

**On the role of serial correlation and field significance in detecting  
changes in extreme precipitation frequency**

Stefano Farris<sup>1</sup>, Roberto Deidda<sup>1</sup>, Francesco Viola<sup>1</sup>, and Giuseppe Mascaro<sup>2</sup>

(1) Dipartimento di Ingegneria Civile, Ambientale e Architettura, Università di Cagliari, Italy

(2) School of Sustainable Engineering and the Built Environment, Arizona State University, USA

---

*Corresponding author address:* Giuseppe Mascaro, School of Sustainable Engineering and the Built Environment, Arizona State University, 126b, 21 E 6th St, Tempe, AZ 85281. *E-mail:* gmascaro@asu.edu

## 27    **Abstract**

28            Statistical trend analyses of observed precipitation (P) time series are key to validate theoretical  
29 arguments and climate projections suggesting that extreme P will increase in a warmer climate. Recent  
30 work warned about possible misinterpretation of trend tests if the presence of serial correlation and field  
31 significance are not considered. Here, we investigate these two aspects focusing on extreme P frequencies  
32 derived from 100-year daily records of 1087 worldwide gauges of the Global Historical Climate  
33 Network. For this aim, we perform Monte Carlo experiments based on count time series generated with  
34 the Poisson integer autoregressive model and characterized by different sample size, level of  
35 autocorrelation, and trend magnitude. The main results are as follows. (1) Empirical autocorrelations are  
36 consistent with those of uncorrelated and stationary or nonstationary count time series, while empirical  
37 trends cannot be explained as the exclusive effect of autocorrelation; incorporating the impact of serial  
38 correlation in trend tests on extreme P frequency has then limited impacts on tests' performance. (2)  
39 Accounting for field significance improves interpretation of test results by limiting type-I errors, but it  
40 also decreases test power; results of local tests could complement field significance outcomes and help  
41 identify weak trend signals where several trends of coherent sign are detected. (3) Based on these  
42 findings, evident patterns of statistically significant increasing (decreasing) trends emerge in central and  
43 eastern North America, northern Eurasia, and central Australia (southwestern America, southern Europe,  
44 and southern Australia). The methodological insights of this work support trend analyses of any  
45 hydroclimatic variable.

46

## 1 Introduction

Extreme precipitation (P) is one of the natural hazards with the most significant socioeconomic impacts. Heavy P is the primary input of floods and flash floods, which cause annually large damages to properties and high numbers of fatalities worldwide (Ashley & Ashley, 2008; Peden et al., 2017). For example, the National Oceanic and Atmospheric Administration (NOAA) estimated that, in United States, flooding and severe storms resulted in \$437 billion damages and 2379 fatalities from 1980 to 2020 (Smith, 2021). In urban regions, intense P storms lead to pluvial flooding with impacts on traffic (Hooper et al., 2014; Bucar & Hayeri, 2020) and occurrence of power outages (Boggess et al., 2014). Extreme P events have also significant consequences on public health by degrading water quality (Gershunov et al., 2018) and increasing outbreaks of waterborne diseases (Cann et al., 2013). Studies have also shown that extreme P events may reduce crop production (Rosenzweig et al., 2004; Li et al., 2019).

Theoretical arguments suggest that the intensity of P extremes is expected to increase in a future warmer climate (Trenberth et al., 2003; Emori & Brown, 2005; Trenberth, 2011; Nie et al., 2018). According to the Clausius-Clapeyron (CC) equation, as surface temperature rises, the atmospheric water-holding capacity should grow at a rate of  $7\% \text{ K}^{-1}$ . Extreme P is held to increase at a rate close to the CC value or even higher if the strength of moisture convergence will rise (Trenberth et al., 2003). Driven by these theoretical arguments and the evidence of increasing global surface temperature over the last five decades (Hansen et al., 2010; Papalexiou et al., 2020), a number of empirical studies have started to investigate temporal changes of magnitude and frequency in observed records of P extremes based on the application of statistical trend tests. Table 1 summarizes some of these efforts conducted at global and regional scales using mainly daily records of rain gages. Conclusions that emerge across all studies are that (i) trends are mainly increasing but statistically significant only at a limited number of sites; (ii) statistically significant trends are more evident in frequency rather than magnitude of extreme P; (iii)

71 increasing trends are mainly located in eastern and Midwestern U.S. and some regions of Eurasia; and  
72 (iv) decreasing trends occur in western U.S. and southern Australia. Despite these common qualitative  
73 outcomes, Table 1 emphasizes how these studies vary widely in terms of duration of the investigated  
74 time period (ranging from 30 to 112 years); spatial aggregation of the information provided by the rain  
75 gages (from point to subcontinental regions); and metrics used to characterize extreme P (targeting  
76 magnitude or frequencies above a threshold). As a result, it is difficult to quantitatively compare their  
77 results, a task that would be highly needed for practical applications including the update of engineering  
78 design standards (Wright et al., 2019).

79 A key step to improve empirical trend studies of extreme P, facilitate their comparison, and  
80 corroborate physical hypotheses on future changes in the driving climate dynamics is to critically assess  
81 power and interpret results of statistical trend tests under the possible conditioning of serial correlation,  
82 if any, and when applied at multiple sites. We argue that these tasks have received limited attention,  
83 likely because these tests are easy to apply numerically via widespread software. These issues have been  
84 also recently highlighted by Serinaldi et al. (2018), who discussed potential causes of misuse and  
85 misinterpretation of statistical trend tests. One of these causes is the presence of autocorrelation in the  
86 analyzed time series, which may occur in hydrologic records as a result of long-term natural climate  
87 variability (Koutsoyiannis, 2011; Sun et al., 2018). Several statistical trend tests evaluate the null  
88 hypothesis  $H_0$  of random ordering in the time series (note that  $H_0$  is more often defined as “the time series  
89 is stationary” or “no trend is present in the time series”). When the time series is autocorrelated while  
90 still being stationary, the ordering is not random and the application of trend tests could result in rejecting  
91  $H_0$  more frequently than expected by the significance level (i.e., the type-I error increases). This problem  
92 has been investigated for time series of real numbers (e.g., P magnitudes), focusing largely on the Mann-  
93 Kendall test (von Storch, 1999; Yue et al., 2002; Hamed, 2009, among others). For this test, the presence  
94 of autocorrelation leads to an increase of the test statistic variance, a phenomenon known as variance

95 inflation. To address this issue, two main methods have been proposed including: (i) applying trend tests  
96 accounting for a proper estimation of the inflated variance (Hamed & Ramachandra Rao, 1998), and (ii)  
97 “prewhitening” the time series, i.e., removing the autocorrelation (Katz, 1988; von Storch, 1999). For  
98 both methods, a serial correlation structure of the process has to be adopted based on, e.g., autoregressive  
99 or fractional Gaussian models (Hamed, 2009).

100 As shown in Table 1, most studies that investigated trends in extreme P have not considered the  
101 presence of autocorrelation at all or found it to be negligible by simply verifying that the lag-1  
102 autocorrelation,  $\rho$ , averaged across all records is close to zero (Groisman et al., 2005; Westra et al., 2013;  
103 Papalexiou & Montanari, 2019). Only a small number of efforts have applied techniques to estimate the  
104 inflated variance (Tramblay et al., 2013; Kunkel & Frankson, 2015) or prewhitening procedures  
105 (Alexander et al., 2006). Unfortunately, several papers have showed that these methods are not easy to  
106 apply, because the interaction between possible trends and autocorrelation leads to biases in the  
107 estimation of their parameters, which could in turn decrease the trend test power (Yue & Wang, 2002;  
108 Bayazit & Önöz, 2007). Moreover, Serinaldi et al. (2018) have demonstrated that the application of  
109 different prewhitening techniques to the same dataset could produce markedly diverse outcomes. We  
110 have also found that, in the literature that investigated the effect of serial correlation on trend tests,  
111 analyses have mainly relied on synthetic experiments in controlled conditions, while observed datasets  
112 have been used only in a limited number of cases. In particular, to our knowledge, no study has  
113 thoroughly investigated this problem focusing on observed extreme P frequencies.

114 Another aspect that deserves careful consideration when conducting statistical trend analyses of  
115 extreme P is test multiplicity or field significance (Livezey & Chen, 1983; Katz & Brown, 1991; Wilks,  
116 1997; Daniel et al., 2012; Serinaldi et al., 2018). This accounts for the fact that, when a test is applied  
117 collectively at  $M$  locations (e.g., rain gages or grid points) with a significance level  $\alpha$ , the null hypothesis  
118 may be rejected, on average, at  $\alpha \cdot M$  sites while holding true for the entire set of locations. If the test

119 outcomes are interpreted locally, one can erroneously conclude that a statistically significant trend exists  
120 at the  $\alpha \cdot M$  sites. This could be even more likely when P records are spatially correlated: in such a case,  
121 local tests are not independent and it may be possible to find spatial clusters where  $H_0$  has been  
122 erroneously rejected that could mistakenly be considered as physically meaningful spatial features.  
123 Results of multiple tests should be instead interpreted globally. To this end, two types of methods have  
124 been proposed, including (i) techniques based on counting the number of  $H_0$  rejections and comparing  
125 them with thresholds derived from the Binomial distribution (Livezey & Chen, 1983) or from  
126 bootstrapping methods (Khaliq et al., 2009; Wilks, 2019), and (ii) methods that minimize the false  
127 discovery rate or FDR (Benjamini & Hochberg, 1995; Wilks, 2006, 2016). Modifications of these  
128 methods have been proposed to account for spatial dependence. The great majority of previous studies  
129 of trend in extreme P have not accounted for field significance, with the exception of Alexander et al.,  
130 (2006) and Westra et al., (2013), who used bootstrapping methods, and Trambly et al. (2013), who  
131 applied a test based on FDR (Table 1). Additional work is then needed to better investigate the importance  
132 of field significance in trend analyses of extreme P records and how its quantification affects power of  
133 statistical trend tests.

134 Driven by these research needs, this study investigates the effect of serial correlation and field  
135 significance on power, errors, and interpretation of trend tests applied to observed records of extreme P  
136 frequencies at multiple sites. We focus on frequencies (i.e., count time series of exceedances above a  
137 threshold) because changes in extreme P have been more effectively detected on counts rather than  
138 magnitudes (Papalexiou & Montanari, 2019; Wright et al., 2019). For our analyses, we use 100-year  
139 daily P records from 1087 gages the Global Historical Climate Network (GHCN)-Daily dataset (Menne  
140 et al., 2012) covering North America, northern and part of southern Europe, northern Asia, and Australia.  
141 The core of our methodological framework is based on Monte Carlo simulations, where stationary and  
142 nonstationary count time series with different levels of autocorrelation and trend magnitude are generated

143 using the Poisson integer autoregressive (INAR) model of order 1 or Poisson-INAR(1). INAR models  
144 were introduced to transfer the structure of autoregressive models for the simulation of integer-valued  
145 time series (e.g., McKenzie, 1985; Al-Osh & Alzaid, 1987; Weiß, 2008; Pedeli et al., 2015) and have  
146 been rarely applied in hydrology. After showing that the Poisson-INAR(1) model adequately reproduces  
147 the autocorrelation structure of most observed count time series, we apply a set of statistical analyses  
148 based on Monte Carlo simulations to gain insights on the impact of serial correlation on trend detection  
149 in the observed records. We then perform additional Monte Carlo experiments to quantify power and  
150 errors of several popular tests (Table 1) conducted locally and at multiple sites, utilizing the FDR test of  
151 Wilks (2006) to account for field significance. Finally, we use the knowledge gained with the analyses  
152 on serial correlation and field significance to apply trend tests to the observed extreme P frequencies and  
153 interpret their results in the studied regions. We repeat the analyses for different sample sizes, ranging  
154 from 30 to 100 years, and thresholds used to define the frequencies. While focused on extreme P, this  
155 work provides methodological insights supporting trend analyses of any hydroclimatic variable.

## 156 **2 Data**

157 We use daily P records from the GHCN dataset, which includes more than 100,000 stations in  
158 180 countries with record lengths ranging from a few years to more than 175 years and has been  
159 previously used in global (Kunkel & Frankson, 2015; Wilks, 2016; Papalexiou & Montanari, 2019) and  
160 regional (Wright et al., 2019; Kunkel et al., 2020) trend analyses. Here, after retaining only records  
161 passing all quality controls (Durre et al., 2010), for each station we label as “complete years” those with  
162 no more than 10% missing daily data and mark as missing all records collected in those years not  
163 satisfying this constraint. Then, we select  $M = 1087$  stations with at least 95 complete years in a common  
164 100-year period from 1916 to 2015. Fig. 1 shows the selected gages that are located in three main regions,  
165 including North America; northern and part of southern Europe; northern Asia; and Australia. For each  
166 record, we derive the count time series of extreme P frequencies  $\{o_t\}$  ( $t = 1, \dots, n$ , with  $n$  being the

number of years), defined as the annual occurrences of daily precipitation exceeding the  $q$ -th quantiles of its empirical cumulative distribution function (including zeros). These count time series are derived for the nonexceedance probabilities  $q = 0.9, 0.925, 0.95, 0.975$  for  $n = 100$  years and the most recent  $n = 30$  and 50 years.

### 3 Methodology

The methodology is described in four subsections. In section 3.1, we briefly illustrate the false FDR test that will be applied to evaluate the field significance in selected statistical tests for trend detection. In section 3.2, we investigate the parent distribution of the observed  $\{o_t\}$  count time series. In section 3.3 we explain the methods used to generate synthetic count time series simulating statistical properties and potential trends of the observed  $\{o_t\}$ . In section 3.4, we describe how Monte Carlo simulations based on these synthetic series are used to apply statistical trend tests under different null hypotheses, including possible presence of trend and autocorrelation.

#### 3.1. Evaluation of field significance

As discussed in the Introduction, results of tests conducted at multiple sites are affected by the problem known as test multiplicity or field significance. To account for this, the global null hypothesis  $H_0$  assuming that  $H_0$  is true at all locations should be investigated with a significance level  $\alpha_{\text{global}}$ . Here, we evaluate the field significance using the FDR test as described in Wilks (2006), since it has been proved more powerful than alternative field significance tests while being computationally efficient (Wilks, 2016). Its application is straightforward; given the  $p$ -values from any local test conducted at  $M$  sites, the FDR test rejects the local null hypothesis in those sites where the corresponding  $p$ -value is lower than a threshold  $p_{FDR}^*$  calculated as:

$$p_{FDR}^* = \max_{i=1, \dots, M} \left[ p_{(i)} : p_{(i)} \leq \left( \frac{i}{M} \right) \cdot \alpha_{FDR} \right] \quad (1)$$



188 where  $p_{(i)}$  is the  $i$ -th value in the sorted sample of the  $M$   $p$ -values, and  $\alpha_{FDR}$  is the significance level of  
189 the local test (see Wilks, 2016 for details). If the  $p$ -value is lower than  $p_{FDR}^*$  at one or more sites, then the  
190 global  $H_0$  is rejected at a level  $\alpha_{global} = \alpha_{FDR}$ . In these sites, the local  $H_0$  is also rejected and the potential  
191 existence of spatial patterns where  $H_0$  is rejected can be explored. A very attractive property of the FDR  
192 test is that it can be easily adapted to the cases of spatial dependence among the gage records. Using  
193 numerical simulations, Wilks (2016) suggests that the FDR test is robust to the presence of spatial  
194 correlation if a value of  $\alpha_{FDR} = 2\alpha_{global}$  is adopted. Unless stated otherwise, for all tests conducted in this  
195 study, we assume  $\alpha_{global} = 0.05$  and  $\alpha_{FDR} = 0.10$ .

### 196 **3.2. Preliminary inference on the parent distribution of exceedance counts**

197 We conduct preliminary analyses to identify a reasonable parent distribution for the observed  
198 exceedance counts  $\{o_t\}$  at the GHCN gages. Specifically, we apply the Chi-Square and Lilliefors (a  
199 generalization of Kolmogorov-Smirnov) goodness-of-fit (GOF) tests to evaluate the null hypothesis  $H_0$   
200 that the Poisson distribution well reproduces the marginal distribution of the observed counts. We do this  
201 for the count series with  $n = 30, 50$ , and  $100$  years. Instead of applying the GOF tests in their traditional  
202 formulation, we build the null distribution of the GOF test statistics through Monte Carlo simulations  
203 (details are provided in Section 3.4), because (i) statistical tables for the Chi-Square null distribution are  
204 usually derived and valid when parameters of the fitted distribution are estimated by minimizing the Chi-  
205 Square statistic (Fisher, 1922); and (ii) performances of GOF tests can be biased when applied to discrete  
206 variables (see e.g. Deidda & Puliga, 2006). We then apply the FDR test for both GOF tests, finding that  
207  $H_0$  cannot be rejected in more than 95% of the gages at  $\alpha_{global} = 0.05$  for all values of  $q$  and  $n$ . Given the  
208 very small number of rejections, the Poisson distribution is adopted as the parent distribution of count  
209 time series.

210

### 211 3.3 Generation of synthetic count time series

212 We conduct several Monte Carlo experiments based on the generation of random Poisson-  
213 distributed count time series that serve two main goals. The first is to gain insights on the open question  
214 raised by several authors (Yue et al., 2002; Hamed, 2009; Serinaldi & Kilsby, 2016) concerning the  
215 influence of serial correlation on trend detection and vice versa. In particular, we investigate (i) the degree  
216 of autocorrelation that can be detected in time series generated under controlled uncorrelated and  
217 nonstationary conditions, and, conversely, (ii) the trend induced by the presence of autocorrelation in  
218 time series generated under stationary conditions. The second goal of the Monte Carlo experiments is to  
219 generate the null distribution for the statistics of the trend tests (as described in Section 3.4) to account  
220 for discretization, sample length, and possible presence of autocorrelation. In such a way, we can also  
221 explore the type-I error and power of trend tests applied locally and at multiple sites. The generation of  
222 the synthetic count time series is described in the next subsections.

#### 223 3.3.1 Nonstationary uncorrelated time series

224 Under the assumption of Poisson distributed counts, we can easily generate synthetic time series  
225 with a controlled trend slope  $\phi$ , applying a linear time-varying relation for the Poisson parameter:

$$\lambda_t = \lambda_0 + \phi \cdot t, \quad t = 1, \dots, n \quad (2)$$

226 where the intercept  $\lambda_0$  is derived by constraining the mean value of  $\{\lambda_t\}$  to be  $\bar{\lambda} = (1 - q) \cdot 365.25$ , with  
227  $q$  being the selected nonexceedance probability. This results in  $\lambda_0 = \bar{\lambda} - \phi \cdot (n + 1)/2$ .

#### 228 3.3.2 Stationary correlated time series

229 We use the INAR(1) model to generate random autocorrelated stationary count time series. INAR  
230 models have been mainly applied in economics and finance (e.g., Blundell et al., 2002; Jung and  
231 Tremayne, 2011), epidemiology (e.g., Allard, 1998; Pascual & Akhundjanov, 2019), and insurance (e.g.,  
232 Gouriéroux & Jasiak, 2004; Boucher et al., 2008), but they have received less attention in hydrology and  
233 climatology. To define the INAR(1) process, we first introduce the binomial thinning operator, “ $\circ$ ”

234 (Steutel & van Harn, 1979). If  $\rho \in [0, 1]$  and  $N$  is a nonnegative integer random variable, this operator is  
 235 defined as:

$$\rho \circ N = \sum_{i=1}^N Y_i, \quad N > 0 \quad (3)$$

236 where  $\{Y_i\}$  are independent and identically distributed (i.i.d.) variates of a Bernoulli distribution  $B(\rho)$ .  
 237 While other thinning operators have been proposed (Weiß, 2008), here the binomial thinning operator is  
 238 used. A process  $\{N_t\}$  is defined INAR(1) if:

$$N_t = \rho \circ N_{t-1} + \epsilon_t \quad (4)$$

239 where  $\{\epsilon_t\}$  is an i.i.d. random process of integer values and the binomial thinning operator with parameter  
 240  $\rho$  is applied to  $N_{t-1}$ . Its lag- $k$  autocorrelation is  $r(k) = \rho^k$ , similar to the AR(1) model for real values.

241 In light of the results discussed in section 3.2, we adopt a Poisson-INAR(1) model to generate  
 242 synthetic correlated count series, where  $\{\epsilon_t\}$  is an i.i.d. random process according to a Poisson  
 243 distribution with parameter  $\mu$ , and the marginal distribution of  $\{N_t\}$  is also a Poisson distribution with  
 244 parameter  $\left(\frac{\mu}{1-\rho}\right)$  (Weiß, 2008). Parameters of the Poisson-INAR(1) model reproducing the statistical  
 245 properties of an observed count time series  $\{o_t\}$  can be estimated as:  $\rho = \rho_{\text{obs}}$ , with  $\rho_{\text{obs}}$  being the observed  
 246 lag-1 autocorrelation of  $\{o_t\}$ ; and  $\mu = (1 - \rho_{\text{obs}}) \bar{\lambda}$ , with  $\bar{\lambda} = (1 - q) \cdot 365.25$  being the expected number  
 247 of annual exceedances above the  $q$ -th quantile. An example of the capability of the Poisson-INAR(1)  
 248 model to reproduce the statistical properties of our observed counts is shown in Fig. 2, where the  
 249 empirical autocorrelation function of two randomly chosen count time series derived from the GHCN P  
 250 records is compared to the 95% confidence intervals (CIs) built from 10,000 model simulations with the  
 251 parameters estimated as just described. Fig. 2 shows that the Poisson-INAR(1) model captures very well  
 252 the empirical autocorrelations at different lags.

253

### 254 3.4 Setup of statistical tests through Monte Carlo simulations

255 To detect empirical trends in our analyses, we focus on three (two) nonparametric (parametric)  
256 statistical tests widely used in trend analyses of P extremes (Table 1). The nonparametric ones include  
257 Mann Kendall (Mann, 1945; Kendall, 1975); Kendall's  $\tau$  (Kendall, 1938; El-Shaarawi & Niculescu,  
258 1992); and Spearman's  $\rho$  (Gauthier, 2001). The parametric tests are based on linear and Poisson  
259 regression (Wilks, 2019). All these tests have been originally devised to investigate the null hypothesis  
260 of trend absence in uncorrelated time series. However, some authors have warned about the possible  
261 degraded test performances due to the possible presence of serial correlation in stationary time series  
262 (e.g., Serinaldi & Kilsby, 2016). To investigate this issue, we use Monte Carlo simulations to build the  
263 distribution of the test statistics under any  $H_0$  that may include uncorrelated and autocorrelated time  
264 series. In such a way, we also reduce potential biases introduced by finite sample sizes and discrete  
265 records (Deidda & Puliga, 2006), as well as by the presence of ties likely found in count time series.

266 In the general case of count time series of length  $n$  affected by serial correlation, a statistical trend  
267 detection test based on Monte Carlo simulations can be applied as follows:

- 268 1. The expected number of exceedances above the  $q$ -th quantile is estimated as  $\bar{\lambda} = (1 - q) \cdot 365.25$ .
- 269 2. Parameters of the Poisson-INAR(1) model in equation (4) are estimated as:  $\rho = \rho_{\text{obs}}$  and  
270  $\mu = (1 - \rho_{\text{obs}}) \cdot \bar{\lambda}$ .
- 271 3. An ensemble of  $n_{\text{ens}}$  (e.g.  $n_{\text{ens}} = 10,000$  in our applications) stationary count time series, each of  
272 length  $n$ , is generated using the Poisson-INAR(1) model with parameters estimated in step (2).
- 273 4. The  $s$  test statistic of interest (e.g.,  $s = \tau$  for Kendall's) is computed for each of the  $n_{\text{ens}}$  count time  
274 series generated in step (3).
- 275 5. The empirical cumulative distribution function (ECDF) of the  $n_{\text{ens}}$  test statistics from step (4) is  
276 used to determine the acceptance region of the null hypothesis. For example, for two-sided tests,  
277 this is the interval of  $s$ -quantiles corresponding to probabilities  $\alpha/2$  and  $(1 - \alpha/2)$ , for any

considered significance level  $\alpha$ . The local null hypothesis is therefore accepted or rejected by comparing the test statistic computed on the time series of interest,  $s_{\text{obs}}$ , with such acceptance region.

6. Similarly, the ECDF of the  $n_{\text{ens}}$  test statistics from step (4) is used to determine the  $p$ -value of  $s_{\text{obs}}$  (note that, for two-sided tests, as those selected here, the corresponding  $p$ -value has to be estimated by doubling the exceedance or nonexceedance probability in the ECDF).

7. If the test is conducted at  $M$  sites, the field significance is taken into account through the FDR test applied with the  $M$   $p$ -values determined at each site, as described in steps (1)-(6).

This procedure is general and can be implemented for any trend test by using the corresponding test statistic in steps (4)-(6) (see Appendix for details on the tests considered here). Moreover, with this method, different null hypotheses can be tested depending on the properties of the synthetic count series generated in step (3). We will use the following compact notation to describe the null hypothesis tested in this study, including:  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = 0$ ” for uncorrelated and stationary signals generated at step (3) from a Poisson distribution with parameter  $\bar{\lambda}$  (in this case, step (2) is skipped); and  $H_0$ : “ $\rho_0 = \rho^*$ ;  $\phi_0 = 0$ ” for serially correlated and stationary signals generated from the Poisson-INAR(1) model with parameter  $\rho = \rho^*$  (e.g.,  $\rho^* = \rho_{\text{obs}}$  in step (2)).

An analogous procedure can also be implemented to test whether a certain degree of autocorrelation detected in a count time series can be reasonably due to the presence of a given trend. In this case: the null hypothesis is  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = \phi^*$ ”; step (2) is skipped; an ensemble of  $n_{\text{ens}}$  nonstationary uncorrelated count time series of length  $n$  is generated in step (3) as described in section 3.3.1, using parameter  $\bar{\lambda}$  from step (1) and a given trend slope  $\phi^*$ ; finally, the lag-1 autocorrelation is used as test statistic in step (4) and the  $n_{\text{ens}}$  estimated lag-1 autocorrelations are utilized to compute the  $p$ -value associated with the observed autocorrelation.

## 4 Results and discussion

### 4.1. Investigation of autocorrelation and its relationship with linear trends

Deciding whether the possible influence of serial correlation in trend detection should be taken into account is not an easy question to answer, because, in principle, there can be a reciprocal feedback between autocorrelation and trend. To investigate this nontrivial issue in our count time series, we use two simple metrics to characterize autocorrelation and trend, namely the lag-1 autocorrelation,  $\rho$ , and the linear trend slope,  $\phi$ , respectively. We first compare the empirical distributions of  $\rho$  and  $\phi$  of the  $M = 1087$  observed count time series with the corresponding 95% CI of  $\rho$  and  $\phi$ , respectively, derived from  $n_{\text{ens}} = 10,000$  Monte Carlo simulations under  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = 0$ ”. In other words, we evaluate whether the observed  $\rho$ 's and  $\phi$ 's can be considered statistically different from those of uncorrelated time series with no trend. Results are shown in Fig. 3 for  $q = 0.95$  and different  $n$  (similar patterns are obtained for the other  $q$ 's; see Figs. S1, S2 and S3 in Supplementary Material). As expected, for both metrics the dispersion of the empirical distributions increases for smaller  $n$ . The simple visual comparison of distributions and 95% CIs suggests that  $H_0$  should be locally rejected for  $\rho$  in a relatively small number of sites (Figs. 3a-c), while the number of rejections appears to be much higher for  $\phi$  (Figs. 3d-f). These visual speculations are confirmed by results for the local test reported in Table 2 and, more importantly, by the application of the FDR test, which reveals that for  $n = 100$  only 3% (or 0% for  $n = 50$  and 30) of the observed  $\rho$ 's can be considered statistically significant at  $\alpha_{\text{global}} = 0.05$  significance level, while the percentage of statistically significant observed  $\phi$ 's is much larger (41%). Results for all considered  $n$  and local and FDR tests are reported in Table 2 and consistently show that, while a large number of sites seem to be affected by significant trend, the same conclusion does not hold for empirical serial correlation.

323 To further explore whether the presence of autocorrelation may introduce bias in the estimation  
 324 of the linear trend slope, we analyze the joint distribution of  $\rho$  and  $\phi$  estimated on the  $M$  observed 100-  
 325 year time series. The scatterplot between these values is plotted in Fig. 4a (grey circles) along with  
 326 estimates derived from  $M$  random stationary and uncorrelated time series (black circles). The visual  
 327 inspection clearly suggests that the observations do not appear consistent with a hypothesis of both no  
 328 autocorrelation and no trend. In particular, the observed counts exhibit more cases with higher slope  
 329 (both positive and negative) that are associated with higher autocorrelation. To gain insights on the  
 330 potential cause-effect relationship of this outcome (i.e., is the autocorrelation causing an artificial trend  
 331 or is the opposite true? Or are these effects independent?), we first evaluate whether the presence of trend  
 332 can artificially induce autocorrelation. For this aim, we generate time series under  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = \phi$ ”,  
 333 with  $\phi$  varying from -0.2 to 0.2 events/yr to cover the whole range of observed trend slopes for  $n = 100$ .  
 334 For each value of  $\phi$ , we produce  $n_{\text{ens}} = 10,000$  samples, estimate  $\rho$  on each time series, and derive the  
 335 95% CI of  $\rho$  (solid lines in Fig. 4b). We find that 95% of the observed  $(\rho, \phi)$  pairs lie within the CI,  
 336 indicating that the observed  $\rho$ 's, even if different from zero, are compatible with those of uncorrelated  
 337 series with trend.

338 Following a similar framework, we then investigate whether the presence of autocorrelation could  
 339 artificially induce significant trends. We do so by computing the 95% CI of  $\phi$  from time series randomly  
 340 generated under  $H_0$ : “ $\rho_0 = \rho$ ;  $\phi_0 = 0$ ”, with  $\rho$  varying from 0 to 0.8 (solid line in Fig. 4c). In this case, a  
 341 large fraction (40%) of observed  $(\rho, \phi)$  pairs lies outside of this CI, implying that several high values of  
 342  $\phi$  cannot be explained solely by the presence of autocorrelation. The same conclusion can be drawn by  
 343 comparing this CI with that obtained under  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = 0$ ” (dotted line in Fig. 4c): the two CIs are  
 344 very close to each other, meaning that accounting or not for the possible presence of serial correlation  
 345 has a very limited impact on the assessment of trend significance. The only region where the trend

significance could be potentially ascribed to the presence of autocorrelation is the area between the two CIs, which includes only a very limited number of observed cases. It is also worth noticing that such a few cases would be certainly less if one rightly considers only the component of autocorrelation that is not ascribed to the presence of trend, which results in a positive overestimation of  $\rho$ , as also clearly reflected in the CIs shown in Fig. 4b (see also Yue & Wang, 2002).

Results presented in Fig. 4 suggest that autocorrelation in observed count time series of extreme P is likely caused by the presence of trends. To complement this conclusion relying on statistical simulations, we provide further evidence based on the physical argument that temporal persistency (if any) in extreme P should significantly decrease after a few years. From each observed time series, we sample the record every four years, thus extracting four sub-series of size  $n = 25$ ; in such a way, we eliminate the effect of potential autocorrelations at lags from 1 to 3 years. For each sub-series, we estimate  $\phi$  and plot it against the slope estimated on the full series. Results are presented in Fig. 5a, which shows that, despite some expected sampling variability, all values are distributed along the 1:1 line. In addition, we randomly generate  $M$  uncorrelated series of duration  $n = 100$  with the same  $M$  slopes estimated on the observed series, and, for each synthetic sample, we repeat the same calculation on four sub-series of size  $n = 25$  sampled every four years. The corresponding outcome, reported in Fig. 5b, is consistent with results for the observed series, thus providing further evidence that statistically significant trends exist in our observed count time series, independently of the possible presence of autocorrelation.

#### 4.2. Performance of local trend tests

After analyzing the relations between trend and possible presence of autocorrelation, we now use Monte Carlo simulations to investigate if accounting or not for autocorrelation can affect the power of local trend tests. To this end, we generate 10,000 nonstationary uncorrelated time series for different values of  $\phi$ ,  $n$  and  $q$  using equation (2) as described in Section 3.3.1. For each combination of  $\phi$ ,  $n$  and  $q$ , we estimate the test power as the fraction of rejections of the null hypotheses  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = 0$ ” and



370  $H_0$ : “ $\rho_0 = \rho_{\text{obs}}$ ;  $\phi_0 = 0$ ”, applying the trend tests as described in Section 3.4. Results are presented in Fig.  
 371 6, where dotted and solid lines are used for the two  $H_0$  settings and colors refer to different tests. For  $n =$   
 372 100 years, Fig. 6a shows that the power of all tests increases in quasi-linear fashion from 0.05 (the test  
 373 significance level) at  $\phi = 0$  to  $\sim 0.9$  at  $\phi = 0.05$  events/yr, reaching 1 for  $\phi > 0.07$  events/yr. As expected,  
 374 for a given  $\phi$ , the test power decreases with  $n$  (Fig. 6b). For  $\phi \leq 0.05$  events/yr, the power is less than 0.5  
 375 for  $n \leq 70$  years, indicating that the statistical tests analyzed here have low ability to detect trends even  
 376 when  $n$  is relatively large. The use of  $H_0$ : “ $\rho_0 = \rho_{\text{obs}}$ ;  $\phi_0 = 0$ ” leads to a slight power reduction compared  
 377 to  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = 0$ ”, a further indication that taking or not taking into account autocorrelation does not  
 378 significantly impact results. Finally, as better shown in Fig. 6b, parametric (linear and Poisson regression)  
 379 and nonparametric (Mann Kendall, Kendall’s  $\tau$  and Spearman  $\rho$ ) tests cluster in two separate groups,  
 380 with the parametric tests exhibiting slightly higher power than the nonparametric ones. Based on these  
 381 findings, we will discuss trends in observed count time series in section 4.4 presenting results only for  
 382 the Poisson regression (PR) and Mann Kendall (MK) tests, which are representative of parametric and  
 383 nonparametric tests, respectively. The difference in power between these two tests as a function of  $\phi$  for  
 384  $n = 100$  years is reported in Fig. 7.

### 385 4.3. Performance of trend tests at multiple sites

386 We gain insights on tests’ performance at multiple locations by conducting synthetic experiments  
 387 on a  $50 \times 100$  grid totaling  $M = 5,000$  sites, where we hypothesize the existence of trend only in an inner  
 388 rectangular domain containing 30% of the grid points. In each site of this region, we generate count time  
 389 series with a given linear trend slope  $\phi$ , while, in the remaining grid points placed in the outer region, we  
 390 generate stationary time series. We do this for  $q = 0.95$  and for  $n = 50$  and 100 years. We discuss here  
 391 results for PR trend tests (results are similar for other tests) applied locally under  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = 0$ ”,  
 392 and globally by accounting for field significance with the FDR test at  $\alpha_{\text{FDR}} = \alpha_{\text{global}}$  (there is no spatial

393 correlation in this experiment). Fig. 8a (Fig. 8b) presents the fraction of  $H_0$  rejections in the inner region  
394 with trends (outer region without trends) as a function of  $\phi$ , quantifying test power (type-I error) in that  
395 part of the domain. For small trend slopes, local tests lead to higher power (differences of up to 0.5  
396 compared to global results), but such discrepancies approach zero as  $\phi$  increases (Fig. 8a). As found for  
397 the local analyses, for a given  $\phi$ , the power is heavily affected by the sample size. For example, for  $\phi =$   
398 0.05 events/yr, the power of the global test drops from 0.8 for  $n = 100$  years to zero for  $n = 50$  years. On  
399 the other hand, applying tests locally without considering field significance leads to much larger type-I  
400 errors in the outer region for any  $\phi$  (Fig. 8b). In other words, the use of local tests leads to several false  
401 rejections of  $H_0$  that the FDR test is able to prevent. In this case, the effect of the sample size is negligible.

402 To visually illustrate performance of tests conducted at multiple sites, we refer to the same  $50 \times$   
403 100 grid with time series in the inner region generated with  $\phi = 0.05$  events/yr, for  $n = 100$  and 50 years.  
404 Figs. 8c-f present maps of test results applied locally and globally, with red (green) colors indicating  $H_0$   
405 rejections for the PR when the trend slope estimated on the generated time series is positive (negative).  
406 We first focus on the maps for  $n = 100$  years (Figs. 8c,d). When tests are performed locally (Fig. 8c),  $H_0$   
407 is rejected, as expected, at  $\sim 5\%$  of the locations in the outer region. This would erroneously indicate  
408 statistically significant trends at sites where trend is not present, inducing wrong physical interpretations  
409 if these sites coincidentally cluster. Accounting for field significance with the FDR test (Fig. 8d) leads  
410 instead to the rejection of  $H_0$  at just a few spurious locations ( $\sim 1\%$  of the points in the outer region). In  
411 this condition, it is more meaningful to interpret these rejections as a result of randomness rather than  
412 physical processes. When considering the inner region with trends, the application of the more  
413 conservative (i.e.,  $H_0$  is rejected less) FDR test returns a higher number of false nonrejections of  $H_0$   
414 compared to the local test (22.8% vs 8.9% of the cases). However, despite the lower power (also  
415 highlighted in Fig. 8a),  $H_0$  is rejected at most locations that are spatially clustered, so that the region with  
416 trend could be readily identified.

417 When  $n = 50$  years, results for the local tests (Fig. 8e) do not change in the outer region, with  
 418 random occurrence of  $H_0$  rejections at  $\sim 5\%$  of the points with positive and negative slopes as found for  
 419  $n = 100$  years. The reduction of test power due to the smaller sample size leads instead to less  $H_0$   
 420 rejections in the inner region. Changes are even more drastic when applying the more conservative FDR  
 421 test, which results in  $H_0$  nonrejections at all sites (Fig. 8f). This outcome suggests that, when the trend  
 422 signal is low, the use of methods accounting for field significance will likely indicate the absence of  
 423 statistically significant trends. In this circumstance, a careful interpretation of results of the more  
 424 powerful local tests could still allow identifying large areas characterized by statistically significant  
 425 trends if the sites exhibit coherent positive or negative trend. This is depicted in the example of Fig. 8e,  
 426 where positive trends are correctly detected at a number of nearby locations that is sufficiently large to  
 427 identify the inner region. In the outer region, the mixture of both positive and negative trends in sites  
 428 close to each other should suggest that no trend signal is detectable in such area. This issue will be further  
 429 discussed in the next section.

#### 430 4.4. Trend analyses of observed count series

431 In light of the insights gained in the previous sections, we now analyze the presence of trends in  
 432 observed count series on the  $M = 1087$  selected stations from the GHCN gage network. Trends are  
 433 investigated applying the PR and MK tests and, then, the FDR test at  $\alpha_{global} = 0.05$  to account for field  
 434 significance. We preliminarily considered two null hypotheses: stationary and uncorrelated signals, and  
 435 stationary and autocorrelated series. Regarding the second null hypothesis, our previous analyses have  
 436 shown that a large portion (or perhaps all) of the lag-1 autocorrelation estimated on the observed sample,  
 437  $\rho_{obs}$ , is likely induced by the presence of trend (see Fig. 4b). As a result, when testing the null hypothesis  
 438 of autocorrelated signals, we should consider only the residual component of  $\rho_{obs}$  that cannot be ascribed  
 439 to the presence of trend (see discussion in Section 4.1). Considering that implementing such an approach  
 440 is not straightforward, the trend tests were preliminary applied under  $H_0: \rho_0 = 0; \phi_0 = 0$  and  $H_0: \rho_0 =$

441  $\rho_{\text{obs}}; \phi = 0$ ", which represent two extreme conditions. Since we found very similar results and patterns  
442 in the two cases (not shown), we hereon discuss only results for  $H_0: \rho = 0; \phi = 0$ .

443 Fig. 9 presents maps of statistically significant trends for  $q = 0.95$  and  $n = 100$  years. Colored  
444 circles and triangles locate significant trends for (i) only PR and (ii) both PR and MK tests, respectively.  
445 We first note that, as suggested by the synthetic experiments, PR detects a larger number of statistically  
446 significant trends than MK, while the opposite never occurs. This is better visualized in the scatterplot  
447 between  $\phi$  and  $\rho$  of Fig. 10, where  $H_0$  rejections by only PR or both PR and MK tests are plotted with  
448 different markers. The occurrence of the different cases is controlled by  $\phi$ , while  $\rho$  is not influential, thus  
449 providing additional evidence on the limited effect of autocorrelation on trend detection. In particular,  
450  $H_0$  is rejected by both tests for  $|\phi| > \sim 0.05$  events/yr, which is a region where the power of all tests is  
451 high for  $n = 100$  years (Fig. 6a).  $H_0$  is rejected only by PR at several sites where  $|\phi|$  is included between  
452 roughly 0.02 and 0.05 events/yr, where the power of both tests decreases (Fig. 6a) but is larger for PR  
453 than MK (Fig. 7). This behavior can, at least partially, explain why the parametric PR rejects  $H_0$  in more  
454 cases than the nonparametric MK test.

455 Despite PR leads to rejection of  $H_0$  at several sites where our synthetic experiments suggest low  
456 test power, Fig. 9 clearly shows that locations where trends are statistically significant are well clusterized  
457 in space, with distinct regions where the trend is either increasing (red symbols) or decreasing (green  
458 symbols). As shown in the synthetic experiments at multiple sites of Fig. 8, the presence of spatial clusters  
459 provides further evidence of trend existence. This empirical result is also supported by the physical  
460 argument that extreme P is often controlled by synoptic processes (Barlow et al., 2019), and that their  
461 occurrence is changing in time (Zhang & Villarini, 2019). As a result, when trends exist, they should  
462 manifest over relatively large regions and, if multiple gages are present, statistical tests should detect  
463 statistically significant trends with the same sign at several of these sites (e.g., Kunkel et al., 2020). In  
464 particular, consistent with previous work with global and regional datasets (Table 1), our analyses reveal

465 that significant trends are mainly increasing in central and eastern North America (Janssen et al., 2014),  
466 northern Europe (Madsen et al., 2014), northern Asia (Zolotokrylin and Cherenkova, 2017), and central  
467 regions of Australia (Gallant et al., 2007). Extreme P exhibit instead negative trends in southwestern  
468 North America (Hoerling et al., 2016), part of southern Europe (Papalexiou & Montanari, 2019), and  
469 southwestern and southeastern regions of Australia close to the coast (Hughes, 2003).

470 The synthetic experiments indicate that the tests' power could be severely reduced when the  
471 sample size decreases. We analyze this issue on the observed count time series by plotting in Fig. 11a  
472 the maps of  $H_0$  rejections by the FDR test applied on PR and MK for  $q = 0.95$  and  $n = 50$  years (results  
473 for  $n = 30$  years are presented in Fig. S6). When compared to Fig. 9, the number of  $H_0$  rejections  
474 dramatically declines. The only regions with a relatively large number of spatially clustered gages that  
475 exhibit statistically significant trends are northern Europe (increasing trend) and southern Australia  
476 (decreasing trend). In North America, there are some gages where  $H_0$  is rejected, but their location is  
477 quite sparse, although there is a relatively clear geographical distinction between increasing and  
478 decreasing trends. In this circumstance where the trend signal might be weak, local test results could be  
479 used to complement results of the more conservative FDR test. As shown in Fig. 11b, local  $H_0$  rejections  
480 have a well-defined spatial pattern with two large regions where the trend sign is the same: central and  
481 northeastern (southwestern) North America, with increasing (decreasing) trend, which are the same  
482 regions identified in Fig. 9 for  $n = 100$  years. To complete our analysis, we investigate the role of the  
483 nonexceedance probability  $q$ , which controls the threshold used to build the count series of extreme P.  
484 Fig. 12 displays maps of global tests results for  $n = 100$  years for  $q = 0.90$  and  $0.975$  (results for  $n = 30$   
485 years are presented in Figs. S4-S7). As  $q$  increases and focus is placed on rarer events, less statistically  
486 significant trends are detected, but the spatial patterns of increasing and decreasing trends in the different  
487 regions of the worlds are always clearly visible.

488

## 5 Summary and conclusions

Increasing evidence and theoretical arguments indicate that global warming is causing and will cause changes in extreme P. Accurate statistical trend analyses of observed and modeled P time series are key to validate hypotheses on the underlying physical mechanisms and improve our ability to predict the magnitude of these changes. In this study, we clarified how autocorrelation and field significance affect application, power, and interpretation of several popular tests for trend detection in count time series. We focused on count time series because stronger trends have been detected in extreme P frequencies rather than magnitudes. We used observed records of extreme P frequency in the 100-year period from 1916 to 2015 collected by 1087 high-quality rain gages of the GHNC network, covering North America, part of Europe and Asia, and Australia. To investigate the role of autocorrelation and field significance and interpret trends in observed records, we designed several Monte Carlo experiments based on the random generation of stationary and nonstationary count time series with different levels of autocorrelation and sample size. The experiments involved the use of the Poisson-INAR(1) model that has been rarely adopted in hydroclimatic applications. Our results can be summarized as follows:

1. Although some observed count time series may exhibit some degree of autocorrelation (quantified through the lag-1 autocorrelation,  $\rho$ ), we proved that such correlations are mainly consistent with those of uncorrelated and either stationary or nonstationary count time series with the same sample size. We observed that records exhibiting stronger trends (quantified through the linear slope,  $\phi$ ) are also characterized by high  $\rho$  values; in these cases, using statistical arguments, we proved that the empirical high  $\rho$  values are compatible with uncorrelated time series with trends of the same observed magnitude. Conversely, we also proved that high trend slopes cannot be interpreted as a spurious outcome of a stationary autocorrelated signals. As a result, autocorrelation in observed count time series of extreme P appears to be caused by the

512 presence of trends, indicating that taking or not taking into account its presence when applying  
513 statistical trend tests does not significantly impact results.

- 514 2. As expected, the power of trend tests is importantly affected by sample size,  $n$ , of the analyzed  
515 series and trend magnitude,  $\phi$ . For example, considering the occurrences of daily precipitation  
516 with nonexceedance probability  $q = 0.95$  and a trend slope  $\phi = 0.05$  events/yr, the power is lower  
517 than 0.5 when  $n \leq 70$  years, which is a relatively long record. The power of parametric tests (linear  
518 and Poisson regression) is slightly larger than that of nonparametric tests (Mann Kendall,  
519 Kendall's  $\tau$  and Spearman  $\rho$ ).
- 520 3. Trend tests are in most cases applied at multiple locations. Here, we confirmed that, if test  
521 multiplicity or field significance is not taken into account, type-I errors could be large and  
522 statistically significant trends could be mistakenly detected at several sites, inducing wrong  
523 physical interpretations when these locations tend to coincidentally cluster. Accounting for field  
524 significance severely reduces this problem. On the other hand, we also showed that the inclusion  
525 of field significance leads to a power reduction compared to local tests. While this issue is  
526 practically irrelevant when the trend signal is moderate and high, it may result in several incorrect  
527 nonrejections of  $H_0$ , especially when the sample size is small. To limit this, the careful  
528 interpretation of results of local tests could help correctly identify trends in large regions where:  
529 (i) several gages are present; (ii) local tests reject  $H_0$  at most locations; and (iii) the trend detected  
530 in close gages has the same sign. These recommendations are supported by the empirical analyses  
531 of observed records presented here, as well as by the physical evidence that extreme P is mainly  
532 driven by large-scale processes whose occurrence has been changing in time. In such a way, the  
533 power of regional trend analyses is expected to increase, a task highly desirable to support  
534 engineering design against natural hazards (Vogel et al., 2013).

4. The application of several trend tests on the selected 1087 rain gages of the GHNC network reveals statistically significant increasing trends in several parts of the world, including central and eastern North America, northern Europe, part of northern Asia, and central regions of Australia. Decreasing trends are instead found in southwestern North America, part of southern Europe, and southwestern and southeastern regions of Australia. These results are largely consistent with previous studies.

Our work provides useful guidance for a more informed application of statistical trend tests in regional and global trend analyses of hydroclimatic extremes, and for a more realistic interpretation of test results.

#### **Acknowledgment**

This work has been supported by Fondazione Banco di Sardegna, funding call 2017, project: Impacts of climate change on water resources and floods, CUP: F71I17000270002. The data used in this study are publicly available via the Global Historical Climatology Network website: <https://www.ncei.noaa.gov/access/metadata/landing-page/bin/iso?id=gov.noaa.ncdc:C00861>.



## 549 Appendix A

550 Given the count time series  $\{o_k\}$  with  $k = 1, \dots, n$ , we investigate the existence of trend through  
 551 three nonparametric tests (Mann Kendall, Kendall's  $\tau$ , and Spearman's  $\rho$ ) and two parametric tests (test  
 552 on linear regression slope and Poisson regression). In the following, we report the statistics of each test,  
 553 which are used to build the null distribution via Monte Carlo simulations as described in section 3.4.

554 The Mann-Kendall test is based on the test statistic  $S$  calculated as:

$$S = \sum_{j=1}^{n-1} \sum_{k=j+1}^n \text{sign}(o_j - o_k) = \sum_{j=1}^{n-1} \sum_{k=j+1}^n \text{sign}(R_j - R_k) \quad (\text{A1})$$

555 where  $o_j$  and  $o_k$  represent  $j$ -th and  $k$ -th values of the count time series,  $R_j$  and  $R_k$  the corresponding ranks,  
 556  $n$  is the length of the series and:

$$\text{sign}(R_j - R_k) = \begin{cases} 1 & \text{for } (R_j - R_k) > 0 \\ 0 & \text{for } (R_j - R_k) = 0 \\ -1 & \text{for } (R_j - R_k) < 0 \end{cases} \quad (\text{A2})$$

557 In the Spearman's rank correlation test, the following  $r_s$  test statistics is used:

$$r_s = 1 - \frac{6 \sum_{i=1}^n (R_i - i)^2}{n(n^2 - 1)} \quad (\text{A3})$$

558 The Kendall's  $\tau$  test is based on a measure of the rank correlation evaluated as follows:

$$\tau = \frac{2}{n(n-1)} \sum_{i=1}^{n-1} \sum_{j=i+1}^n \text{sign}(R_i - R_j) \text{sign}(i - j) \quad (\text{A4})$$

559 The trend test on linear regression slope is based on the regression between  $\{o_k\}$  and  $k = 1, \dots,$   
 560  $n$  as follows:

$$\mu_k = b_0 + b_1 k, \quad (\text{A5})$$

561 where  $\mu_k$  is the predicted value, and  $b_0$  and  $b_1$  are parameters estimated through the least squares  
 562 approach. Here, we apply the trend test using  $b_1$  as test statistic.

563           The Poisson regression is a generalized linear model that links a Poisson-distributed variable with  
564 a set of predictors. Here, we consider only one predictor. The model relates the logarithm of the  $\mu$   
565 parameter of the parent Poisson distribution of the predictand with the predictor as:

$$\ln(\mu_k) = b_0 + b_1 k \quad (\text{A6})$$

566   The statistics used to apply the test under the proposed modification with Monte Carlo simulations is  $b_1$ .

567

## 568    **References**

- 569    Alexander, L. V., Zhang, X., Peterson, T. C., Caesar, J., Gleason, B., Klein Tank, A. M. G., et al.  
570        (2006). Global observed changes in daily climate extremes of temperature and precipitation.  
571        *Journal of Geophysical Research*, 111(D5), D05109. <https://doi.org/10.1029/2005JD006290>
- 572    Allard, R. (1998). Use of time-series analysis in infectious disease surveillance. *Bulletin of the World*  
573        *Health Organization*, 76(4), 327–333.
- 574    Al-Osh, M. A., & Alzaid, A. A. (1987). First-Order Integer-Valued Autoregressive (INAR (1)) Process.  
575        *Journal of Time Series Analysis*, 8(3), 261–275. <https://doi.org/10.1111/j.1467->  
576        9892.1987.tb00438.x
- 577    Alpert, P., Ben-Gai, T., Baharad, A., Benjamini, Y., Yekutieli, D., Colacino, M., et al. (2002). The  
578        paradoxical increase of Mediterranean extreme daily rainfall in spite of decrease in total values.  
579        *Geophysical Research Letters*, 29(11), 31–1. <https://doi.org/10.1029/2001GL013554>
- 580    Asadieh, B., & Krakauer, N. Y. (2015). Global trends in extreme precipitation: climate models versus  
581        observations. *Hydrol. Earth Syst. Sci.*, 19(2), 877–891. <https://doi.org/10.5194/hess-19-877->  
582        2015
- 583    Ashley, S. T., & Ashley, W. S. (2008). Flood Fatalities in the United States. *Journal of Applied*  
584        *Meteorology and Climatology*, 47(3), 805–818. <https://doi.org/10.1175/2007JAMC1611.1>
- 585    Barlow, M., Gutowski, W. J., Gyakum, J. R., Katz, R. W., Lim, Y.-K., Schumacher, R. S., et al. (2019).  
586        North American extreme precipitation events and related large-scale meteorological patterns: a  
587        review of statistical methods, dynamics, modeling, and trends. *Climate Dynamics*, 53(11),  
588        6835–6875. <https://doi.org/10.1007/s00382-019-04958-z>
- 589    Bayazit, M., & Önöz, B. (2007). To prewhiten or not to prewhiten in trend analysis? *Hydrological*  
590        *Sciences Journal*, 52(4), 611–624. <https://doi.org/10.1623/hysj.52.4.611>
- 591    Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful

592 Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B*  
593 *(Methodological)*, 57(1), 289–300.

594 Blundell, R., Griffith, R., & Windmeijer, F. (2002). Individual effects and dynamics in count data  
595 models. *Journal of Econometrics*, 108(1), 113–131. [https://doi.org/10.1016/S0304-](https://doi.org/10.1016/S0304-4076(01)00108-7)  
596 4076(01)00108-7

597 Boggess, J. M., Becker, G. W., & Mitchell, M. K. (2014). Storm & flood hardening of electrical  
598 substations. In *2014 IEEE PES T&D Conference and Exposition* (pp. 1–5).  
599 <https://doi.org/10.1109/TDC.2014.6863387>

600 Boucher, J.-P., Denuit, M., & Guillen, M. (2008). Models of Insurance Claim Counts with Time  
601 Dependence Based on Generalisation of Poisson and Negative Binomial Distributions.  
602 *Variance*, 2(1), 135–162.

603 Bucar, R. C. B., & Hayeri, Y. M. (2020). Quantitative assessment of the impacts of disruptive  
604 precipitation on surface transportation. *Reliability Engineering & System Safety*, 203, 107105.  
605 <https://doi.org/10.1016/j.ress.2020.107105>

606 Cann, K. F., Thomas, D. R., Salmon, R. L., Wyn-Jones, A. P., & Kay, D. (2013). Extreme water-  
607 related weather events and waterborne disease. *Epidemiology and Infection*, 141(4), 671–686.  
608 <https://doi.org/10.1017/S0950268812001653>

609 Daniel, J. S., Portmann, R. W., Solomon, S., & Murphy, D. M. (2012). Identifying weekly cycles in  
610 meteorological variables: The importance of an appropriate statistical analysis. *Journal of*  
611 *Geophysical Research: Atmospheres*, 117(D13). <https://doi.org/10.1029/2012JD017574>

612 Deidda, R., & Puliga, M. (2006). Sensitivity of goodness-of-fit statistics to rainfall data rounding off.  
613 *Physics and Chemistry of the Earth, Parts A/B/C*, 31, 1240–1251.  
614 <https://doi.org/10.1016/j.pce.2006.04.041>

615 Durre, I., Menne, M. J., Gleason, B. E., Houston, T. G., & Vose, R. S. (2010). Comprehensive

Automated Quality Assurance of Daily Surface Observations. *Journal of Applied Meteorology and Climatology*, 49(8), 1615–1633. <https://doi.org/10.1175/2010JAMC2375.1>

El-Shaarawi, A. H., & Niculescu, S. P. (1992). On kendall's tau as a test of trend in time series data. *Environmetrics*, 3(4), 385–411. <https://doi.org/10.1002/env.3170030403>

Emori, S., & Brown, S. J. (2005). Dynamic and thermodynamic changes in mean and extreme precipitation under changed climate. *Geophysical Research Letters*, 32(17). <https://doi.org/10.1029/2005GL023272>

Fisher, R. A. (1922). On the Interpretation of  $\chi^2$  from Contingency Tables, and the Calculation of P. *Journal of the Royal Statistical Society*, 85(1), 87–94. <https://doi.org/10.2307/2340521>

Gallant, A. E., Hennessy, K., & Risbey, J. S. (2007). Trends in rainfall indices for six Australian regions: 1910–2005. *Australian Meteorological Magazine*, 56, 223–239.

Gauthier, T. D. (2001). Detecting Trends Using Spearman's Rank Correlation Coefficient. *Environmental Forensics*, 2(4), 359–362. <https://doi.org/10.1006/enfo.2001.0061>

Gershunov, A., Benmarhnia, T., & Aguilera, R. (2018). Human health implications of extreme precipitation events and water quality in California, USA: a canonical correlation analysis. *The Lancet Planetary Health*, 2, S9. [https://doi.org/10.1016/S2542-5196\(18\)30094-9](https://doi.org/10.1016/S2542-5196(18)30094-9)

Gourieroux, C., & Jasiak, J. (2004). Heterogeneous INAR(1) model with application to car insurance. *Insurance: Mathematics and Economics*, 34(2), 177–192. <https://doi.org/10.1016/j.insmatheco.2003.11.005>

Groisman, P. Y., Knight, R. W., Easterling, D. R., Karl, T. R., Hegerl, G. C., & Razuvaev, V. N. (2005). Trends in Intense Precipitation in the Climate Record. *Journal of Climate*, 18(9), 1326–1350. <https://doi.org/10.1175/JCLI3339.1>

Hamed, K. H. (2009). Enhancing the effectiveness of prewhitening in trend analysis of hydrologic data. *Journal of Hydrology*, 368(1), 143–155. <https://doi.org/10.1016/j.jhydrol.2009.01.040>

640 Hamed, K. H., & Ramachandra Rao, A. (1998). A modified Mann-Kendall trend test for autocorrelated  
641 data. *Journal of Hydrology*, 204(1), 182–196. [https://doi.org/10.1016/S0022-1694\(97\)00125-X](https://doi.org/10.1016/S0022-1694(97)00125-X)

642 Hansen, J., Ruedy, R., Sato, M., & Lo, K. (2010). Global surface temperature change. *Reviews of*  
643 *Geophysics*, 48(4). <https://doi.org/10.1029/2010RG000345>

644 Hennessy, K. J., Suppiah, R., & Page, C. M. (1999). Australian rainfall changes, 1910–1995. *Aust. Met.*  
645 *Mag.*, 48, 1–13.

646 Hoerling, M., Eischeid, J., Perlwitz, J., Quan, X.-W., Wolter, K., & Cheng, L. (2016). Characterizing  
647 Recent Trends in U.S. Heavy Precipitation. *Journal of Climate*, 29(7), 2313–2332.  
648 <https://doi.org/10.1175/JCLI-D-15-0441.1>

649 Hooper, E., Chapman, L., & Quinn, A. (2014). The impact of precipitation on speed–flow relationships  
650 along a UK motorway corridor. *Theoretical and Applied Climatology*, 117(1), 303–316.  
651 <https://doi.org/10.1007/s00704-013-0999-5>

652 Hughes, L. (2003). Climate change and Australia: Trends, projections and impacts. *Austral Ecology*,  
653 28(4), 423–443. <https://doi.org/10.1046/j.1442-9993.2003.01300.x>

654 Janssen, E., Wuebbles, D. J., Kunkel, K. E., Olsen, S. C., & Goodman, A. (2014). Observational- and  
655 model-based trends and projections of extreme precipitation over the contiguous United States.  
656 *Earth's Future*, 2(2), 99–113. <https://doi.org/10.1002/2013EF000185>

657 Jung, R. C., & Tremayne, A. R. (2011). Convolution-closed models for count time series with  
658 applications. *Journal of Time Series Analysis*, 32(3), 268–280. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-9892.2010.00697.x)  
659 [9892.2010.00697.x](https://doi.org/10.1111/j.1467-9892.2010.00697.x)

660 Katz, R. W. (1988). Statistical Procedures for Making Inferences about Climate Variability. *Journal of*  
661 *Climate*, 1(11), 1057–1064. [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0442(1988)001<1057:SPFMIA>2.0.CO;2)  
662 [0442\(1988\)001<1057:SPFMIA>2.0.CO;2](https://doi.org/10.1175/1520-0442(1988)001<1057:SPFMIA>2.0.CO;2)

663 Katz, R. W., & Brown, B. G. (1991). The problem of multiplicity in research on teleconnections.

664 *International Journal of Climatology*, 11(5), 505–513. <https://doi.org/10.1002/joc.3370110504>

665 Kendall, M. G. (1938). A New Measure of Rank Correlation. *Biometrika*, 30(1/2), 81–93.

666 <https://doi.org/10.2307/2332226>

667 Kendall, M. G. (1975). *Rank correlation methods*. London: Griffin.

668 Khaliq, M. N., Ouarda, T. B. M. J., Gachon, P., Sushama, L., & St-Hilaire, A. (2009). Identification of

669 hydrological trends in the presence of serial and cross correlations: A review of selected

670 methods and their application to annual flow regimes of Canadian rivers. *Journal of Hydrology*,

671 368(1), 117–130. <https://doi.org/10.1016/j.jhydrol.2009.01.035>

672 Koutsoyiannis, D. (2011). Hurst-Kolmogorov Dynamics and Uncertainty1. *Journal of the American*

673 *Water Resources Association*, 47(3), 481–495. <https://doi.org/10.1111/j.1752->

674 1688.2011.00543.x

675 Kruger, A., & Nxumalo, M. (2017). Historical rainfall trends in South Africa: 1921–2015. *Water SA*,

676 43, 285. <https://doi.org/10.4314/wsa.v43i2.12>

677 Kunkel, K. E., & Frankson, R. M. (2015). Global Land Surface Extremes of Precipitation: Data

678 Limitations and Trends. *Journal of Extreme Events*, 02(02), 1550004.

679 <https://doi.org/10.1142/S2345737615500049>

680 Kunkel, K. E., Karl, T. R., Squires, M. F., Yin, X., Stegall, S. T., & Easterling, D. R. (2020).

681 Precipitation Extremes: Trends and Relationships with Average Precipitation and Precipitable

682 Water in the Contiguous United States. *Journal of Applied Meteorology and Climatology*,

683 59(1), 125–142. <https://doi.org/10.1175/JAMC-D-19-0185.1>

684 Li, Y., Guan, K., Schnitkey, G. D., DeLucia, E., & Peng, B. (2019). Excessive rainfall leads to maize

685 yield loss of a comparable magnitude to extreme drought in the United States. *Global Change*

686 *Biology*, 25(7), 2325–2337. <https://doi.org/10.1111/gcb.14628>

687 Livezey, R. E., & Chen, W. Y. (1983). Statistical Field Significance and its Determination by Monte

688 Carlo Techniques. *Monthly Weather Review*, 111(1), 46–59. <https://doi.org/10.1175/1520->  
689 0493(1983)111<0046:SFSaid>2.0.CO;2

690 Madsen, H., Lawrence, D., Lang, M., Martinkova, M., & Kjeldsen, T. R. (2014). Review of trend  
691 analysis and climate change projections of extreme precipitation and floods in Europe. *Journal*  
692 *of Hydrology*, 519, 3634–3650. <https://doi.org/10.1016/j.jhydrol.2014.11.003>

693 Mann, H. B. (1945). Nonparametric Tests Against Trend. *Econometrica*, 13(3), 245–259.  
694 <https://doi.org/10.2307/1907187>

695 McKenzie, Ed. (1985). Some Simple Models for Discrete Variate Time Series. *Journal of the American*  
696 *Water Resources Association*, 21(4), 645–650. <https://doi.org/10.1111/j.1752->  
697 1688.1985.tb05379.x

698 Menne, M. J., Durre, I., Vose, R. S., Gleason, B. E., & Houston, T. G. (2012). An Overview of the  
699 Global Historical Climatology Network-Daily Database. *Journal of Atmospheric and Oceanic*  
700 *Technology*, 29(7), 897–910. <https://doi.org/10.1175/JTECH-D-11-00103.1>

701 New, M., Hewitson, B., Stephenson, D. B., Tsiga, A., Kruger, A., Manhique, A., et al. (2006).  
702 Evidence of trends in daily climate extremes over southern and west Africa. *Journal of*  
703 *Geophysical Research: Atmospheres*, 111(D14). <https://doi.org/10.1029/2005JD006289>

704 Nie, J., Sobel, A. H., Shaevitz, D. A., & Wang, S. (2018). Dynamic amplification of extreme  
705 precipitation sensitivity. *Proceedings of the National Academy of Sciences*, 115(38), 9467.  
706 <https://doi.org/10.1073/pnas.1800357115>

707 Papalexiou, S. M., & Montanari, A. (2019). Global and Regional Increase of Precipitation Extremes  
708 Under Global Warming. *Water Resources Research*, 55(6), 4901–4914.  
709 <https://doi.org/10.1029/2018WR024067>

710 Papalexiou, S. M., Rajulapati, C. R., Clark, M. P., & Lehner, F. (2020). Robustness of CMIP6  
711 Historical Global Mean Temperature Simulations: Trends, Long-Term Persistence,



712 Autocorrelation, and Distributional Shape. *Earth's Future*, 8(10), e2020EF001667.  
 713 <https://doi.org/10.1029/2020EF001667>

714 Pascual, F. G., & Akhundjanov, S. B. (2019). Monitoring a bivariate INAR(1) process with application  
 715 to Hepatitis A. *Communications in Statistics - Theory and Methods*, 1–23.  
 716 <https://doi.org/10.1080/03610926.2019.1645856>

717 Pedeli, X., Davison, A. C., & Fokianos, K. (2015). Likelihood Estimation for the INAR(p) Model by  
 718 Saddlepoint Approximation. *Journal of the American Statistical Association*, 110(511), 1229–  
 719 1238. <https://doi.org/10.1080/01621459.2014.983230>

720 Peden, A., Franklin, R., Leggat, P., & Aitken, P. (2017). Causal Pathways of Flood Related River  
 721 Drowning Deaths in Australia. *PLoS Currents*, 1.  
 722 <https://doi.org/10.1371/currents.dis.001072490b201118f0f689c0fbe7d437>

723 Rosenzweig, C., Tubiello, F., Goldberg, R., Mills, E., & Bloomfield, J. (2004). Increased crop damage  
 724 in the US from excess precipitation under climate change. *Global Environmental Change*, 12,  
 725 197–202. [https://doi.org/10.1016/S0959-3780\(02\)00008-0](https://doi.org/10.1016/S0959-3780(02)00008-0)

726 Serinaldi, F., & Kilsby, C. G. (2016). The importance of prewhitening in change point analysis under  
 727 persistence. *Stochastic Environmental Research and Risk Assessment*, 30(2), 763–777.  
 728 <https://doi.org/10.1007/s00477-015-1041-5>

729 Serinaldi, F., Kilsby, C. G., & Lombardo, F. (2018). Untenable nonstationarity: An assessment of the  
 730 fitness for purpose of trend tests in hydrology. *Advances in Water Resources*, 111, 132–155.  
 731 <https://doi.org/10.1016/j.advwatres.2017.10.015>

732 Smith, A. (2021). *2020 U.S. Billion-Dollar Weather and Climate Disasters—In Historical Context*.  
 733 <https://doi.org/10.13140/RG.2.2.25871.00166/1>

734 Steutel, F. W., & van Harn, K. (1979). Discrete Analogues of Self-Decomposability and Stability. *The*  
 735 *Annals of Probability*, 7(5), 893–899.

736 von Storch, H. (1999). Misuses of Statistical Analysis in Climate Research. In H. von Storch & A.  
737 Navarra (Eds.), *Analysis of Climate Variability* (pp. 11–26). Berlin, Heidelberg: Springer Berlin  
738 Heidelberg.

739 Sun, F., Roderick, M. L., & Farquhar, G. D. (2018). Rainfall statistics, stationarity, and climate change.  
740 *Proceedings of the National Academy of Sciences*, 115(10), 2305.  
741 <https://doi.org/10.1073/pnas.1705349115>

742 Tramblay, Y., El Adlouni, S., & Servat, E. (2013). Trends and variability in extreme precipitation  
743 indices over Maghreb countries. *Natural Hazards and Earth System Sciences Discussions*, 1.  
744 <https://doi.org/10.5194/nhessd-1-3625-2013>

745 Trenberth, K. E. (2011). Attribution of climate variations and trends to human influences and natural  
746 variability. *WIREs Climate Change*, 2(6), 925–930. <https://doi.org/10.1002/wcc.142>

747 Trenberth, K. E., Dai, A., Rasmussen, R. M., & Parsons, D. B. (2003). The Changing Character of  
748 Precipitation. *Bulletin of the American Meteorological Society*, 84(9), 1205–1218.  
749 <https://doi.org/10.1175/BAMS-84-9-1205>

750 Vogel, R. M., Rosner, A., & Kirshen, P. H. (2013). Brief Communication: Likelihood of societal  
751 preparedness for global change: trend detection. *Natural Hazards and Earth System Sciences*,  
752 13(7), 1773–1778. <https://doi.org/10.5194/nhess-13-1773-2013>

753 Weiß, C. H. (2008). Serial dependence and regression of Poisson INARMA models. *Journal of*  
754 *Statistical Planning and Inference*, 138(10), 2975–2990.  
755 <https://doi.org/10.1016/j.jspi.2007.11.009>

756 Westra, S., Alexander, L., & Zwiers, F. (2013). Global Increasing Trends in Annual Maximum Daily  
757 Precipitation. *Journal of Climate*, 26, 7834. <https://doi.org/10.1175/JCLI-D-12-00502.1>

758 Wilks, D. S. (1997). Resampling Hypothesis Tests for Autocorrelated Fields. *Journal of Climate*,  
759 10(1), 65–82. [https://doi.org/10.1175/1520-0442\(1997\)010<0065:RHTFAF>2.0.CO;2](https://doi.org/10.1175/1520-0442(1997)010<0065:RHTFAF>2.0.CO;2)

760 Wilks, D. S. (2006). On “Field Significance” and the False Discovery Rate. *Journal of Applied*  
761 *Meteorology and Climatology*, 45(9), 1181–1189. <https://doi.org/10.1175/JAM2404.1>

762 Wilks, D. S. (2016). “The Stippling Shows Statistically Significant Grid Points”: How Research  
763 Results are Routinely Overstated and Overinterpreted, and What to Do about It. *Bulletin of the*  
764 *American Meteorological Society*, 97, 160309141232001. [https://doi.org/10.1175/BAMS-D-15-](https://doi.org/10.1175/BAMS-D-15-00267.1)  
765 00267.1

766 Wilks, D. S. (2019). Chapter 7 - Statistical Forecasting. In *Statistical Methods in the Atmospheric*  
767 *Sciences (Fourth Edition)* (pp. 235–312). Elsevier. [https://doi.org/10.1016/B978-0-12-815823-](https://doi.org/10.1016/B978-0-12-815823-4.00007-9)  
768 4.00007-9

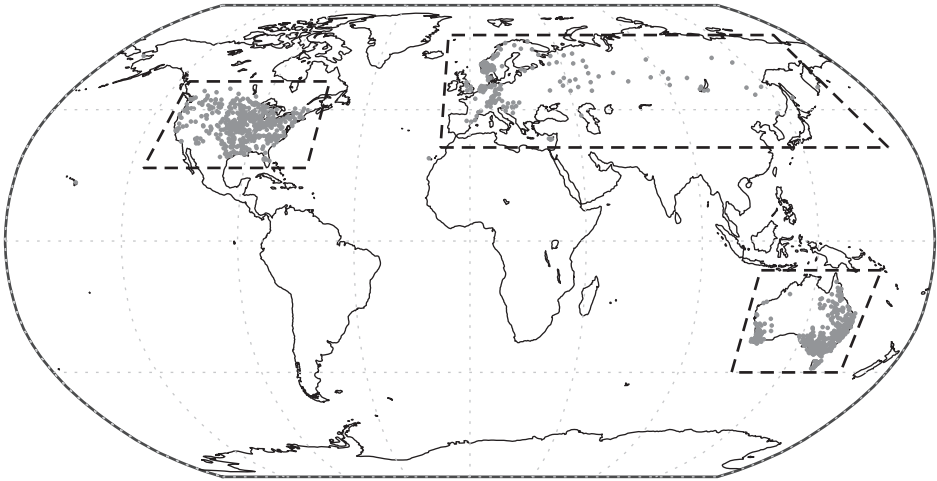
769 Wright, D. B., Bosma, C. D., & Lopez-Cantu, T. (2019). U.S. Hydrologic Design Standards  
770 Insufficient Due to Large Increases in Frequency of Rainfall Extremes. *Geophysical Research*  
771 *Letters*, 46(14), 8144–8153. <https://doi.org/10.1029/2019GL083235>

772 Yue, S., & Wang, C. Y. (2002). Applicability of prewhitening to eliminate the influence of serial  
773 correlation on the Mann-Kendall test. *Water Resources Research*, 38(6), 4–1.  
774 <https://doi.org/10.1029/2001WR000861>

775 Yue, S., Pilon, P., Phinney, B., & Cavadias, G. (2002). The influence of autocorrelation on the ability  
776 to detect trend in hydrological series. *Hydrological Processes*, 16(9), 1807–1829.  
777 <https://doi.org/10.1002/hyp.1095>

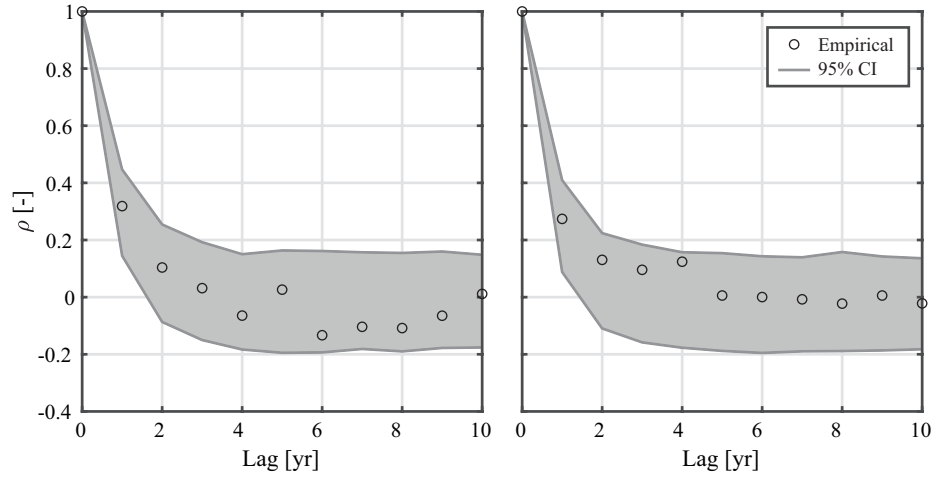
778 Zhang, W., & Villarini, G. (2019). On the weather types that shape the precipitation patterns across the  
779 U.S. Midwest. *Climate Dynamics*, 53(7), 4217–4232. [https://doi.org/10.1007/s00382-019-](https://doi.org/10.1007/s00382-019-04783-4)  
780 04783-4

781 Zolotokrylin, A., & Cherenkova, E. (2017). Seasonal changes in precipitation extremes in russia for the  
782 last several decades and their impact on vital activities of the human population. *Geography,*  
783 *Environment, Sustainability*, 10, 69–82. <https://doi.org/10.24057/2071-9388-2017-10-4-69-82>

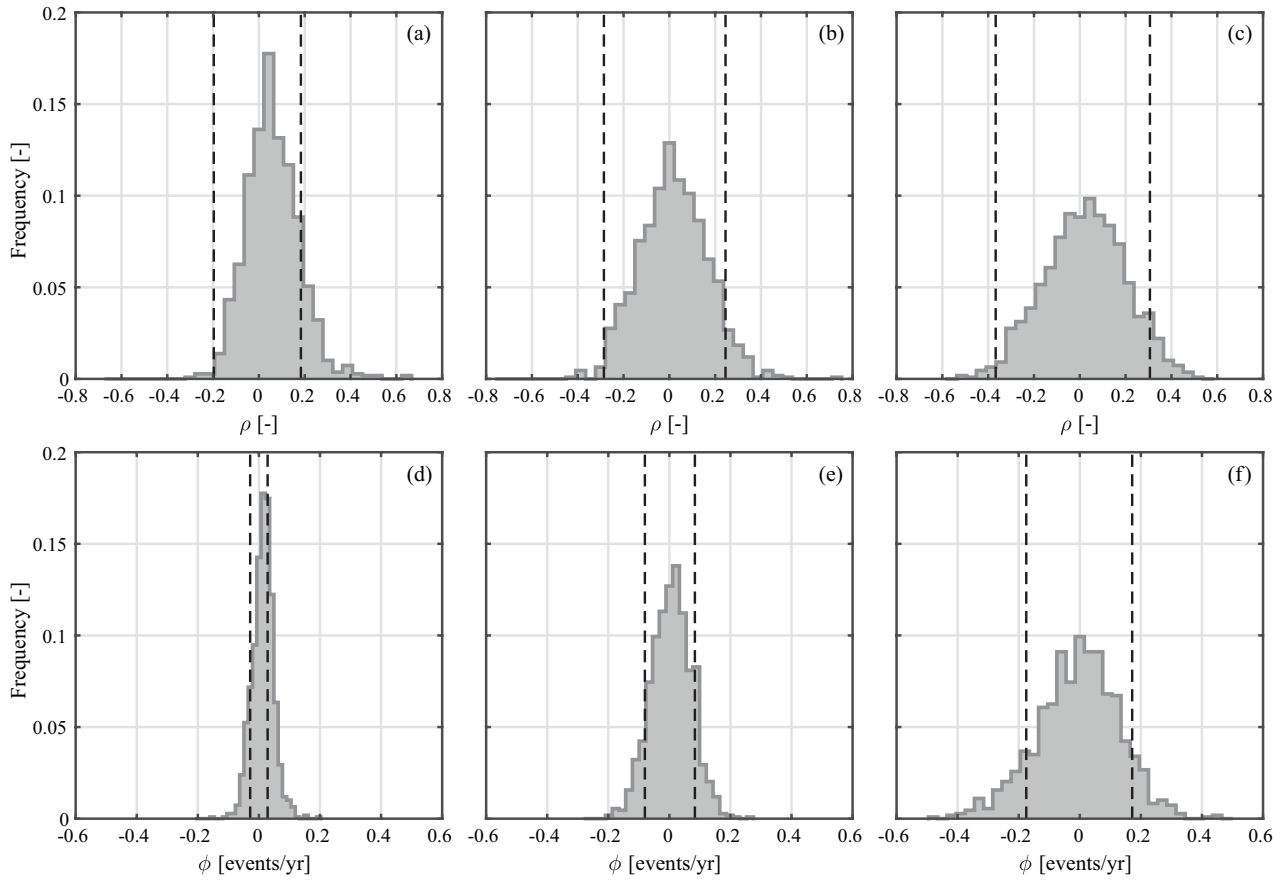


785  
786    **Figure 1.** GHCN rain gauges selected for this study with indication of the three regions of (i) North  
787    America, (ii) Europe and Asia, and (iii) Australia displayed in subsequent figures.

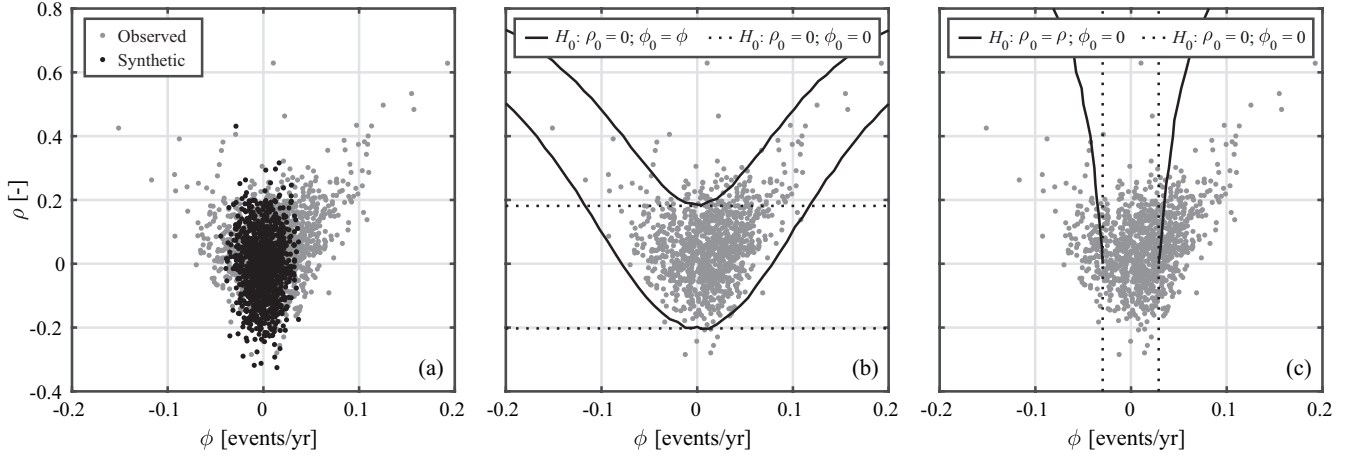
788



789  
790 **Figure 2.** Examples of empirical autocorrelation function of two randomly chosen observed count time  
791 series ( $q = 0.95$ ,  $n = 100$  years) of the GHCN dataset along with 95% confidence interval (CI) derived  
792 from 10,000 synthetic time series generated with the Poisson-INAR(1) model. For both series, the slope  
793 of the linear trend is smaller than 0.02 events/yr.

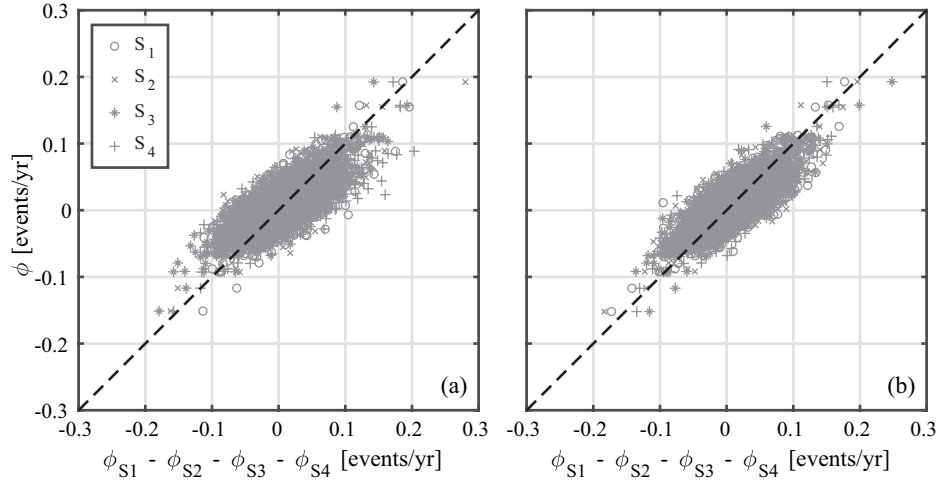


794  
 795 **Figure 3.** Histograms of (a)-(c) lag-1 autocorrelation  $\rho$  and (d)-(f) linear trend slope  $\phi$  estimated on the  
 796  $M = 1087$  observed count time series for  $q = 0.95$  and sample size  $n = 100, 50$ , and  $30$  years (from left  
 797 to right). Vertical lines depict the 95% confidence intervals obtained from 10,000 synthetic uncorrelated  
 798 and stationary time series ( $H_0: \rho_0 = 0; \phi_0 = 0$ ).



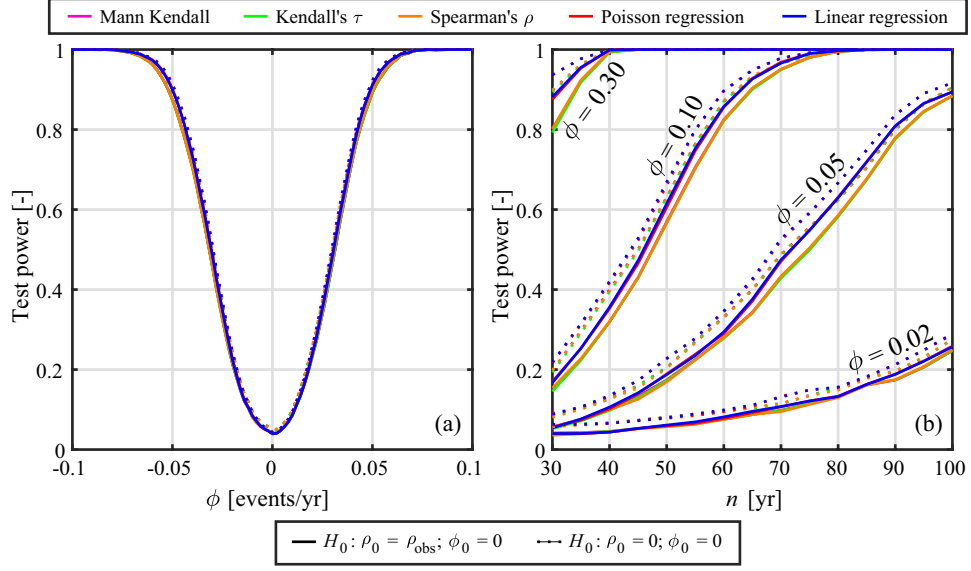
799

800 **Figure 4.** Scatterplot between  $\phi$  and  $\rho$  computed on  $M = 1087$  observed count time series for  $q = 0.95$   
 801 and  $n = 100$  years (grey circles) along with: (a) scatterplot between  $\phi$  and  $\rho$  calculated on synthetic counts  
 802 with  $H_0: \rho_0 = 0; \phi_0 = 0$  (black circles); (b) 95% CIs of  $\rho$  computed under  $H_0: \rho_0 = 0; \phi_0 = \phi$  (solid  
 803 line) with  $\phi$  being the value in the  $x$ -axis, and  $H_0: \rho_0 = 0; \phi_0 = 0$  (dashed line); (c) 95% CIs of  $\phi$  computed  
 804 under  $H_0: \rho_0 = \rho; \phi_0 = 0$  (solid line) with  $\rho$  being the value in the  $y$ -axis, and  $H_0: \rho_0 = 0; \phi_0 = 0$  (dashed  
 805 line).

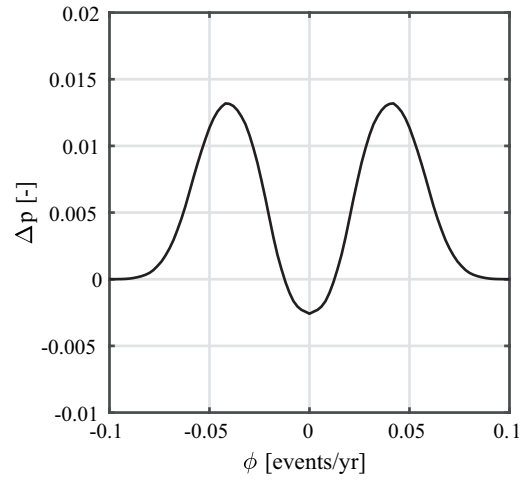


806  
807 **Figure 5.** (a) Scatterplot of linear slopes  $\phi$  estimated on  $M = 1087$  observed count time series for  $q =$   
808  $0.95$  and  $n = 100$  years versus linear slopes  $\phi_{S*}$  ( $* = 1, 2, 3, 4$ ) estimated on the corresponding  $4 \times M$  sub-  
809 series of  $n = 25$  years extracted from each full series by sampling one record every four years and denoted  
810 with S1-S4. (b) Same as (a) but for synthetic time series generated under  $H_0$ : “ $\rho_0 = 0$ ;  $\phi_0 = \phi_{\text{obs}}$ ”, with  $\phi_{\text{obs}}$   
811 being the observed slope.

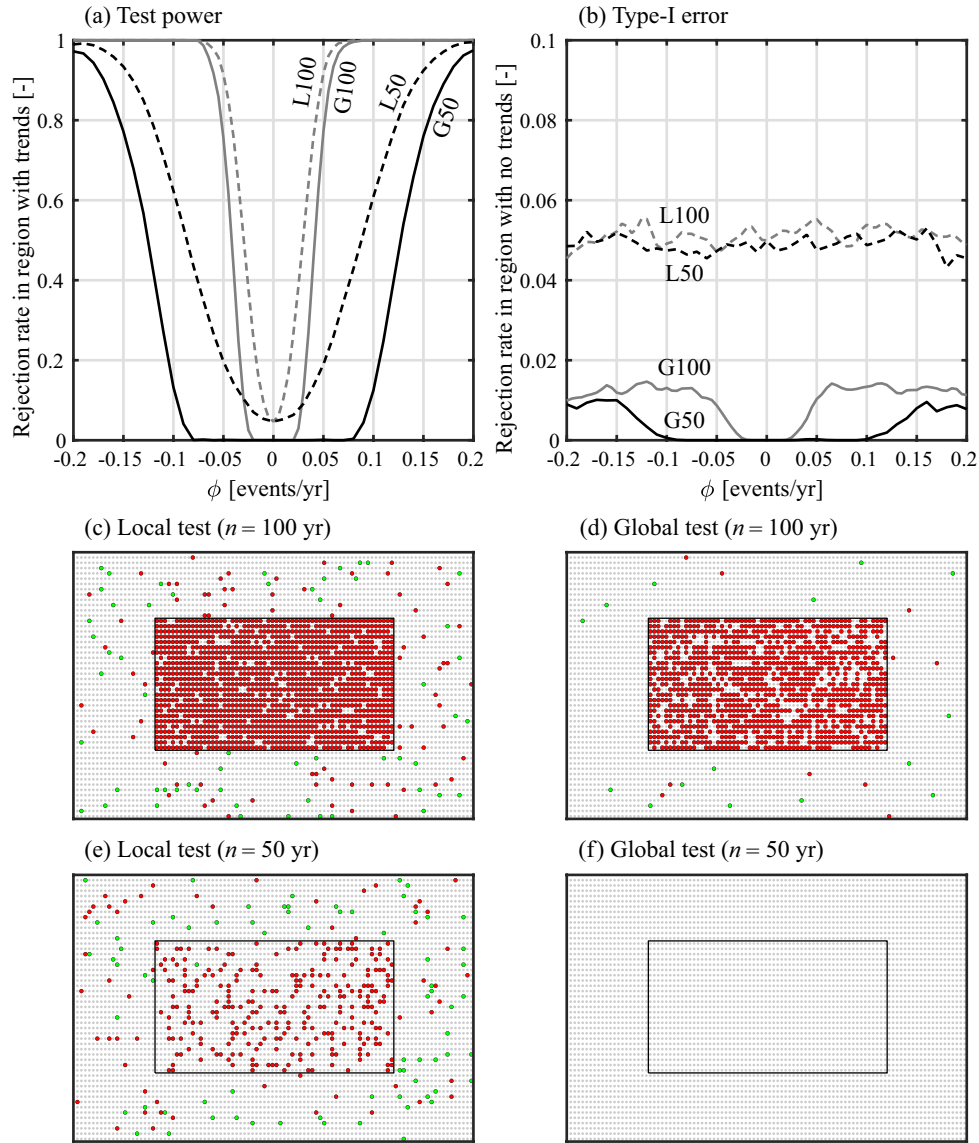




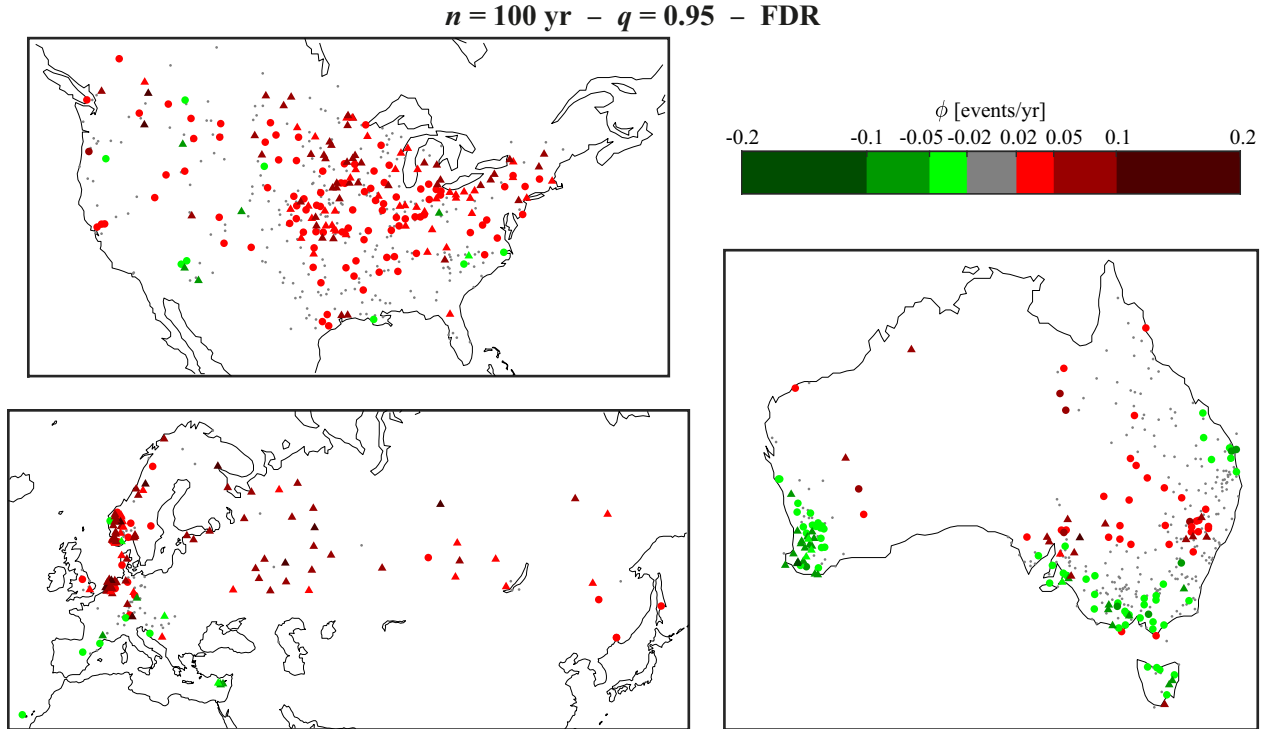
812  
 813 **Figure 6.** Performances of several trend tests with the null distribution of the test statistics built under  
 814  $H_0: \rho_0 = 0; \phi_0 = 0$  (dashed line) and  $H_0: \rho_0 = \rho_{\text{obs}}; \phi_0 = 0$  (solid line), evaluated on synthetic count  
 815 time series relative to  $q = 0.95$ . (a) Power of tests as a function of  $\phi$  for uncorrelated nonstationary time  
 816 series for length  $n = 100$  years. (b) Power of tests as a function of  $n$  for uncorrelated nonstationary time  
 817 series for  $\phi = 0.02, 0.05, 0.10$ , and  $0.30$  events/yr.



818  
819 **Figure 7.** Gaussian-weighted moving average of the differences between power of PR and MK tests  
820 (indicated with  $\Delta p$ ) reported in Fig. 6a for  $q = 0.95$  and  $n = 100$  years.



821  
 822 **Figure 8.** Performance of trend test at multiple sites quantified through a synthetic experiment in a 50 x  
 823 100 grid points (see text for details). (a) Fraction of local (L) and global (G) rejections of  $H_0$  as a function  
 824 of  $\phi$  in the inner region with trend (test power) for  $n = 100$  and 50 years. (b) Same as (a) but for the outer  
 825 region with no trend (type-I error). (c) Map of local rejections of  $H_0$  for the case where an increasing  
 826 trend with slope  $\phi = 0.05$  events/yr is assumed in the inner region and  $n = 100$  years. (d) Same as (c), but  
 827 for global rejections of  $H_0$  after applying the FDR test. (e)-(f) same as (c)-(d), but for  $n = 50$  years. In  
 828 (c)-(f), red (green) dots represent rejections of  $H_0$  with increasing (decreasing) trend, while grey dots are  
 829 used when  $H_0$  is not rejected.



830

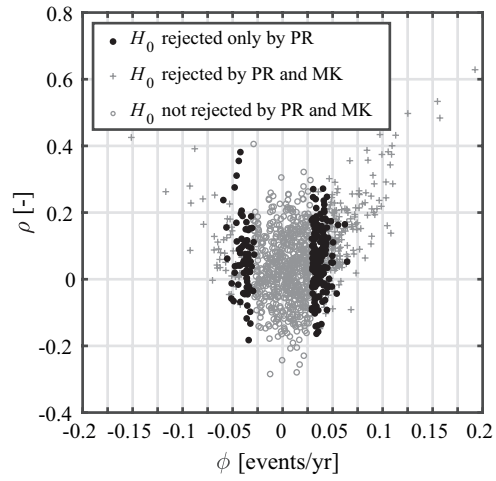
831

832

833

834

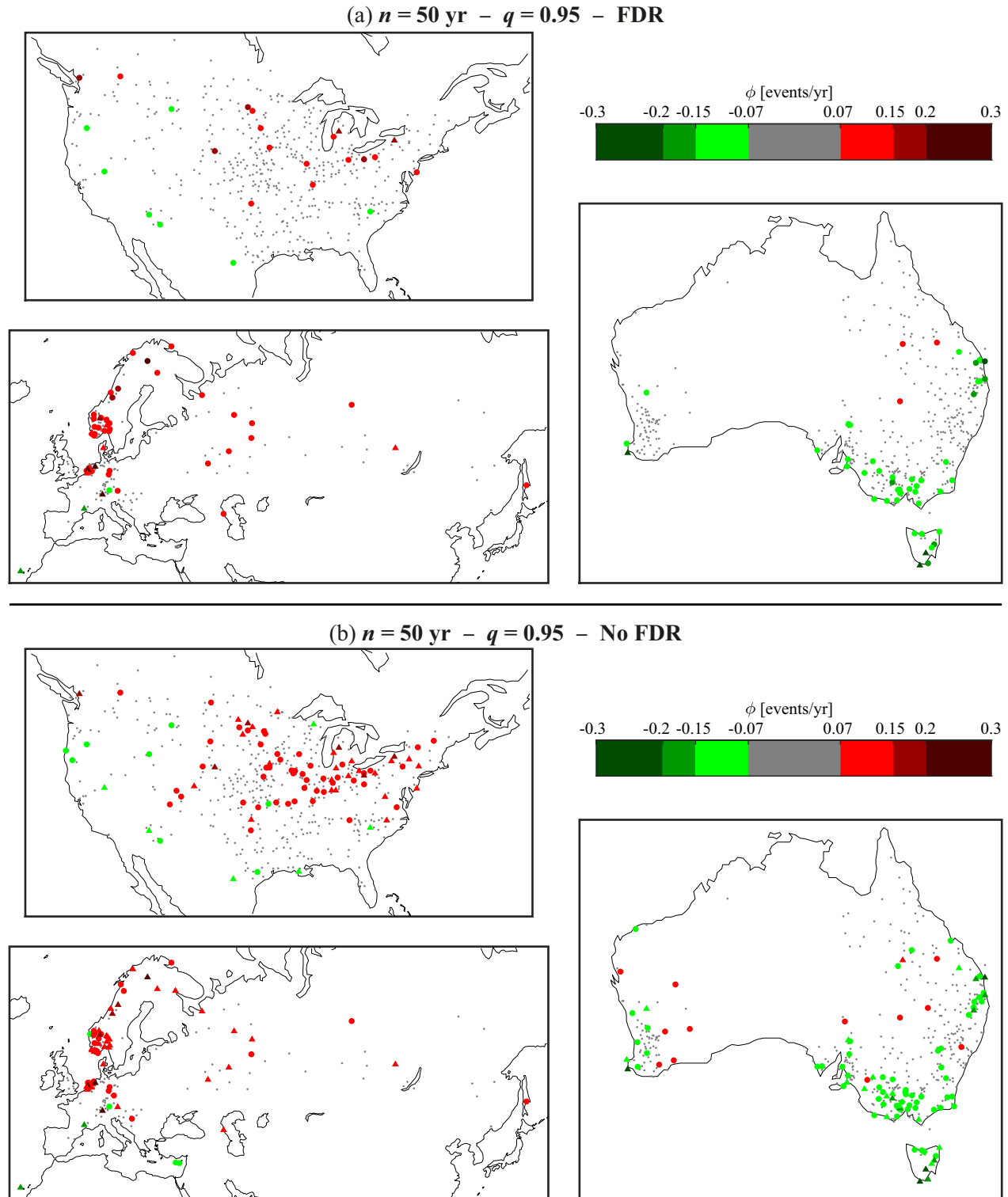
**Figure 9.** Statistically significant trends at the GHCN gages after applying the FDR tests at  $\alpha_{\text{global}} = 0.05$  for  $n = 100$  and  $q = 0.95$ . Larger circles (triangles) are used when  $H_0$  is rejected by PR only (PR and MK), with colors based on the trend slope value and sign. Smaller grey dots are used when  $H_0$  is not rejected by both tests.



835

836 **Figure 10.** Scatterplot between  $\phi$  and  $\rho$  computed on  $M = 1087$  observed count time series for  $q = 0.95$

837 and  $n = 100$  years, with different markers visualizing possible combined outcomes of PR and MK tests.



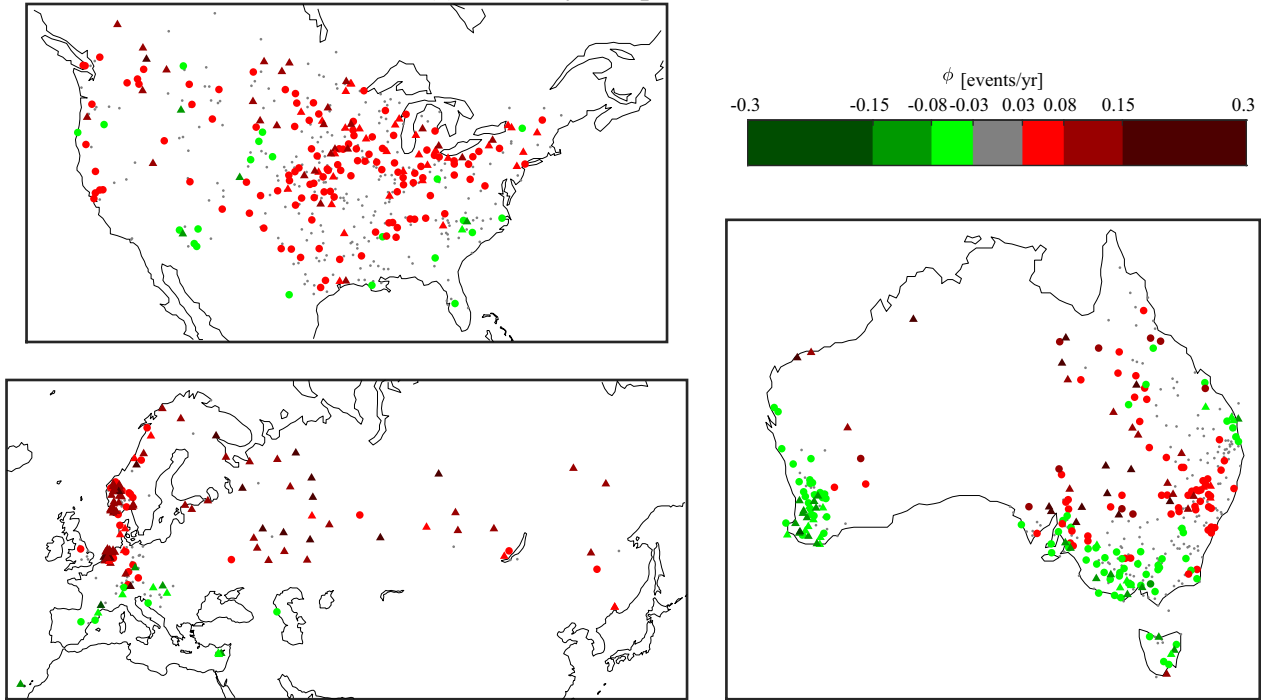
838

839

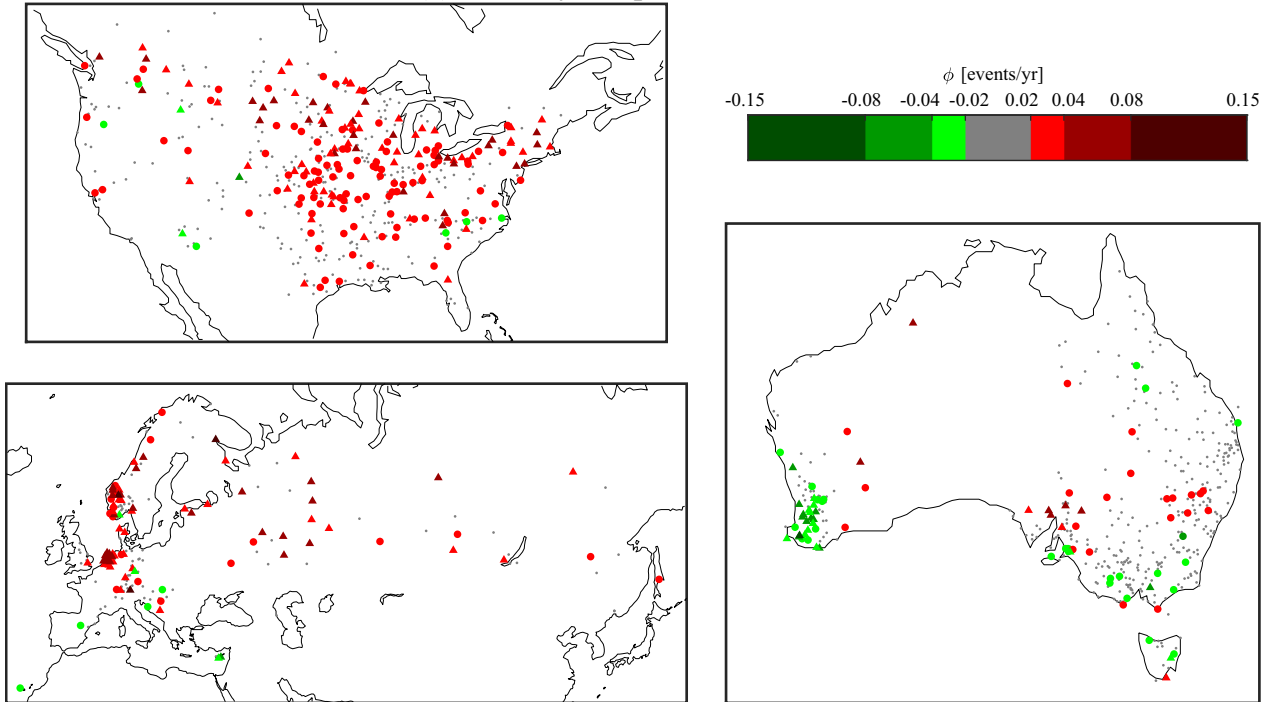
840

**Figure 11.** (a) As in Fig. 9 but for  $n = 50$  yr. (b) As in (a) but for local results without the application of the FDR test.

(a)  $n = 100 \text{ yr} - q = 0.90 - \text{FDR}$



(b)  $n = 100 \text{ yr} - q = 0.975 - \text{FDR}$



841

842 **Figure 12.** As Fig. 9 but for (a)  $q = 0.90$  and (b)  $q = 0.975$ .

Reference	Spatial coverage	Time Period	Number of years	Data	Spatial aggregation	Metrics	Statistical techniques	Account for field significance	Account for serial correlation	Main results
Groisman et al., (2005)	Global	1893-2002	32 to 110	NCDC rain gages	Sub-continental regions	AF <sub>POT</sub> and AT <sub>POT</sub>	SP	No	No significant observed $\rho$ 's	Statistically significant increasing (decreasing) trends in Europe, America, South-Africa (southwestern Australia).
Alexander et al., (2006)	Global	1901-2003	52 to 102	GHCN, ECA and GSN rain gages (5948)	Point and 5° x 5° grids	11 PCCIs	MK	Bootstrapping	Prewhitening	Trends on extreme events mainly increasing, but statistically significant on 13%–37% of the stations depending on the metric.
Westra et al., (2013)	Global	1900-2009	30 to 110	HadEX2 rain gages (8326)	Point	AM	MK	Bootstrapping	Negligible mean $\rho$	- Increasing trends in 2/3 of rain gages, but only 8.6% statistically significant; - No evident spatial patterns of significant trends.
Kunkel & Frankson, (2015)	Global	1951-2014	64	GHCN rain gages (6619)	10° x 10° grids	AF <sub>POT</sub> and AT <sub>POT</sub>	KT	No	Trend test estimating variance inflation	- Most trends not statistically significant. - Increasing trends in most of the world except western North America, southern Europe, northern Eurasia and western and eastern coasts of Australia.
Asadieh & Krakauer, (2015)	Global	1901-2010	30 to 110	HadEX2 rain gages (~11,600); CMIP5 outputs	2.54° x 3.75° grids	AM	MK	No	No	- Increasing (decreasing) trends in 66.2% (33.8%) of grid cells, but only 18% (4%) statistically significant; - Consistent results with CMIP5 outputs but with trend underestimation.
Papalexiou & Montanari, (2019)	Global	1964-2013	50	GHCN rain gages (8730)	Point and 5° x 5° grids	AF <sub>POT</sub> and AM <sub>POT</sub>	MC	No	AF: negligible mean $\rho$ AMM: use of AR(1) in Monte Carlo simulations	- Coherent spatial patterns of trends more evident in frequency (AF <sub>POT</sub> ) than magnitude (AM <sub>POT</sub> ); - For AF <sub>POT</sub> : increasing trends in central and eastern USA, Europe, eastern Russia and most of China; - For AM <sub>POT</sub> : increasing trends in western and northern Europe, and eastern and central USA; - Ratio of increasing/decreasing statistically significant trends: 2.4 (1.3) for AF <sub>POT</sub> (AM <sub>POT</sub> ).
Janssen et al., (2014)	USA	1901-2012	112	HadEX2 rain gages (726) CMIP5 outputs	Sub-continental regions and 1° x 1° grids	AF <sub>POT</sub>	PD	No	No	- Statistically significant increasing (decreasing) trends in central and eastern (western) USA; - Consistent results from CMIP5 but with a trend underestimation.
Hoerling et al., (2016)	USA	1901-2013; 1979-2013.	35 to 114	GHCN rain gages (~10,000)	Sub-continental regions	AT <sub>POT</sub> , AF <sub>POT</sub> and AM <sub>POT</sub>	ND	No	No	- Statistically significant increasing (decreasing) trends in 1901-2013 in the northeastern (southwestern) USA; - Similar patterns for the 1979-2013 period.
Wright et al., (2019)	USA	1950-2017	68	GHCN rain gages (911)	Sub-continental regions	AF <sub>POT</sub>	NB	No	No	Statistically significant increasing trends in eastern USA, smaller and less significant changes in western parts



Kunkel et al., (2020)	USA	- 1949-2016; - 1979-2016.	37 to 68	GHCN rain gages (3098)	Sub-continental regions	AF <sub>POT</sub> and AT	KT	No	No	- In 1949-2016, statistically significant increasing trends in several areas of USA; statistically significant decreasing trends in western USA, but in lower number; - Similar patterns in 1979-2016 but lower number of significant trends.
Kruger & Nxumalo, (2017)	South Africa	1921-2015	95	Rain gages (60)	Point and rainfall districts	11 PI	<i>t</i> -test	No	No	Statistically significant increasing trends in indices related to extreme events in southern and middle regions.
New et al., (2006)	South Africa	1961-2000	30 to 40	Rain gages (63)	Point and regions	10 PI	KT	No	No	- Increasing trends in indices related to extreme events at regional scale; - Few statistically significant trends and no evident spatial patterns at local scale.
Tramblay et al., (2013)	North-Africa	1950-2008	33 to 59	Rain gages (22)	Point	11 PI	MK	FDR test	Trend test accounting for variance inflation	- Few decreasing trends on indices related to extreme events; - Statistically significant trends with local tests, but in much lower number after applying the FDR test.
Hennessy et al., (1999)	Australia	1910-1995	86	Rain gages (379)	Countries	Several PI	KT	No	No	Statistically significant increasing (decreasing) trends on indices focused on extreme events in South Australia and New South Wales (Western Australia) regions.
Hughes, (2003)	Australia	-	-	-	-	Review	-	-	-	Increasing trends in most areas, decreasing trends in southwestern and southeastern regions.
Gallant et al., (2007)	Australia	- 1910-2005; - 1950-2005.	56 to 96	Rain gages (92)	Six regions	Several PI	KT	No	No	Statistically significant increasing (decreasing) trends on indices related to extreme events in central (southwestern and southeastern) regions.
Alpert et al., (2002)	Mediterranean Basin	1951-1995	45	Rain gages (265)	Countries	AT <sub>POT</sub>	SP	FDR test	No	Statistically significant increasing trends in two of the four considered countries.
Madsen et al., (2014)	Europe	-	-	-	-	Review	-	-	-	Overall increase both in frequency and magnitude of extreme precipitation, especially in northern parts.
Zolotokrylin & Cherenkova, (2017)	Russia	1961-2013	53	Rain gages (527)	Point	SF <sub>POT</sub> and ST <sub>POT</sub>	Not defined	No	No	Statistically significant increasing trends in 2/3 of rain gages for all seasons.

843 **Table 1.** Summary of several empirical trend analyses of precipitation extremes performed at global and regional scale with daily records.  
844 Acronyms for datasets, metrics and statistical tests are defined as follows. (1) Datasets: NCDC = National Climatic Data Center; GHCN =  
845 Global Historical Climate Network; ECA = European Climate Assessment; GCN = GCOS (Global Observing System for Climate) Surface  
846 Network; HadEX2 = Hadley Center Global Climate Extremes Index 2; CMIP5 = Coupled Model Intercomparison Project 5. (2) Metrics:  
847 AM = Annual maxima;  $AF_{POT}$  ( $SF_{POT}$ ) = Annual (seasonal) frequencies in peak-over-threshold (POT) series;  $AT_{POT}$  ( $ST_{POT}$ ) = Annual totals  
848 of exceedances in POT series;  $AM_{POT}$  = Annual average magnitude of exceedances in POT series; AT (ST) = Annual (seasonal) totals; PI  
849 = Precipitation-based indices. (3) Statistical tests: MK = Mann-Kendall test; KT = Kendall's  $\tau$  test; SP = Spearman's  $\rho$  test; PD = Poisson's  
850 distribution-based test; NB = Negative binomial regression; MC = Monte Carlo simulations.

	$n = 100$	$n = 50$	$n = 30$
Significant $\rho$ 's for $H_0: " \rho_0 = 0; \phi_0 = 0 "$			
Local test	156 (14%)	77 (7%)	80 (7%)
FDR test	29 (3%)	3 (0%)	0 (0%)
Significant $\phi$ 's for $H_0: " \rho_0 = 0; \phi_0 = 0 "$			
Local test	467 (43%)	244 (22%)	193 (18%)
FDR test	451 (41%)	114 (10%)	94 (9%)

851 **Table 2.** Number and percentage of count series derived for  $q = 0.95$  with significant  $\rho$  and/or  $\phi$  for local  
852 and FDR tests assuming  $H_0: " \rho_0 = 0; \phi_0 = 0 "$ .