daily transitions among our precipitation patterns are represented by a nonhomogeneous Markov model influenced by large-scale atmospheric variables like NCEP reanalyses. Third, conditionally on these regional patterns, precipitation occurrence and intensity distributions are modeled as statistical mixtures. Precipitation amplitudes are assumed to follow our mixture of Gamma and GP densities. The proposed downscaling approach is applied to 37 weather stations in Illinois and compared to various possible parameterizations and to a direct modeling. Model selection procedures show that choosing one GP distribution shape parameter per pattern for all stations provides the best rainfall representation amongst all tested models. This work highlights the importance of EVT distributions to improve the modeling and downscaling of local extreme precipitations.

**Citation:** Vrac, M., and P. Naveau (2007), Stochastic downscaling of precipitation: From dry events to heavy rainfalls, *Water Resour. Res.*, 43, W07402, doi:10.1029/2006WR005308.

### 1. Introduction

[2] In recent decades, the accuracy of general circulation models (GCM) to simulate the large-scale behavior of the atmosphere has greatly improved. Still, such models have difficulties capturing small-scale intermittent processes, for example, local precipitation. To better understand and represent these sub-grid-scale meteorological characteristics, Regional Climate Models (RCM) offer an elegant way to integrate local processes through physical and dynamical equations. However, they can be extremely computer-intensive and their spatial resolution, generally from 5 to 50 km, does not always provide the required information needed in impact studies. Again, local precipitation can be considered as the archetypical example of such limitations. While advances in computer sciences may give the necessary computer power to resolve these smaller

scales in the future, practitioners (flood planners, insurance companies, etc.) need to make decisions locally with the current information today.

[3] In order to link our large-scale knowledge supplied by today's GCM, RCM and reanalysis outputs with measurements recorded at weather stations, statistical downscaling techniques offer a computationally attractive and ready-to-use route. This statistical approach consists of inferring significant relationships among large-, regional- and local-scale variables. How to estimate, apply and test such relationships in order to have accurate representations of local features constitutes the so-called group of statistical downscaling questions. Three categories of methods are usually given to answer such questions: transfer functions, stochastic weather generators and weather typing methods. The first category is a direct approach. The relationships between large-scale variables and location-specific values are directly estimated via either parametric, nonparametric, linear or nonlinear methods such as the analog method [e.g., Barnett and Preisendorfer, 1978; Zorita and von Storch, 1998], multiple linear regressions [e.g., Wigley et al., 1990; Huth, 2002], kriging [e.g., Biau et al., 1999] and neural networks [e.g., Snell et al., 2000; Cannon and Whitfield, 2002]. The second category focuses on weather generators in which GCM outputs drive stochastic models of precipitation [e.g., Wilks, 1999; Wilks and Wilby, 1999].

<sup>1</sup>Center for Integrating Statistical and Environmental Science, University of Chicago, Chicago, Illinois, USA.

<sup>2</sup>Now at Laboratoire des Sciences du Climat et de l'Environnement, IPSL-CNRS, CEA Saclay, Gif-sur-Yvette, France.

<sup>3</sup>Laboratoire des Sciences du Climat et de l'Environnement, IPSL-CNRS, CEA Saclay, Gif-sur-Yvette, France.

They are particularly of interest to assess local climate change [e.g., *Semenov and Barrow*, 1997; *Semenov et al.*, 1998]. The weather typing approach, the third and last category, encapsulates a wide range of methods that have in common an algorithmic step in which recurrent large-scale and/or regional atmospheric patterns are identified. These patterns are usually obtained from clustering and classification algorithms applied to geopotential height, pressure or other meaningful atmospheric variables over a large spatial area. These clustering and classification algorithms can be of different types: CART (Classification and Regression Trees) [see *Breiman et al.*, 1984; *Schnur and Lettenmaier*, 1998], “K-means” methods [e.g., *Huth*, 2001; *Yiou and Nogaj*, 2004], hierarchical clustering approaches [e.g., *Davis et al.*, 1993; *Bunkers et al.*, 1996], fuzzy-rules-based procedures [e.g., *Pongracz et al.*, 2001], neural networks [e.g., *Bárdossy et al.*, 1994] or mixture of copula functions [*Vrac et al.*, 2005]. Introducing such an intermediate layer (the weather patterns) in a downscaling procedure provides a strong modeling flexibility. For example, linking directly the relationships between large-scale atmospheric variables and precipitation recorded at a few weather stations may be too complex in most inhabited regions. In comparison, it may be easier and more efficient to first model the dependences between large-scale data and weather patterns, the latter representing the recurrent atmospheric structures corresponding to a kind of summary of the large scale. Then we can focus on the coupling between weather patterns and local measurements. Obviously, such a strategy will only be successful if the weather patterns are carefully chosen, i.e., if they capture relevant recurrent summary information. From a probabilistic point of view, the coupling step of a weather typing approach can be viewed as deriving the following conditional probability density function (pdf):

$$f_{\mathbf{R}_t | \mathbf{S}_t}, \quad (1)$$

which corresponds to the probability of observing local rainfall intensities, say  $\mathbf{R}_t$ , given the current weather state, say  $\mathbf{S}_t$ , at time  $t$ . In addition to providing a simple mathematical framework that can easily integrate various uncertainties, this probabilistic definition of statistical downscaling is wide enough to cover many case studies. In this work, to get more realistic precipitation variability than with a model only conditional on weather patterns, the pdf (1) is also defined conditionally on a vector of large-scale atmospheric variables, say  $\mathbf{X}_t$ , at time  $t$ ,

$$f_{\mathbf{R}_t | \mathbf{X}_t, \mathbf{S}_t}. \quad (2)$$

[4] In this paper, our main application is to downscale precipitation over the region of Illinois. Consequently, we would like to address the following questions: how to find adequate regional weather patterns for  $\mathbf{S}_t$ ? How to model the coupling between large atmospheric variables  $\mathbf{X}_t$  and  $\mathbf{S}_t$ ? What is an appropriate form for the conditional density defined by (2)? The last question is the central one for the practitioner.

[5] To our knowledge, none of the statistical downscaling methods discussed previously in this section has been developed to address the issue of modeling both common

and extreme values. Nevertheless, although, for example, hydrologists and flood planners are interested in mean precipitation, they also have a particular interest in modeling extreme local precipitation because of its human, economical and hydrological impacts where large-scale information may help at modeling such extreme events. Past studies [*Katz et al.*, 2002; *Naveau et al.*, 2005] have illustrated how Extreme Value Theory (EVT), a statistical theory developed over the past 80 years, provides the mathematical foundation for appropriately modeling extreme precipitation. Hence another important objective in this paper is to integrate EVT models within a weather typing approach, i.e., throughout the density (2). To perform such a task, we extend the original work on the nonhomogeneous stochastic weather typing approach by *Vrac et al.* [2007].

[6] The paper is organized as follows. In the first part of section 2, we recall three classical distribution candidates that have been proposed to fit rainfall and we also introduce a mixture model inspired by *Frigessi et al.* [2003]. A comparison and a discussion about the performance of these four distributions is undertaken. In section 3 the full data sets are presented. Regional precipitation-related patterns are obtained by applying a hierarchical ascending clustering (HAC) algorithm to observed precipitation. Then our statistical downscaling model is explained. Section 4 contains results about our application and many different diagnostics are computed to assess the quality of the models and to select the most appropriate one. All along this section, instead of “pure” GCM outputs as large-scale atmospheric variables, we take advantage of reanalysis data from the National Centers for Environmental Prediction (NCEP). Indeed, not only are NCEP reanalyses constrained GCM outputs, but also, using NCEP is necessary to assess our daily downscaling method in a present climate, before fitting the method to (pure) GCM outputs to project local change in precipitation. Hence, because the motivation is driven by the scale transformation of large-scale atmospheric variables (GCM outputs or reanalysis data), working on reanalyses is a first essential step. Lastly, in section 5, we conclude and give some future research directions.

## 2. Modeling Rainfall Locally

[7] There exists a wide range of distribution families to statistically model rainfall intensities. For example, *Katz* [1977], *Wilks* [1999], *Bellone et al.* [2000], *Vrac et al.* [2007], and *Wilks* [2006] argued that most of the precipitation variability can be approximated by a Gamma distribution. However, it is also well known [e.g., *Katz et al.*, 2002] that the tail of this distribution can be too light to capture heavy rainfall intensities. This leads to the underestimation of return levels and other quantities linked to high percentiles of precipitation amounts. Consequently, the societal and economical impacts associated with heavy rains (e.g., floods) can be miscalculated. To solve this issue, an increasingly popular approach in hydrology [*Katz et al.*, 2002] is to disregard small precipitation values and to focus only on the largest rainfall amounts. The advantage of this strategy is that an elegant mathematical framework called Extreme Value theory (EVT) developed in 1928 [*Fisher and Tippett*, 1928] and regularly updated during the last decades [e.g., *Coles*, 2001] dictates the distribution of heavy precip-

itation. More specifically, EVT states that rainfall exceedances, i.e., amounts of rain greater than a given threshold  $u$ , can be approximated by a Generalized Pareto Distribution (GPD) if the threshold and the number of observations are large enough. In other words, the probability that the rainfall amount, say  $R$ , is greater than  $r$  given that  $R > u$  is given by

$$P(R > r | R > u) = \left(1 + \xi \frac{r - u}{\sigma}\right)_+^{-1/\xi}, \quad (3)$$

where  $a_+ = \max(a, 0)$  and  $\sigma > 0$  represents the scale parameter. The shape parameter  $\xi$  describes the GPD tail behavior. If  $\xi$  is negative, the upper tail is bounded. If  $\xi$  is zero, this corresponds to the case of an exponential distribution (all moments are finite). If  $\xi$  is positive, the upper tail is still unbounded but higher moments eventually become infinite. These three cases are termed “bounded”, “light-tailed”, and “heavy-tailed”, respectively. The flexibility of the GPD to describe three different types of tail behavior makes it a universal tool for modeling exceedances. Although this GPD approach has been very successful to model heavy rains, it has the important drawback of overlooking small precipitation. Recently, *Wilson and Toumi* [2005] proposed a new probability distribution for heavy rainfall by invoking a simplified water balance equation. They claimed that the stretched exponential distribution tail defined by

$$P(R > r) = \exp\left[-\left(\frac{r}{\psi}\right)^\nu\right], \quad (4)$$

where  $\psi > 0$  and  $\nu > 0$  correspond to the scale and shape parameter. The latter should be equal to  $\nu = 2/3$ . This was justified by physical arguments that take into account of the distributions probabilities of quantities like the upward wind velocity  $\mathbf{w}$  (although the distribution of  $\mathbf{w}$  is much more unknown than the distribution of  $R$ ). Note also that, although the parameter  $\nu$  is expected to be equal to  $2/3$  in theory, *Wilson and Toumi* did not say that in practice this parameter has to be equal to  $2/3$ . Indeed, they estimated the shape parameter from different weather station precipitation measurements over the world. They found that, in practical applications, the estimated shape parameter is usually different from the  $2/3$  constant. Despite its drawbacks, such a type of model is promising because it tries to combine probabilistic reasoning with physical arguments. Still, it is not designed to model small precipitation amounts. For their main example, *Wilson and Toumi* [2005] estimated the parameter  $(\psi, \nu)$  in (4) for “heavy precipitation defined as daily totals with probability less than 5%”. Hence one may wonder how to deal with the remaining 95% and what is the justification for working with 5% of the data and not 10%, 3% or any small percentages (this later problem also exists with a classical EVT approach). Because our final objective is to downscale the full range of precipitation values and because we do not want to choose an arbitrarily preset threshold (or percentage), we follow a different direction and opt for the method proposed by *Frigessi et al.* [2003]. These authors introduced the following mixture model:

$$h_\beta(r) = c(\beta) \times \left[ (1 - w_{m,\tau}(r)) \times f_{\beta_0}(r) + w_{m,\tau}(r) \times g_{\xi,\sigma}(r) \right], \quad (5)$$

where  $c(\beta)$  is a normalizing constant,  $\beta = (m, \tau, \beta_0, \xi, \sigma)$  encapsulates the vector of unknown parameters,  $f_{\beta_0}$  corresponds to a light-tailed density with parameters  $\beta_0$ , the function  $g_{\xi,\sigma}$  represents the GPD density that can be obtained from deriving the tail defined by (3) and  $w_{m,\tau}(\cdot)$  is a weight function that depends on two parameters,

$$w_{m,\tau}(r) = \frac{1}{2} + \frac{1}{\pi} \arctan\left(\frac{r - m}{\tau}\right). \quad (6)$$

[8] Note that this weight function is nondecreasing, takes values in  $(0,1]$  and tends to 1 as  $r$  goes to  $\infty$ ; that is, heavy rains are represented by the GPD density  $g_{\xi,\sigma}(r)$  in the mixture  $h_\beta(r)$  for large  $r$ . Conversely, small precipitation values are mostly captured by the light-tailed density  $f_{\beta_0}(r)$ . Hence the idea behind equations (5) and (6) is rather simple: the mixing function  $w_{m,\tau}(r)$  provides a smooth transition from a light-tailed density (small and medium precipitation) to the GPD density (heavy rainfalls). The parameters  $m$  and  $\tau$  in  $w_{m,\tau}(r)$  correspond to the location and the speed of the transition from  $f_{\beta_0}$  to  $g_{\xi,\sigma}$  in (5), respectively. In 2003, *Frigessi et al.* [2003] applied their model to Danish fire loss data and opted for a Weibull distribution as a light-tailed density in (5). In the context of precipitation modeling, past works [*Bellone et al.*, 2000; *Vrac et al.*, 2007; *Wilks*, 2006] indicate that a Gamma density, i.e.,

$$f_{\beta_0}(x) = \frac{1}{\lambda^\gamma \Gamma(\gamma)} x^{\gamma-1} \exp(-x/\lambda), \text{ with } \beta_0 = (\gamma, \lambda), \quad (7)$$

should fit appropriately the bulk of the precipitation values (heavy rains excluded). This hypothesis could be challenged if the variable of interest was different, for example, temperature. In addition, one may be puzzled by the “absence” of a threshold in equation (5). Indeed, the threshold  $u$  in equation (3) is forced to be equal to zero in (5). However, introducing the weight function  $w_{m,\tau}(r)$  and fixing the GPD threshold to zero brings two important benefits. First, the difficult threshold selection problem is replaced by a simpler unsupervised estimation procedure, i.e., finding  $m$  and  $\tau$  in  $w_{m,\tau}(r)$  from the data. This strategy is particularly relevant to large data sets analysis because it would be very time-consuming to find an adequate threshold for a large number of weather stations. Second, allowing for nonzero thresholds in (5) would impose an unwelcome discontinuity in  $h_\beta(r)$ . From a physical point of view, such a discontinuity represents an unrealistic feature in precipitation.

[9] In summary, we have four candidates for modeling local rainfall distribution: (1) the Gamma density that works well for the main rainfall range but not for large values; (2) the recently introduced stretched-exponential distribution function defined by equation (4), constructed on a physical foundation but only designed for heavy rainfall and not for small precipitation values; (3) the GPD function that works for extreme precipitation but not for small values, that is mathematically sound and universal, in the sense that it can also fit temperature, winds extremes, etc.; and (4) and our new mixture model defined by (5) and (7) that combines the advantages of the Gamma and GPD densities, and consequently can fit small and heavy rainfall.

[10] To compare the performances of these four distributions, we implement the following procedure. We simulate



**Table 1.** Frequencies of Selections of the four candidate Distributions by the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) Values Obtained From 100 Samples of 1000 Simulated Data for Each Given Density<sup>a</sup>

| True Density          | Fitted Density                      |                                    |                                    |                                    |
|-----------------------|-------------------------------------|------------------------------------|------------------------------------|------------------------------------|
|                       | Gamma                               | GP                                 | Mixture                            | Stretched                          |
| Gamma                 | <b>AIC = 90</b><br><b>BIC = 100</b> | AIC = 0<br>BIC = 0                 | AIC = 10<br>BIC = 0                | AIC = 0<br>BIC = 0                 |
| Generalized-Pareto    | AIC = 0<br>BIC = 0                  | <b>AIC = 86</b><br><b>BIC = 96</b> | AIC = 11<br>BIC = 1                | AIC = 3<br>BIC = 3                 |
| Mixture: GP + Gamma   | AIC = 7<br>BIC = 36                 | AIC = 0<br>BIC = 0                 | <b>AIC = 93</b><br><b>BIC = 64</b> | AIC = 0<br>BIC = 0                 |
| Stretched exponential | AIC = 3<br>BIC = 3                  | AIC = 0<br>BIC = 0                 | AIC = 10<br>BIC = 0                | <b>AIC = 87</b><br><b>BIC = 97</b> |

<sup>a</sup>The bold fonts correspond to the highest frequencies with respect to the AIC and the BIC.

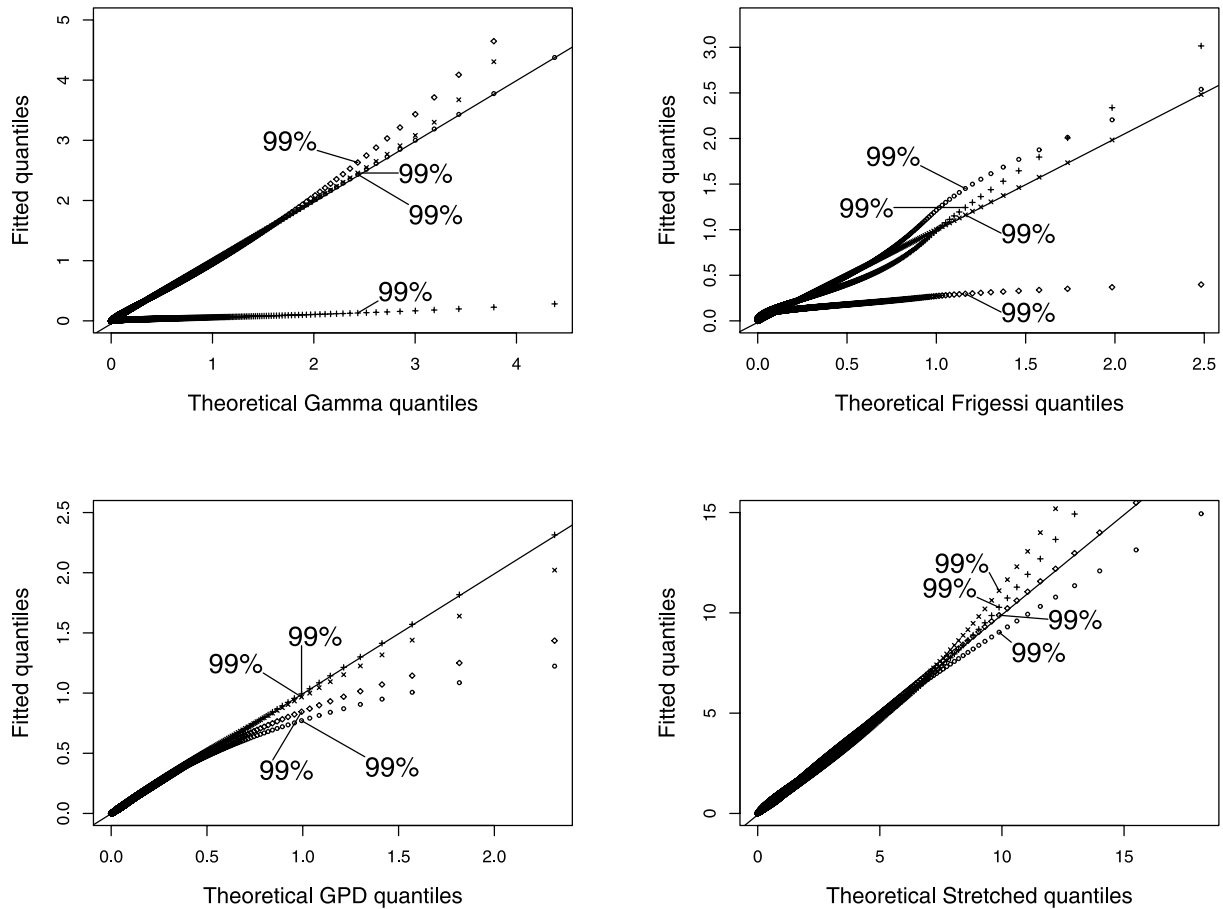
100 samples of 1000 iid realizations of each density with:  $\lambda = 1$  and  $\gamma = 0.25$  for the Gamma distribution (see equation (7)),  $u = 0$ ,  $\xi = 0.3$ ,  $\sigma = 0.1$  for the GPD (see equation (3)),  $u = 0$ ,  $m = 1$ ,  $\tau = 0.1$  for the mixture of the two previous distributions (see equation (5)), and  $\nu = 2/3$  and  $\psi = 1$  for the stretched exponential (see equation (4)), respectively. Such parameter values were chosen because they correspond to reasonable estimates for precipitation data. In particular,  $\nu = 2/3$  is recommended by *Wilson and Toumi [2005]*. As a second step, we fit each distribution to each of the four simulated samples by using the maximum likelihood approach to compute the “optimal” parameters for each distribution. To be consistent with Wilson and Toumi’s paper, the parameter  $\nu$  in (4) is not considered as a constant; that is, we assume that this shape parameter has to be estimated. This has also the advantage that we do not penalize the stretched exponential distribution with respect to the other distributions we test and for which the shape parameter is also not fixed but estimated. The last step is to compare the qualities of the fit with respect to the given density. Classically, one can compute the Akaike information criterion (AIC) [*Akaike, 1974*], defined by  $-2 \log(L) + 2p$ , and the Bayesian information criterion (BIC) [*Schwarz, 1978*], defined by  $-2 \log(L) + p \log(n)$ , where  $L$  is the likelihood of the model fitted to the data,  $p$  is the number of parameters, and  $n$  is the number of data. Minimizing AIC and BIC helps to select the model with a good fit to the data (i.e., high likelihood) while penalizing a model with too many parameters. The BIC tends to add a larger parameter cost than the AIC. For our simulations, the frequencies of selection of the four candidate distributions by the AIC and BIC values are summarized in Table 1. As expected, the best AIC and BIC (in bold) are majoritarily obtained along the diagonal of the table; that is, the simulated samples are best fitted by the density from which they were generated. We can remark that about one time every third, the BIC indicates a Gamma fit when the true density is a mixture; that is, the BIC penalizes too much. In comparison, the AIC largely selects the correct distribution for all four cases.

[11] Hence, for these simulations, the AIC appears to perform reasonably and will be used in the subsequent analyses. Still, we cannot solely rely on these two criteria

to discriminate among models. In particular, these criteria may not be well adapted for extreme values. Concerning the fit quality of the largest values, Figure 1 displays four quantile-quantile type plots (QQplots). The circles, crosses, pluses, and diamonds correspond to the analytically fitted Gamma, mixture, GP and stretched exponential densities, respectively. The  $y = x$  black line represents the “true” distribution that can either be a Gamma (Figure 1, top left), a mixture (Figure 1, top right), a GP (Figure 1, bottom left) and a stretched exponential (Figure 1, bottom right) density. This graph mainly tells us that the mixture distribution (crosses) appears to provide a very good fit in all cases. As expected, a Gamma fit (circles) does not work very well when the true trail is heavy. The stretched exponential (diamonds) is somehow limited because it only provides a good fit when the true tail is stretched exponential. The worst case is the GPD (pluses), but this is expected because the threshold  $u$  was set to zero and it is well known that the GPD only works well for very large values. An alternative would be to select a high threshold, but then the main part of rainfall cannot be statistically modeled (and consequently, be compared with the other densities). Still, it is very interesting to see that, despite also having a GPD threshold set to zero, the mixture density provides very good results. This reveals that the weight function  $w_{m,\tau}$  in (6) can bring enough flexibility even if the mixture threshold is equal to zero. One may argue that the mixture density has too many parameters, but the AIC and BIC summarized in Table 1 do not show much cases of overfitting. Even more importantly, Figure 1 shows that the other three classical distributions for rainfall (Gamma, stretched exponential and GPD) do not offer the necessary latitude to model the full spectrum of precipitation distribution.

[12] Although the scope of this small simulation study is very limited and a more thorough investigation would be welcome to review the arguments and problems related to local rainfall distributions, Table 1 and Figure 1 strongly suggest that our mixture model could provide a competitive probabilistic foundation. Consequently, this model will be used in the rest of this paper. Concerning the choice between the AIC and BIC, only the AIC will be presented in the remainder of this paper. In most cases, the BIC provides similar results and does not change the meaning of the main findings that will be presented in section 3.

[13] With respect to real data, our goal is to analyze daily observations that were recorded at 37 weather stations in Illinois from 1980 to 1999. Those stations correspond to the complete data set of precipitation provided for Illinois by the co-op observational program. The stations are found to be uniformly distributed over Illinois. To reduce seasonal influences, we only consider three winter months, December, January and February (DJF). To illustrate the fit between our mixture model and real rainfall observations and also to show the difference of fit to the data between a Gamma distribution and our mixture, we select one station (Aledo) and apply a maximum likelihood estimation procedure to derive the parameters of each distribution. Figure 2 shows the resulting quantile-quantile plots. Figure 2 (top) displays the fit obtained using a Gamma distribution, while Figure 2 (bottom) shows the result for our mixture distribution. As already seen in our simulation study, this latter model provides a gain at capturing extreme values behavior.



**Figure 1.** Quantiles-quantiles plots (i.e., theoretical versus fitted quantiles). The circles, crosses, pluses, and diamonds correspond to the QQplots from the Gamma, mixture, GP, and stretched exponential densities, respectively. Each distribution is analytically fitted by (top left) a Gamma, (top right) our mixture, (bottom left) a GP, and (bottom right) a stretched exponential density. The 99% quantile is indicated for each fitted distribution. These graphs mainly tell us that the mixture distribution (crosses) appears to provide a very good fit in all cases.

At this stage, one could be satisfied by this type of station-per-station analysis. However, from a statistical and physical point of views, we prefer to go a step further in our statistical analysis by relating local precipitation with large-scale variables through an extension of our mixture model. This is the object of the following section.

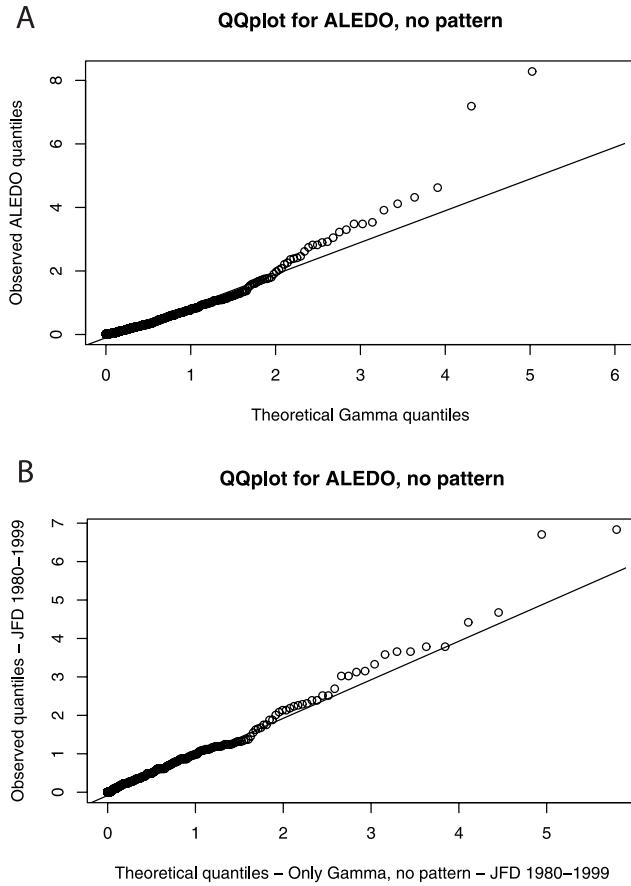
### 3. Our Downscaling Procedure

[14] To develop a statistical model capable of downscaling precipitation, we need large-scale atmospheric variables and local observed precipitation measurements. The latter are provided here by daily observations described in section 2. Large-scale atmospheric variables are given by NCEP reanalysis data, with a  $2.5^\circ \times 2.5^\circ$  spatial resolution and at 850 mb. Three NCEP variables are considered in our analysis: geopotential height denoted  $Z_{850}$ , specific humidity,  $Q_{850}$ , and dew point temperature depression  $\Delta T_{d850}$  defined as  $T_{850} - T_{d850}$ , where  $T_{850}$  and  $T_{d850}$  are the temperature and dew point temperature at 850 mb, respectively.

#### 3.1. Modeling Regional-Scale Precipitation Patterns

[15] Classically, weather typing methods are based on circulation-related patterns. A number of studies [e.g.,

*Mamassis and Koutsoyiannis, 1996*] showed that, according to the studied region, large-scale atmospheric patterns can be efficient to explain and characterize local precipitation variability. However, to better represent precipitation behaviors, we follow the approach of *Vrac et al. [2007]*. Instead of defining upper air circulation patterns, these authors recently constructed precipitation-related patterns, directly obtained from a subset of observed local precipitations, and showed that, for Illinois, these patterns are more efficient than classical upper air circulation patterns to characterize and simulate local precipitation. These precipitation patterns were derived from a hierarchical ascending clustering (HAC) algorithm with Ward criterion [*Ward, 1963*], applied to the observed precipitation of the 1980–1999 winter months (DJF). Instead of the common Euclidean distance, a special metric tailored to precipitation was developed to take account of the spatio-temporal rain features. The details of this clustering algorithm are given by *Vrac et al. [2007]*. Figure 3 shows the four precipitation patterns over the region of Illinois. It is clear that pattern 1 represents the smallest rainfall intensities whereas pattern 4 corresponds to the most intense precipitation. Patterns 2 and 3 show moderate precipitation, with opposite South/North



**Figure 2.** QQplot for Aledo with (a) Gamma distribution and (b) our mixture.

and North/South gradients respectively. The North/South gradient (drier in the north and wetter in the south) that is also perceptible in pattern 4, is a classical recurrent feature of winter precipitation in Illinois.

### 3.2. Relating Regional Precipitation Patterns With Large-Scale NCEP Outputs

[16] At this stage, precipitation-related structures  $S_t$  have been derived (see Figure 3) and represent the regional scale. How to link them to the larger scale (the NCEP reanalysis) and how to connect them to the smaller scale (the weather stations) are the two remaining questions we have to address in this paper. In this section, we focus on answering the first one. To perform this task, we model the day-to-day probability transitions from the given weather state at day  $t$ , say  $S_t$ , to the state of the following day,  $S_{t+1}$  as a function of the current large atmospheric variables, say  $\mathbf{X}_t$ , from the NCEP reanalysis. More precisely, a nonhomogeneous Markov model [e.g., Bellone *et al.*, 2000] is fitted to our NCEP data and our states by applying the following temporal dependence structure,

$$P(S_t = s | S_{t-1} = s', \mathbf{X}_t) \propto \gamma_{s's} \cdot \exp \left[ -\frac{1}{2} (\mathbf{X}_t - \mu_{s's}) \Sigma^{-1} (\mathbf{X}_t - \mu_{s's})' \right], \quad (8)$$

where the symbol  $\propto$  means “proportional to” and where  $\gamma_{s's}$  is the baseline transition probability from pattern  $s'$  to pattern  $s$ , corresponding to the observed transition probability from  $s'$  to  $s$ , i.e., the proportion of transitions from  $s'$  to  $s$  over the total number of transitions. In the above formula, we can recognize a weight represented by the exponential term that is proportional to a normal density whose mean  $\mu_{s's}$  and variance matrix  $\Sigma$  are directly representing the influence of the large atmospheric variable  $\mathbf{X}_t$ . Equation (8) comes from Bayes's theorem, saying that

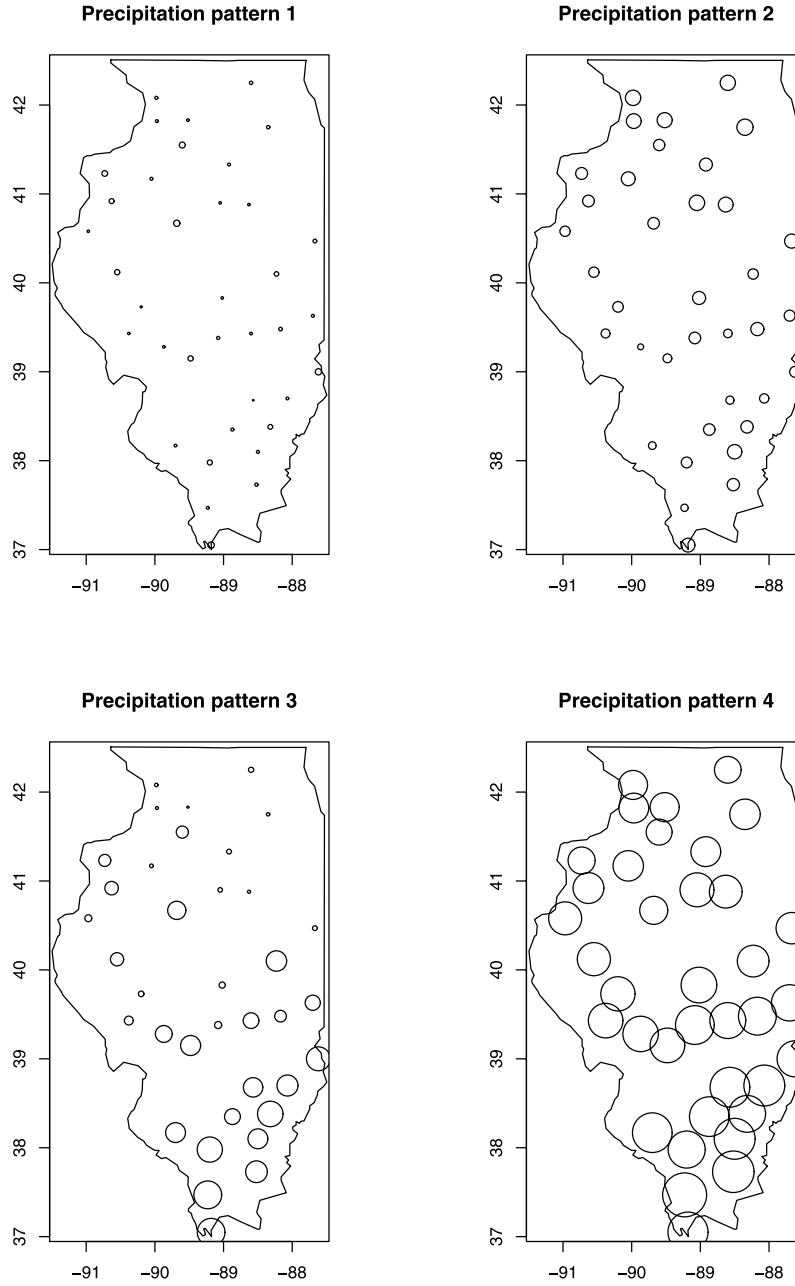
$$\begin{aligned} P(S_t = s | S_{t-1} = s', \mathbf{X}_t) &= \frac{P(S_t = s | S_{t-1} = s') P(\mathbf{X}_t | S_t = s, S_{t-1} = s')}{P(\mathbf{X}_t | S_{t-1} = s')} \\ &= \frac{\gamma_{s's} P(\mathbf{X}_t | S_t = s, S_{t-1} = s')}{\sum_k \gamma_{ks} P(\mathbf{X}_t | S_t = s, S_{t-1} = k)}. \end{aligned} \quad (9)$$

[17] By assuming in equation (8) that  $\mathbf{X}_t$  is multivariate normal, equation (8) is easily derived. In equation (8),  $\mu_{s's}$  corresponds to the mean vector of the atmospheric variables at time  $t$  when transitioning from  $S_{t-1} = s'$  to  $S_t = s$ . The four precipitation patterns defined in section 3.1 imply a reasonable number of 16 possible transition. Hence the 16  $\mu_{s's}$  and  $\gamma_{s's}$  to be computed can be estimated very fast. As for  $\Sigma$ , it is the variance-covariance matrix for the whole data set of large-scale atmospheric data (centered around their mean). Indeed, as in work by Charles *et al.* [1999], Bellone *et al.* [2000] or Vrac *et al.* [2007], for stability reasons, a single covariance matrix is preferred over one matrix per transition. In contrast to the exponential part of equation (8), the baseline transition probability  $\gamma_{s's}$  in (8) is time invariant and corresponds to the transition probabilities that one would have if large-scale features did not bring any information. This case corresponds to the homogeneous Markov model. Hence allowing a nonhomogeneity in our Markov modeling brings the necessary flexibility to mathematically integrate large-scale information at the intermediate level of the regional precipitation patterns.

### 3.3. Linking Regional Precipitation Patterns to Local Precipitation

[18] In order to implement an efficient downscaling precipitation scheme, we also need to model accurately the distributional properties of precipitation at the smallest scale, i.e., the ones recorded at rain gauges. We now assume that, given the current weather state  $s$ , all the rainfall intensities for station  $i$  follow the density  $h_{\beta_{si}}$  given in (5) with state- and site-specific parameters. This gives us the last ingredient to determine our main density defined by (2): the probability of observing local rainfall intensities at day  $t$ , say  $\mathbf{R}_t = (R_{t,1}, \dots, R_{t,N})$ , given the current weather state, say  $S_t = s$ , and large-scale atmospheric variables, say  $\mathbf{X}_t$ . To compute  $f_{\mathbf{R}_t | \mathbf{X}_t, S_t}$ , we follow Bellone *et al.* [2000] who considered that each rain gauge is spatially independent given the state  $S_t$ . Mathematically, this assumption translates into the following equality:

$$f_{\mathbf{R}_t | \mathbf{X}_t, S_t}(r_{t1}, \dots, r_{tN}) = \prod_{i=1}^N f_{R_{ti} | \mathbf{X}_t, S_t}(r_{ti}). \quad (10)$$



**Figure 3.** Four station-based precipitation patterns over Illinois derived by the *Vrac et al.* [2007] HAC method, with area proportional to mean rainfall for each cluster.

[19] To give an explicit form for the density  $f_{R_{ti}|X_t, S_t}$ , we take advantage of work by *Vrac et al.* [2007], who suggested the following form:

$$f_{R_{ti}|X_t, S_t}(r_{ti}) = [p(X_t; \alpha_{si})h_{\beta_{si}}(r_{ti})]^{\mathbb{1}_{\{r_{ti}>0\}}} \times [1 - p(X_t; \alpha_{si})]^{\mathbb{1}_{\{r_{ti}=0\}}}, \quad (11)$$

where  $h_{\beta_{si}}$  is given by (5),  $\mathbb{1}_{\{a\}} = 1$  if  $a$  is true and 0 if false, and  $p(X_t; \alpha_{si})$  represents the probability of rain occurrence for weather station  $i$  in state  $s$ . Equation (11) may look complex at first sight. Basically, it is composed of three elements (1) the indicator function  $\mathbb{1}_{\{r_{ti}=0\}}$  is necessary to take into account that the rain gauge  $i$  can record no precipitation during day  $t$ ; (2)  $1 - p(X_t; \alpha_{si})$  provides the

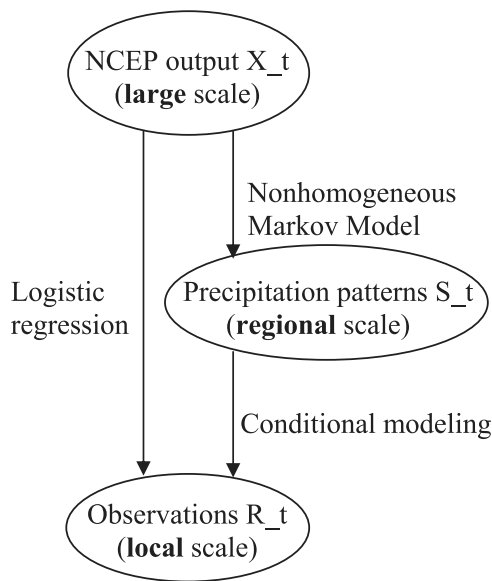
probability of such a dry day and it depends on the atmospheric variables  $X_t$  through a logistic regression with parameters  $\alpha_{si}$ , as suggested by *Jeffries and Pfeiffer* [2000],

$$p(X_t; \alpha_{si}) = P(R_{ti} > 0 | S_t = s, X_t) = \frac{\exp(\mathbf{X}_t' \alpha_{si})}{1 + \exp(\mathbf{X}_t' \alpha_{si})}; \quad (12)$$

and (3) the density  $h_{\beta_{si}}(r_{ti})$  corresponds to positive rainfall values.

[20] Combining equations (8), (5), (10) and (11) constitutes the main components of our stochastic weather typing approach. Through the variables  $\mathbf{R}_t$ ,  $S_t$  and  $\mathbf{X}_t$ , it integrates three scales (small, regional and large) summarized in Figure 4. In addition, the full spectrum of precipitation





**Figure 4.** Schematic graph explaining the main components of our downscaling scheme.

values (dry events, medium precipitation, heavy rainfall) is modeled.

#### 4. A Case Study: Precipitation in Illinois

[21] As previously mentioned, Figure 3 displays our four selected regional precipitation patterns over the region of Illinois. From these four patterns, the nonhomogeneous Markov model is parameterized, and the parameters of the conditional distributions of precipitation are estimated by Maximum Likelihood Estimation (MLE), given each observed (i.e., predefined) pattern. In the following simulation process, the precipitation patterns are stochastically simulated, for each  $t$ , according to the parameterized NMM, influenced by the large-scale atmospheric variables. In other words, in the simulation step, we do not use the patterns defined previously by HAC but we generate new ones according to  $X_t$  and our model. Conditionally on the four patterns, equations (10) and (11) offer a wide range of modeling possibilities. For example, one may wonder if it is better to have a unique GPD shape parameter  $\xi$  for all precipitation patterns and at all rain gauges or if a better statistical fit can be obtained by allowing this shape parameter to vary from station to station, while taking into account the risk of overparameterization. Before presenting the seven different models that we have tested and compared, we note that the parameter  $\tau$  in equation (6) cannot be null. For this reason, from the limit of equation (6) when  $\tau$  goes to 0, we extend equation (6) to

$$w_{m,0}(r) = \begin{cases} 0, & \text{if } r < m \\ 0.5, & \text{if } r = m \\ 1, & \text{if } r > m \end{cases} \quad (13)$$

for  $\tau = 0$ , whenever we do not wish to estimate  $\tau$  and we think that the transition from the Gamma to the GPD distribution is very fast in the mixture defined by (10). Our seven models are the following ones: model (0), Gamma and

GPD mixtures whose parameters vary with location and precipitation pattern; model (i), only Gamma distributions (no GPD in the model) whose parameters vary with location and precipitation pattern; model (ii), Gamma and GPD mixtures with one  $\xi$  parameter per pattern (i.e., given the weather pattern, the weather stations have the same  $\xi$ ); model (iii), the same as model (ii) with  $\tau$  set to be equal to 0; model (iv), Gamma and GPD mixtures with one common  $\xi$  for all stations and all patterns; model (v), same as model (iv) with  $\tau$  set to be equal to 0; and model (iii)\*, same as model (iii), one  $\xi$  parameter per pattern with  $\tau = 0$ , except that only Gamma distributions are used in pattern 1. Indeed, since this pattern corresponds to small or null intensities of rainfall, a modeling of the extreme events could have no sense here.

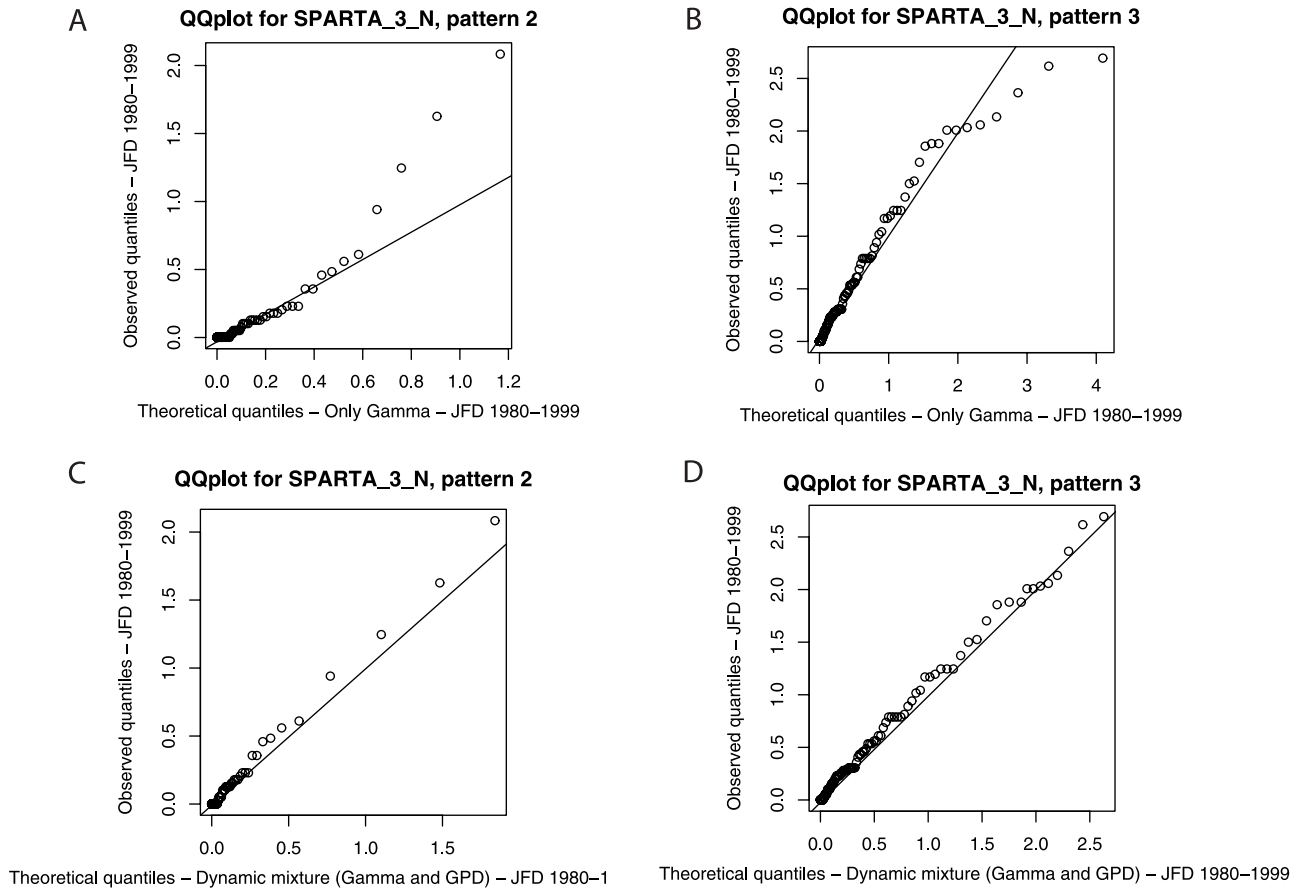
[22] From a statistical point of view, the GPD shape parameters are very difficult to estimate (wide confidence intervals). Hence diminishing the number of  $\xi$  parameters to estimate like in model (iii) reduces the overall variability. In addition, interpreting four  $\xi$  parameters (one per pattern, see models (ii) and (iii)) instead of  $37 \times 4$  is much easier for the hydrologist. Besides these two general guidelines, we need a more objective “measure” to compare our seven models. As in section 2, we opt for minimizing the classical AIC criterion (similar results are obtained with the BIC).

[23] Our seven models’ differences primarily focus on the degree of flexibility allowed for  $\xi$  and  $\tau$ . Concerning the other parameters ( $\sigma$ ,  $m$ , ...), we allow them to vary across stations and across patterns because they mainly represent local variability.

[24] For each model, we estimate its parameters by implementing a maximum likelihood estimation method. To illustrate the quality and drawbacks of our approach, we will comment on five example stations in this section: Aledo (northwest of Illinois), Aurora (northeast), Fairfield (southeast), Sparta (southwest), and Windsor (center-east of Illinois). This subset was picked because we believe that it represents a large range of cases and space limitations make it impossible to provide plots and tables for all 37 stations.

[25] Concerning the large-scale atmospheric variables  $X_t$ , we assume that only the NCEP grid cells over Illinois have the potential to influence local precipitation and transition probabilities. Consequently, we only work with the six grid cells that cover Illinois. According to the studied region, it is possible that taking more NCEP grid cells into account could improve the modeling and the simulation process. A few attempts have been made to enlarge the NCEP area influencing local precipitation and patterns transitions. The associated results, not presented here, did not show any clear improvement for the Illinois region, compared to the results obtained from the six grid cells. Moreover, the more grid cells we work on, the more parameters we have (with a risk of overparameterization). Hence from a computational point of view, it is better to restrict the large-scale influence to a reasonable number of NCEP grid cells over Illinois. On the basis of these two considerations, we then limit the application presented here to the six NCEP grid cells over Illinois to influence local precipitation and patterns transitions.

[26] Instead of working directly with the raw variables,  $Z_{850}$ ,  $Q_{850}$ , and  $\Delta T_{d850}$ , corresponding to  $6 \times 3 = 18$  variables, we perform a Singular Value Decomposition [Von Storch



**Figure 5.** QQplots of precipitation patterns 2 and 3 for station “Sparta”, for (a, b) function  $h_{\beta}$  in equation (11) as a Gamma distribution and (c, d)  $h_{\beta}$  as a mixture (equation (5)). Units are centimeters.

and Zwiers, 1999; Vrac et al., 2007; Wilks, 2006]. This has the advantage of reducing significantly the dimensionality of the NCEP data, while keeping the main part of information brought by the reanalysis. The SVD operation gives us the following summary: the SVD explains 93.6%, 98.6%, and 97.5% of the correlation for  $Z_{850}$ ,  $Q_{850}$ , and  $\Delta T_{d850}$  respectively.

[27] A central theme in this paper is how to capture the full range of precipitation, extremes included. To determine if the addition of a GPD to a Gamma density is worthwhile, Figure 5 displays QQplots (empirical quantiles versus modeled quantiles) for the Sparta station for two precipitation patterns (see the left and right panels) and in two models: (0) and (i), see Figure 5, top and bottom plots, respectively. In contrast to histograms, the QQplots are, by design, capable of representing the quality of the estimated fit at the end of the distribution tail; that is, they can show the capacity of our mixture model to represent extreme precipitation.

[28] Figure 5 indicates that a fitted Gamma has the tendency to either underestimate (Figure 5a) or overestimate (Figure 5b) the largest precipitation for this station, respectively to the precipitation patterns. Figures 5a and 5c show that, for pattern 2, our mixture can model heavier rainfall than the gamma distribution alone (i.e., characterizes stronger intensities for this pattern/station). To explain how the Gamma model can overestimate large precipitation in Figure 5b, we have to keep in mind that the whole rainfall

range is fitted and the Gamma distribution does not have a shape parameter for the tail of the distribution. In the presence of a heavy tail, it is not clear how the estimation procedure is going to compensate the facts that the gamma distribution is not heavy tailed and that the whole distribution has to be fitted. Either the Gamma-scale parameter can be largely overestimated (by the largest values) or underestimated (depending on the spread and the size of the sample). Applying a robust estimator to find the Gamma-scale parameter should remove the problem of overestimation, but then heavy tailed values will even be more disregarded. Consequently, a possible solution is to allow a distribution (like the GPD) with a shape parameter. More generally, Figure 5 clearly indicates that integrating a GPD improves the fit of “large” rainfalls for this station, as the closer the estimated quantiles are to the empirical quantiles the better. Of course, this does not mean that this is true for all stations and all patterns. Instead, this shows that our mixture defined by (5) provides the necessary modeling flexibility to describe heavy-tailed behaviors when needed. If no heavy rainfalls are observed at a given station, the estimated weight defined by (6) should take small values to favor the Gamma distribution, i.e.,  $m$  large for this station.

[29] Concerning the model selection, Table 2 compares models (0) and (i) with respect to the Akaike Information Criterion (AIC) for our five selected stations and for each precipitation pattern. Because the BIC values gave us equivalent results, they are not provided in this table,

**Table 2.** Akaike Information Criterion (AIC) Values Obtained Pattern by Pattern for Five Weather Stations<sup>a</sup>

| Station   | Model | Pattern 1      | Pattern 2      | Pattern 3      | Pattern 4     |
|-----------|-------|----------------|----------------|----------------|---------------|
| Aledo     | (0)   | <b>-351.15</b> | -486.98        | -162.21        | 86.72         |
|           | (i)   | -349.15        | <b>-493.05</b> | <b>-163.29</b> | <b>84.86</b>  |
| Aurora    | (0)   | -948.62        | -663.43        | -228.28        | 272.63        |
|           | (i)   | <b>-954.48</b> | <b>-670.41</b> | <b>-235.40</b> | <b>265.43</b> |
| Fairfield | (0)   | -367.72        | <b>-513.05</b> | <b>57.15</b>   | <b>499.93</b> |
|           | (i)   | <b>-375.99</b> | -282.21        | 97.42          | 741.63        |
| Sparta    | (0)   | <b>-131.34</b> | <b>-488.52</b> | <b>-128.03</b> | <b>613.23</b> |
|           | (i)   | -129.61        | -466.06        | -123.57        | 766.54        |
| Windsor   | (0)   | <b>-632.25</b> | -982.22        | -321.01        | <b>441.08</b> |
|           | (i)   | -613.92        | <b>-985.26</b> | <b>-325.78</b> | 579.47        |

<sup>a</sup>The bold values correspond to the optimal criteria for either model (0) or (i).

illustrating that the optimal choice between model (0) and model (i) varies greatly across stations and across patterns. For example, introducing a GPD seems to be a good choice for Sparta, while a simpler Gamma model appears to be sufficient for Aurora.

[30] Table 3 contain the AIC values obtained for the seven models. The bold values correspond to the optimal criterion of each row. Taking model (iii)\* ( $\tau = 0$ , a Gamma distribution for pattern 1 and one  $\xi$  parameter per pattern for patterns 2–4) provides the best AIC for Sparta, while setting one overall  $\xi$  parameter gives the best AIC for the four other stations. For any of the five stations, we can remark that setting  $\tau = 0$  in model (ii), i.e., going from model (ii) to model (iii), brings an improvement of the AIC. This means that restricting the number of  $\xi$  parameters generally provides better criteria. Models (iii)\* and (iv) seem to be the most competitive ones in general (i.e., for most of the stations separately), while the preferred model tends to be (iii)\* for the set of the five selected weather stations altogether (last row of Table 3). Consequently, model (iii)\*, i.e., pattern 1 associated with Gamma distributions and patterns 2–4 to mixtures with one  $\xi$  parameter per pattern with the constant  $\tau = 0$ , is chosen as the most efficient model, as it provides the best overall criterion for the set of these five stations. Hence this model can well represent both common and extreme precipitation values with an acceptable number of parameters and has the overall preference.

[31] Table 4 shows the values of the  $\xi$  parameters and the values of the  $m$  parameters (when applicable) for the five example stations for model (iii)\*. The three  $\xi$  parameters are clearly positive. These positive values indicate that the heavy tail component in our mixture pdf is essential to model heavy rainfalls for precipitation patterns 2–4, while the Gamma distributions (with light tails) are sufficient in pattern 1 corresponding to small precipitation events. Unsurprisingly, the  $m$  parameters tend to increase from pattern 2 (with the smallest rainfall intensities among patterns 2–4) to pattern 4 (with the strongest rainfalls among all patterns).

[32] To visually evaluate the fit between our model (iii)\* and the observed precipitation, a QQplot is plotted for the Aledo station in Figure 6. The agreement between observed

**Table 3.** Akaike Information Criterion (AIC) Values Obtained for our Five Selected Weather Stations and for Our Seven Models<sup>a</sup>

| Station           | Model (0) $p = 24n$ | Model (i) $p = 8n$ | Model (ii) $p = 20n + 4$ | Model (iii) $p = 16n + 4$ | Model (iv) $p = 20n + 1$ | Model (v) $p = 16n + 1$ | Model (iii)* $p = 12n + 5$ |
|-------------------|---------------------|--------------------|--------------------------|---------------------------|--------------------------|-------------------------|----------------------------|
| Aledo             | AIC = -796.52       | AIC = -816.58      | AIC = -795.76            | AIC = -809.79             | AIC = -819.46            | AIC = -816.18           | AIC = -816.79              |
| Aurora            | AIC = -1137.47      | AIC = -1149.99     | AIC = -1256.53           | AIC = -1293.89            | AIC = -1358.48           | AIC = -1152.51          | AIC = -1299.89             |
| Fairfield         | AIC = 14.36         | AIC = 103.07       | AIC = 22.45              | AIC = 22.37               | AIC = -76.81             | AIC = -10.21            | AIC = 16.37                |
| Sparta            | AIC = 277.10        | AIC = 372.92       | AIC = 235.65             | AIC = 228.35              | AIC = 231.91             | AIC = 251.44            | AIC = <b>222.35</b>        |
| Windsor           | AIC = -1014.80      | AIC = -920.68      | AIC = -1016.25           | AIC = -1017.59            | AIC = -1069.99           | AIC = -1028.91          | AIC = -1023.59             |
| All five stations | AIC = -4433.18      | AIC = -4422.27     | AIC = -4479.50           | AIC = -4515.13            | AIC = -4425.06           | AIC = -4423.78          | AIC = <b>-4553.13</b>      |

<sup>a</sup>The bold values correspond to the optimal criterion per row. Below each model's name, the number  $p$  of parameters for  $n$  stations is provided.

**Table 4.** Values of the  $\xi$  and  $m$  Parameters for the Five Example Stations for Model (iii)\*<sup>a</sup>

|                   | Pattern 1 | Pattern 2 | Pattern 3 | Pattern 4 |
|-------------------|-----------|-----------|-----------|-----------|
| $\xi$             | NA        | 0.3       | 0.13      | 0.26      |
| $m$ for Aledo     | NA        | 0.73      | 0.81      | 1.06      |
| $m$ for Aurora    | NA        | 0.28      | 0.48      | 1.38      |
| $m$ for Fairfield | NA        | 1.61      | 1.24      | 1.84      |
| $m$ for Sparta    | NA        | 0.46      | 1.01      | 1.83      |
| $m$ for Windsor   | NA        | 0.56      | 0.81      | 0.96      |

<sup>a</sup>Nonapplicable (NA) is indicated for pattern 1, since this pattern is associated with Gamma distributions in this model.

and theoretical quantiles (even for high quantiles) is clearly good. Figure 6 has to be compared to Figure 2. This allows us to conclude that, not only the AIC is better for model (iii)\* than for a “no pattern” modeling, but also that model (iii)\* improves the QQplot.

[33] Besides heavy rainfalls, an important characteristic of precipitation modeling is the representation of the so-called wet and dry spell periods, fundamental quantities in agriculture. Note that none of the following results concerning wet and dry spells and local precipitation probabilities, presented and shown from Figure 7, depends on the Gamma or mixture models. Indeed, they are only related to the nonhomogeneity introduced in the Markov model (8), that characterizes pattern transitions, and to the probabilities of local rain occurrence modeled as logistic regressions (see equations (11) and (12)). So, the following results are directly derived from the model developed by *Vrac et al.* [2007] and allow us to compare some precipitation appearance characteristics obtained from the “four precipitation patterns” and those obtained from the alternative “no pattern” approach.

[34] In this context, we have noticed that the four precipitation patterns have to be included in order to obtain adequate wet and dry spell probabilities. For example, Figure 7 shows such probabilities (in log-scale) at two stations, respectively Fairfield and Windsor. Figures 7a and 7b display these probabilities when the four precipitation patterns are included in our analysis. In contrast,

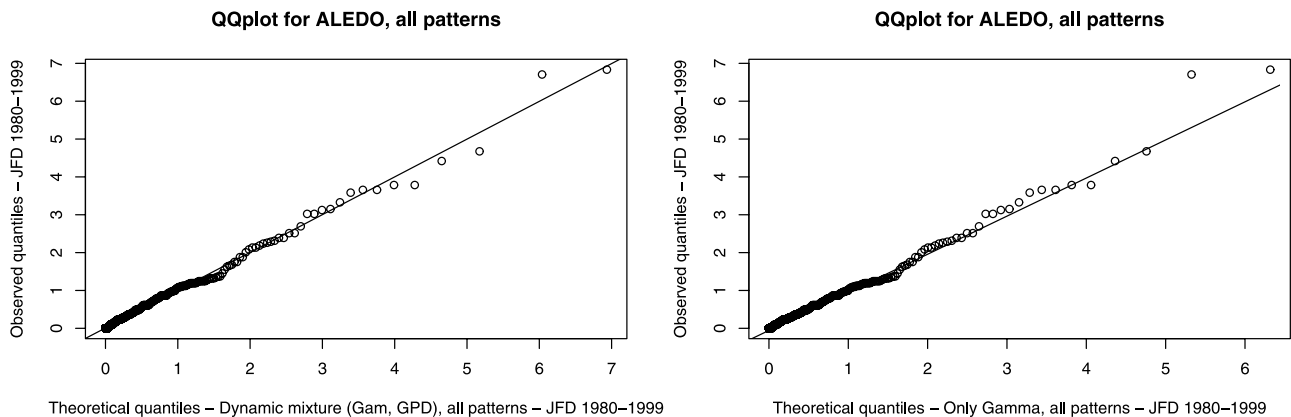
Figures 7c and 7d show the results when no patterns are introduced. From these graphs, one can see that the “no pattern” option is not completely satisfying, it tends to underestimate the probabilities for long spells, above all for dry spells.

## 5. Conclusion

[35] We presented here a nonhomogeneous stochastic weather typing method to downscale the full spectrum of precipitation distributional behaviors. Our downscaling technique is based on a nonhomogeneous Markov model that characterizes the transitions amongst different precipitation patterns obtained from a hierarchical ascending clustering algorithm. Conditionally on these precipitation patterns, the precipitation distribution is modeled by a mixture model that integrates heavy rainfalls, medium precipitation and no rain occurrences, and that depends on large-scale features given from a SVD applied to NCEP reanalysis.

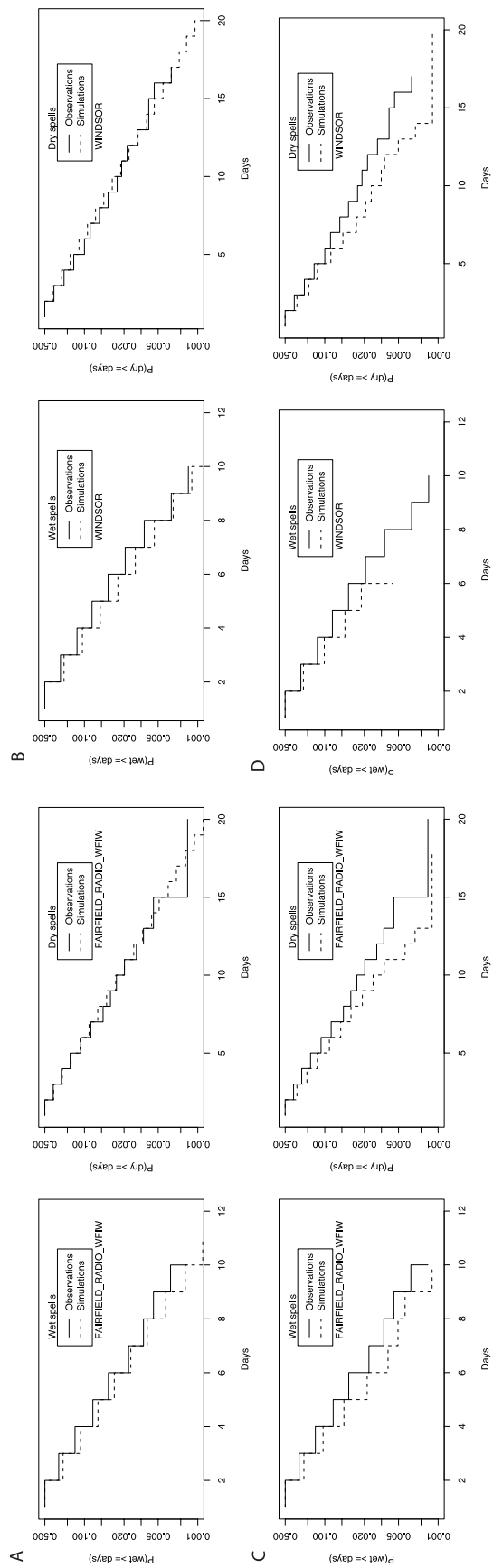
[36] After applying our approach to the region of Illinois, it appears that a specific subclass of our model (the one with Gamma distributions for pattern 1 and mixture models with a single GPD shape parameter per pattern for patterns 2–4) produces the best fit with respect to the AIC criterion for this region. In terms of extreme precipitation, this model corresponds to a very fast transition from the Gamma distribution to the GPD for patterns 2–4. It is also worthwhile to highlight that introducing four precipitation patterns produces better precipitation characteristics than a direct “no pattern” approach does.

[37] As possible improvements, spatial dependence modeling could be introduced in this model to better represent the correlation between stations. In that context, Bayesian hierarchical methods could provide an additional flexibility. A possible application of our downscaling procedure could be the projection of future local precipitation based on large-scale climate change simulated by GCMs. While the estimation step requires both present large- and local-scale data, the local projection of future climate scenarios can be done by using only the GCM outputs describing future time periods. On the basis of the NMM previously fitted, the future large-scale outputs are first used to influence the simulation of future precipitation patterns through equation (8). No local



**Figure 6.** QQplot for Aledo with four patterns and model (iii)\*, i.e., Gamma distributions for pattern 1 and mixtures for patterns 2–4 with one  $\xi$  per pattern and  $\tau = 0$ .





**Figure 7.** Wet and dry spells probabilities (in log-scale) obtained for Fairfield and Windsor. (a, b) The “4 patterns” approach. (c, d) The “no pattern” approach.

precipitation is needed for this step, since it is obviously not even available. Conditionally on the generated future patterns, probabilities of local rainfall events can be computed, influenced by the large-scale GCM outputs, through equation (12) for rain appearances and through equation (11) for intensities. These local projections would then allow economic impact studies of extreme precipitation.

[38] **Acknowledgments.** Although this research has been funded in part by the United States Environmental Protection Agency through STAR Cooperative Agreement R-82940201 to the University of Chicago, it has not been subjected to the Agency's required peer and policy review and therefore does not necessarily reflect the views of the Agency, and no official endorsement should be inferred. P. Naveau's research work is supported by the European E2-C2 grant, the National Science Foundation (grant NSF-GMC (ATM-0327936)) and by The Weather and Climate Impact Assessment Science Initiative at the National Center for Atmospheric Research (NCAR).

## References

- Akaike, H. (1974), A new look at the statistical model identification, *IEEE Trans. Autom. Control*, *19*, 716–723.
- Bárdossy, A., H. Muster, L. Duckstein, and I. Bogardi (1994), Knowledge-based classification of circulation patterns for stochastic precipitation modelling, in *Time Series Analysis in Hydrology and Environmental Engineering*, vol. 3, edited by K. W. Hipel et al., pp. 19–32, Kluwer Acad., Dordrecht, Netherlands.
- Barnett, T., and R. Preisendorfer (1978), Multifield analog prediction of short-term climate fluctuations using a climate state vector, *J. Atmos. Sci.*, *35*, 1771–1787.
- Bellone, E., J. P. Hughes, and J. P. Guttorp (2000), A hidden markov model for downscaling synoptic atmospheric patterns to precipitation amounts, *Clim. Res.*, *15*, 1–12.
- Biau, G., E. Zorita, H. von Storch, and H. Wackernagel (1999), Estimation of precipitation by kriging in the EOF space of the sea level pressure field, *J. Clim.*, *12*, 1070–1085.
- Breiman, L., J. Friedman, R. Olshen, and C. Stone (1984), *Classification And Regression Trees (CART)*, Chapman and Hall, New York.
- Bunkers, M. J., J. R. Miller, and A. T. DeGaetand (1996), Definition of climate regions in the northern plains using an objective cluster modification technique, *J. Clim.*, *9*, 130–146.
- Cannon, A. J., and P. H. Whitfield (2002), Downscaling recent streamflow conditions in British Columbia, Canada using ensemble neural network models, *J. Hydrol.*, *259*, 136–151.
- Charles, S. P., B. C. Bates, and J. P. Hughes (1999), A spatio-temporal model for downscaling precipitation occurrence and amounts, *J. Geophys. Res.*, *104*, 31,657–31,669.
- Coles, S. (2001), *An Introduction to Statistical Modeling of Extreme Values*, Springer, London.
- Davis, R. E., R. Dolan, and G. Demme (1993), Synoptic climatology of Atlantic coast northeasters, *Int. J. Climatol.*, *13*, 171–189.
- Fisher, R. A., and L. H. C. Tippett (1928), Limiting forms of the frequency distribution of the largest or smallest member of a sample, *Proc. Cambridge Philos. Soc.*, *24*, 180–190.
- Frigessi, A., O. Haug, and H. Rue (2003), A dynamic mixture model for unsupervised tail estimation without threshold selection, *Extremes*, *5*, 219–235.
- Huth, R. (2001), Disaggregating climatic trends by classification of circulation patterns, *Int. J. Climatol.*, *21*, 135–153.
- Huth, R. (2002), Statistical downscaling of daily temperature in central Europe, *J. Clim.*, *15*, 1731–1742.
- Jeffries, N., and R. Pfeiffer (2000), A mixture model for the probability distribution of rain rate, *Environmetrics*, *12*, 1–10.
- Katz, R. W. (1977), Precipitation as a chain-dependent process, *J. Appl. Meteorol.*, *16*, 671–676.
- Katz, R., M. B. Parlange, and P. Naveau (2002), Statistics of extremes in hydrology, *Adv. Water Resour.*, *25*, 1287–1304.
- Mamassis, N., and D. Koutsoyiannis (1996), Influence of atmospheric circulation types on space-time distribution of intense rainfall, *J. Geophys. Res.*, *101*, 26,267–26,276.
- Naveau, P., M. Nogaj, C. Ammann, P. Yiou, D. Cooley, and V. Jomelli (2005), Statistical methods for the analysis of climate extremes, *C. R. Geosci.*, *337*, 1013–1022.
- Pongracz, R., J. Bartholy, and I. Bogardi (2001), Fuzzy rule-based prediction of monthly precipitation, *Phys. Chem. Earth*, *9*, 663–667.
- Schnur, R., and D. Lettenmaier (1998), A case study of statistical downscaling in Australia using weather classification by recursive partitioning, *J. Hydrol.*, *212–213*, 362–379.
- Schwarz, G. (1978), Estimating the dimension of a model, *Ann. Stat.*, *6*, 461–464.
- Semenov, M. A., and E. M. Barrow (1997), Use of a stochastic weather generator in the development of climate change scenarios, *Clim. Res.*, *35*, 397–414.
- Semenov, M. A., R. J. Brooks, E. M. Barrow, and C. W. Richardson (1998), Comparison of the WGEN and the LARS-WG stochastic weather generators in diverse climates, *Clim. Res.*, *10*, 95–107.
- Snell, S. E., S. Gopal, and R. K. Kaufmann (2000), Spatial interpolation of surface air temperatures using artificial neural networks: Evaluating their use for downscaling GCMs, *J. Clim.*, *13*, 886–895.
- Von Storch, H., and F. W. Zwiers (1999), *Statistical Analysis in Climate Research*, Cambridge Univ. Press, Cambridge, U.K.
- Vrac, M., A. Chédin, and E. Diday (2005), Clustering a global field of atmospheric profiles by mixture decomposition of copulas, *J. Atmos. Oceanic Technol.*, *22*, 1445–1459.
- Vrac, M., M. Stein, and K. Hayhoe (2007), Statistical downscaling of precipitation through a nonhomogeneous stochastic weather typing approach, *Clim. Res.*, in press.
- Ward, J. H. (1963), Hierarchical grouping to optimize an objective function, *J. Am. Stat. Assoc.*, *58*, 236–244.
- Wigley, T. M., P. Jones, K. Briffa, and G. Smith (1990), Obtaining subgrid scale information from coarse resolution general circulation output, *J. Geophys. Res.*, *95*, 1943–1953.
- Wilks, D. S. (1999), Multisite downscaling of daily precipitation with a stochastic weather generator, *Clim. Res.*, *11*, 125–136.
- Wilks, D. (2006), *Statistical Methods in the Atmospheric Sciences*, 2nd ed., Elsevier, Oxford, U.K.
- Wilks, D. S., and R. L. Wilby (1999), The weather generation game: A review of stochastic weather models, *Prog. Phys. Geogr.*, *23*, 329–357.
- Wilson, P. S., and R. Toumi (2005), A fundamental probability distribution for heavy rainfall, *Geophys. Res. Lett.*, *32*, L14812, doi:10.1029/2005GL022465.
- Yiou, P., and N. Nogaj (2004), Extreme climatic events and weather regimes over the North Atlantic: When and where?, *Geophys. Res. Lett.*, *31*, L07202, doi:10.1029/2003GL019119.
- Zorita, E., and H. von Storch (1998), The analog method as a simple statistical downscaling technique: Comparison with more complicated methods, *J. Clim.*, *12*, 2474–2489.

---

P. Naveau and M. Vrac, Laboratoire des Sciences du Climat et de l'Environnement, IPSL-CNRS, CEA Saclay, F-91191 Gif-sur-Yvette, France. (mathieu.vrac@cea.fr)