





24

## Abstract

25       With the development of refined numerical forecasts, the problems such as the score distortion  
26 due to the division of precipitation thresholds in both traditional and improved scoring methods for  
27 precipitation forecast and the increasing subjective risk arisen from the scale setting of the  
28 neighbourhood spatial verification method have become increasingly prominent. To solve this issue,  
29 a general comprehensive evaluation method (GCEM) has been developed for cross-scale  
30 precipitation forecasts by directly analysing the proximity of precipitation forecasts and observations  
31 in this study. In addition to the core element of the precipitation forecast accuracy score (PAS) index,  
32 the GCEM system also includes score indices for insufficient precipitation forecasts, excessive  
33 precipitation forecasts, precipitation forecast biases and clear/rainy forecasts. The PAS does not  
34 distinguish the magnitude of precipitation and delimit the area of influence, it constitutes a fair  
35 scoring formula with objective performance and can be suitable for evaluating the rainfall events  
36 such as general and extreme precipitation. The PAS can be used to calculate the accuracy of  
37 numerical models or quantitative precipitation forecasts, enabling the quantitative evaluation of the  
38 comprehensive capability of various refined precipitation forecasting products. Based on the GCEM,  
39 comparative experiments between the PAS and TS are conducted for two typical precipitation  
40 weather processes. The results show that relative to TS, the PAS aligns with subjective expectations  
41 much more, indicating that the PAS is more reasonable than the TS. In addition, other indices of the  
42 GCEM are utilized to analyse the range and extent of both insufficient and excessive forecasts of  
43 precipitation, as well as the precipitation forecast ability in two weather processes. These indices not  
44 only provide overall scores for individual cases similar to the TS but also offer two-dimensional  
45 score distribution plots, which can comprehensively reflect the performance and characteristics of



46 precipitation forecasts. Both theoretical and practical applications demonstrate that the GCEM  
47 exhibits distinct advantages and potential promotion and application value compared to the various  
48 mainstream precipitation forecast verification methods.

49

## 50 **1. Introduction**

51 Precipitation is one of the most important forecasting elements in weather forecasting (Bi et al.,  
52 2016; Han al., 2023). Short duration heavy rainfall often lead to flooding and geological disasters,  
53 causing widespread and severe impacts (Zhong et al., 2022; Yang et al., 2023). Precipitation forecasts,  
54 as focuses and challenges in meteorological department operations, have drawn widespread attention  
55 from governments, societies and the public (Bi et al., 2016; Hao et al., 2023). Scientifically  
56 evaluating precipitation forecasts helps people gain a clear understanding of the current precipitation  
57 forecast levels and maintain appropriate psychological expectations for such forecasts. Moreover,  
58 such evaluations assist forecasters in rationally analysing the quality and characteristics of  
59 quantitative precipitation forecast systems and aid researchers in understanding the level, strengths  
60 and weaknesses of various types of forecasting systems, which in turn, offers valuable insights to  
61 improve these systems (Zhong et al., 2022; Zhang et al., 2023; Liu et al., 2022; Gofa et al., 2018).  
62 However, there are several shortcomings in current precipitation verification approaches. For  
63 instance, traditional scoring methods often fail to reflect model performance improvements, and  
64 small errors in the location or timing of convective features can lead to false alarms and missed  
65 events. A model's utility is limited by diagnosing model errors such as a displaced forecast feature or  
66 an incorrect mode of convective organization (Ahijevych et al., 2009). For high-resolution  
67 precipitation forecasts, even if the spatial distribution and intensity of precipitation are consistent



68 with the observations, slight spatial and temporal deviations between forecasts and observations may  
69 still result in a large false alarm ratio and missed alarm ratio, leading to lower forecast scores (Zhao  
70 et al., 2018). With the rapid development of seamless fine quantitative precipitation forecasts, the  
71 need for objective and rational evaluations of the accuracy and characteristics of precipitation  
72 forecasts have become increasingly important and urgent (Chen et al., 2021).

73       Precipitation forecast verification involves various methods, including traditional contingency  
74 table-based classification verification and spatial verification methods. The traditional verification  
75 method can be traced back to 1884 when Finley introduced a dichotomous contingency table for  
76 tornado forecasts and evaluated these forecasts using the proportion correct scoring method (Finley,  
77 1884). Subsequently, systematic attention was given to the evaluation of forecast classification  
78 methods, and Finley's forecast verification method became a classic example of the discussion of  
79 forecast scoring methods (Murphy, 1996). Shortly thereafter, Gilbert (1984) proposed two scoring  
80 methods, namely, the verification and success ratios in forecasting. The verification ratio later  
81 became known as the threat score (TS) (Palmer et al., 1949) or the critical success index (Donaldson  
82 et al., 1975; Mason, 1989). In forecasting, the success ratio is referred to as the Gilbert skill score  
83 (GSS) (Schaefer, 1990) or the equitable threat score (ETS) (Doswell et al., 1990; Gandin et al., 1992).  
84 The TS encourages correct event forecasts (hits) and accounts for the impacts on the false alarm and  
85 missed alarm ratios, which can better guide forecasters or research and development personnel in  
86 making reasonable subjective and objective predictions compared to relying solely on simple  
87 "accuracy". Meanwhile, the ETS eliminates the influence of random forecasts on the score, resulting  
88 in a fairer skill score (Liu et al., 2022).

89       For rare-event forecast verification, in addition to the TS and ETS, the methods such as the



90 Peirce skill score (PSS) (Peirce, 1884; Hanssen et al., 1965; Murphy et al., 1985; Flueck, 1987) and  
91 the Heidke skill score (HSS) (Doolittle, 1885; Doolittle, 1888; Heidke, 1926) can be used. The PSS  
92 is a fair score index that is equal to the hit rate minus the false detection probability; the HSS  
93 eliminates the influence of random forecasts, and the results can reflect the forecast skill (Liu et al.,  
94 2022). Many studies have reviewed and compared these two scoring methods (Doswell et al., 1990;  
95 Schaefer, 1990; Marzban, 1998; Mason, 2003). In extreme weather event verification (including  
96 severe convective weather such as short duration heavy rainfall), the traditional scoring methods  
97 (such as the TS and ETS) for dichotomous events often yield scores of zero when the occurrence  
98 probability of the object being verified is very low. Therefore, Stephenson proposed the extreme  
99 dependency score (EDS) for evaluating extreme events. The EDS has the advantage that different  
100 forecast systems converge to different values and has no explicit dependence on the bias of the  
101 prediction system (Stephenson et al., 2008; Casati et al., 2008).

102 It has been more than a century since Gilbert proposed the two scoring concepts, i.e., the  
103 verification and success ratios in forecasting (later known as the TS and ETS). The TS and ETS have  
104 been widely used for the performance evaluation of threshold-based event forecasts despite their  
105 evident shortcomings (Stephenson et al., 2008). Today, in various forecast verification applications,  
106 including high-resolution quantitative precipitation and extreme weather forecast verification, the TS  
107 and ETS remain mainstream approaches (Tang et al., 2017; Wei et al., 2019; Chen et al., 2021; Liu et  
108 al., 2023). Of course, with the continuous introduction of new scoring methods, some of the  
109 problems in traditional verification have been solved. However, the advantageous position of the TS  
110 remains unchallenged. The reasons for which, although varied, are worthy of attention, but include  
111 its objectivity and practicality.



112           The traditional TS categorizes precipitation according to thresholds and performs verification  
113 using a dichotomous contingency table. The TS can be viewed as a measure of the forecast accuracy  
114 that excludes the hit forecasts for “non-occurrence” precipitation events (referred to as no  
115 precipitation), and its calculation formula being simple, objective and standardized. However, there  
116 are two main limitations of the TS. First, precipitation is categorized by thresholds based on the  
117 contingency table, which has limitations in terms of classification. The drawback of artificially  
118 dividing precipitation into different threshold ranges is that it may not guarantee that two adjacent  
119 precipitation values fall within the same threshold range and slightly different precipitation values  
120 may not be within the same threshold, which can lead to precipitation score distortion. The second  
121 limitation is related to the so-called "double penalty" issue. With the development of high-resolution  
122 numerical weather forecasting and the shortening of the spacing between model grid points, some  
123 medium- and small-scale phenomena have been portrayed by models. However, it is difficult for  
124 high-resolution numerical forecasts to match the characteristics of the observed medium- and  
125 small-scale forecasts. In addition, traditional scoring methods often cannot reflect these  
126 improvements in terms of the model performance. Assuming a constant forecast area, when there is a  
127 small deviation in the timing and location of events between a forecast and an observation, be both  
128 "false alarms" and "missed alarms" will occur and is referred to as the "double penalty" phenomenon.  
129 This phenomenon leads to a score lower than the subjective expected result, making it difficult to  
130 obtain appropriate verification scores when a forecast that “looks good” is not as good as one that  
131 “looks bad” (Ahijevych et al., 2009; Wilks, 2006; Ebert, 2008; Chen et al., 2021). For  
132 low-probability events with limited sample size for verification, such as torrential rain and short-term  
133 heavy rainfall, the “double penalty” issue becomes more prominent. The TS and ETS for torrential



134 rain tends to align closely the unskillful portion of the scoring values (Chen et al., 2019). In recent  
135 years, new mainstream scoring methods have mainly addressed the abovementioned limitations but  
136 still have shortcomings. Such methods include the improved gradient decreasing method, which still  
137 results in poor scores for good forecasts, and the neighbourhood spatial verification method, which  
138 has too many subjective components and may miss medium and small scale information.

139 To address the limitations of threshold-based precipitation classification, and improve the  
140 verification effect, e.g., the gradient decreasing method (hereinafter referred to as the magnitude  
141 improved TS) is used to verify the accuracy of rainstorm forecasts (Yang Dong et al., 2017), and  
142 appropriate weights are assigned to close forecast values to avoid scores of zero (Table 1). However,  
143 the magnitude-improved TS still has limitations. For example, if the observed 24-hour accumulated  
144 precipitation is 50 mm, when forecast A is 48 mm while forecast B is 98 mm, it is evident that  
145 forecast A is better than forecast B. According to the original TS, forecast B scores 1 point, while  
146 forecast A does not score any points. For the magnitude-improved TS, forecast B scores 1 point, as it  
147 still falls within the same magnitude category as the observed precipitation, while forecast A only  
148 scores 0.4 points, which still fails to reflect the fact that forecast A is superior to forecast B (Table 2).  
149 By employing the new scoring method, i.e., the precipitation forecast accuracy score (PAS), which  
150 will be discussed later, forecast A scores 0.998 points, while forecast B scores 0.398 points,  
151 confirming the rationality and validity of this new method.

152 To address the “double penalty” issue, a common approach is to employ the neighbourhood  
153 spatial verification method (also known as the fuzzy method), which has two specific processing  
154 forms. The first form is simple upscaling, which uses a certain method (including value averaging,  
155 maximum, value weighting, etc.) to select values within the scale range, adjusting the high-resolution



156 forecast and observation information to a larger scale to reduce the accidental information of  
157 high-resolution data, and then using the traditional skill score. Although this upscaling method  
158 provides traditional skill scores at different scales (Yates et al., 2006; Weygandt et al., 2004), it  
159 cannot address the issue of excessive smoothness of the precipitation fields during the upscaling  
160 process (Zhao et al., 2018), which may result in the omission of some small- to medium-scale  
161 information (Zepeda et al., 2000). The other form is the improved neighbourhood spatial verification  
162 method proposed by Roberts and Lean (2008). By referring to the Murphy skill score, this method  
163 obtains comprehensive evaluation information by comparing the occurrence frequency (probability)  
164 of precipitation within different scale windows. If the forecasted occurrence frequency closely  
165 approximates the observed occurrence frequency, the forecast is considered value (Zhao et al., 2018).  
166 From the perspective of the precipitation occurrence probability within the analysis region, the  
167 precipitation occurrence probability for observations and forecasts is the ratio of the precipitation  
168 area to the analysed area of the region, which is referred to as the fraction skill score (FSS). This  
169 method also effectively improves the “double penalty” problem.

170       The neighbourhood spatial verification method considers values that are spatially and  
171 temporally adjacent between forecasts and observations during the matching process, thus relaxing  
172 the strict requirements for spatiotemporal matching (Ebert, 2008; Casati et al., 2008). However, since  
173 the determination of the neighbourhood range is a rather subjective process, it hinders the  
174 standardization of verification scores and lacks comparability, which may negatively affect objective  
175 quantitative verification. Numerous experiments have shown that there is an obvious improvement in  
176 the scoring values after adopting the neighbourhood spatial verification method (Chen et al., 2019),  
177 particularly for forecasts of large-magnitude precipitation. Nevertheless, the purpose of scoring is not





178 to achieve a monotonous increase in scoring values but rather to follow the principle of objectivity as  
179 much as possible. Errors are errors and cannot be solved by simply lowering the standard. Instead,  
180 reasonable and fair criteria should be utilized to reflect the true extent of errors.

181 Currently, numerical weather forecasts and intelligent gridded forecasts have been developed to  
182 output high-resolution precipitation products, while precipitation observations, whether in the form  
183 of gridded or station data, are already high-resolution. Staying at the dichotomous classification level  
184 for precipitation verification not only wastes existing data resources but also fails to meet the  
185 evaluation requirements of refined forecasts. Therefore, to adapt to the development of refined  
186 forecasts, a new scoring method is needed. In light of this, a comprehensive verification index for  
187 precipitation forecasts is designed, and the following five aspects are considered. (1) The impact of  
188 categorical events on rainstorm forecasts should be reduced. In particular, high-resolution forecasts  
189 can refer to continuous variables for scoring methods. (2) The design of the scoring method should  
190 aim to minimize subjective factors such as the artificial range division and condition settings,  
191 ensuring the scoring objectivity and comparability. (3) The designed scoring performance indices  
192 should possess ideal attributes such as fairness, base rate independence, suitability for extreme events,  
193 and boundedness as much as possible. (4) The devised scoring method should be easy to promote,  
194 concise and efficient, with clear concepts and scientific rationality. (5) Different comprehensive  
195 verification indices for precipitation forecasts should reflect the forecast performance and  
196 characteristics of high-resolution quantitative precipitation from various perspectives.

197 In this study, on the basis of analysing the limitations of traditional verification methods as well  
198 as improved methods, a new general comprehensive evaluation method (GCEM) for cross-scale  
199 precipitation prediction is proposed. This method is applied and verified through practical examples.



200 The remainder of this paper is organized as follows. Section 2 provides an overview of various  
201 scoring indices and their attributes in the GCEM and introduces the optimization processing method  
202 for the PAS index in the application. Through ideal experiments, the characteristics of the scoring  
203 methods are analysed based on the score curves described in Section 3. Section 4 presents the  
204 comparative experiments conducted between the new scoring method and the traditional scoring  
205 method based on typical cases. Finally, the summary and discussion are given in Section 5.

## 206 **2 Cross-scale general comprehensive evaluation method**

### 207 **2.1 General comprehensive evaluation method overview**

208 To address the issues of “distorted scores due to the division of precipitation thresholds and  
209 increased subjective risks brought about by the setting of the neighbourhood spatial verification  
210 method” in traditional and improved precipitation scoring methods, referring to the verification  
211 method for heavy rainfall forecasts based on predictability (Chen et al., 2019) and combining the  
212 advantages of the relative and absolute errors in this study, a GCEM is constructed by directly  
213 analysing the proximity of forecasted precipitation to observed precipitation. It primarily includes the  
214 PAS, and the expression of its core scoring function as follows.

$$215 \quad \text{PAS} = \begin{cases} \sin\left(\frac{\pi}{2} \cdot \frac{x}{u}\right), & 0 \leq x < u \\ e^{-\left(\frac{x-u}{u}\right)^2}, & 0 < u \leq x \end{cases}, (1)$$

216 where PAS represents the scoring value,  $x$  is the forecasted precipitation (mm), and  $u$  is the  
217 observed precipitation (mm). The PAS falls between 0 and 1, where a higher score indicates a better  
218 precipitation forecast effect. When  $\text{PAS} = 1$ , it signifies a perfect forecast, indicating that the  
219 forecasted and observed precipitation match entirely. When the observed precipitation is not  
220 forecasted,  $\text{PAS} = 0$ . When the forecasted precipitation amount is sufficiently large,  $\text{PAS} \rightarrow 0$  (Fig. 1).



221 Additionally, considering the characteristic large fluctuations in the function curve when the  
222 observed precipitation is less than 10 mm, a smoothing optimization is applied to Eq. (1) (see Section  
223 2.2 for details).

224 The GCEM system also includes the following indices:

225 (1) Accuracy score of the insufficient precipitation forecast (IPS), whose core scoring function  
226 expression is as follows.

$$227 \quad \text{IPS} = \sin\left(\frac{\pi}{2} \cdot \frac{x}{u}\right) - 1, \quad 0 \leq x < u \quad , \quad (2)$$

228 where IPS represents the scoring value, reflecting the degree of underestimation in precipitation  
229 forecasts when the forecasted value is less than the observed value. The IPS falls between  $[-1, 0)$ ,  
230 where a scoring value closer to 0 indicates a lower degree of underestimation.

231 (2) Accuracy score of the excessive precipitation forecast (EPS), whose core scoring function  
232 expression is as follows.

$$233 \quad \text{EPS} = 1 - e^{-\left(\frac{x-u}{u}\right)^2}, \quad 0 < u < x \quad , \quad (3)$$

234 where EPS represents the scoring value, reflecting the degree of overestimation in precipitation  
235 forecasts when the forecasted value exceeds the observed value. The EPS falls between  $(0, 1]$ , where  
236 a scoring value closer to 0 indicates a lower degree of overestimation.

237 (3) Accuracy score of the insufficient and excessive precipitation forecast (IEPS), whose core  
238 scoring function expression is as follows.

$$239 \quad \text{IEPS} = \begin{cases} \sin\left(\frac{\pi}{2} \cdot \frac{x}{u}\right) - 1, & 0 \leq x < u \\ 0, & x = u \\ 1 - e^{-\left(\frac{x-u}{u}\right)^2}, & 0 < u < x \end{cases} \quad , \quad (4)$$

240 where IEPS represents the scoring value, reflecting the degree of deviation of the forecasted



241 precipitation from the observed precipitation. The IEPS falls between  $[-1, 1]$ , where a scoring value  
242 closer to 0 indicates a lower degree of deviation. An IEPS less (greater) than 0 indicates an  
243 insufficient (excessive) forecast, and an IEPS equal to 0 represents an unbiased forecast.

244 (4) The PAS clear/rainy forecast accuracy score (PASC), whose scoring function expression is  
245 as follows.

$$246 \quad \text{PASC} = \frac{1}{m+n} (\sum_m \text{PAS}_{|ux0.1} + \sum_n \text{PASN}) \quad , (5)$$

247 where PASN denotes the score of the correctly forecasted non-precipitation event, and its formula is  
248 given as:

$$249 \quad \text{PASN} = \begin{cases} 1, & 0 \leq u < 0.1 \text{ and } 0 \leq x < 0.1 \\ 0, & u \geq 0.1 \text{ or } x \geq 0.1 \end{cases} . \text{ PASC represents the PAS scoring value for}$$

250 clear/rainy forecasts, and  $\text{PAS}_{|ux0.1}$  denotes the overall PAS for precipitation forecasts under  
251 specific conditions where the observed precipitation  $u \geq 0.1$  mm or the forecasted precipitation  $x \geq$   
252  $0.1$  mm.  $m$  is the number of stations or grid points for  $\text{PAS}_{|ux0.1}$ , and  $n$  is the number of stations or  
253 grid points where non-precipitation events are correctly forecasted.

254 The discussion below pertains to the characteristics of the PAS scoring method. As an ideal  
255 performance indicator, the PAS has the attributes of boundedness, fairness, sensitivity disparity,  
256 suitability for extreme events and moderate symmetry.

257 (1) Boundedness. The PAS scoring values range between 0 and 1. A PAS score of 1 represents  
258 an ideal forecast, while a score of 0 indicates that there is observed precipitation but no forecasted  
259 precipitation or that the forecasted precipitation is sufficiently large. The scoring range is consistent  
260 with that of traditional TS, making it easy to compare and evaluate the scoring methods and suitable  
261 for practical forecast verification applications.

262 (2) Fairness. The PAS scoring method constitutes a scoring formula in an objective form



263 without a subjective boundary definition. Precipitation forecasts are verified without magnitude or  
264 delimitation of the area of influence, and the closer to the observed situation the forecast is, the  
265 higher the score, which is fair.

266 (3) Sensitivity disparity. From the Chinese national standard GB/T 28592—2012 “Grade of  
267 precipitation” on the classification of precipitation grades, the public is more sensitive to low-grade  
268 precipitation forecasts. As rainfall intensity increases, the public's sensitivity gradually decreases;  
269 that is, the public has a higher tolerance for errors in response to heavier rainfall forecasts. In other  
270 words, large errors in the forecasts of heavy rainfall events may be considered equivalent to smaller  
271 errors in weaker rainfall events in terms of forecast scoring. As shown in Fig. 1, the intersection point  
272 on the PAS scoring curves for the observed precipitation amounts of 25 mm and 100 mm  
273 corresponds to a forecasted amount of 42.4 mm. That is, the forecast errors are 17.4 mm and 57.6  
274 mm for the observed 24-hour accumulated precipitation of 25 mm and 100 mm, respectively, while  
275 the scores are both 0.62. From the perspective of forecast service effectiveness, this aligns with the  
276 general public perception.

277 (4) Suitability for extreme events. From the PAS scoring curves for forecasts corresponding to  
278 different observed precipitation amounts ( $u = 10, 25, 50$  and  $100$  mm) (Fig. 1), it is evident that the  
279 PAS scoring method performs well in evaluating precipitation events forecasts at the level of  
280 torrential rain and above. For example, when the observed precipitation is 100 mm, with forecasted  
281 amounts of 59 mm and 147.2 mm, the PASs are both 0.8, whereas the TSs are 0 and 1, and the  
282 improved TSs are 0.8 and 1, respectively. This result indicates that the PAS is suitable for scoring  
283 heavy rainfall events, meeting the general applicability requirements as a scoring method that does  
284 not degrade due to extreme events.



285 (5) Moderate symmetry. In Eq. (1), observed precipitation is the independent variable  $x$ , and  
286 forecasted precipitation is the parameter  $u$ . Then, the equation is rewritten as:

$$287 \quad \text{PAS} = \begin{cases} e^{-\left(\frac{u-x}{x}\right)^2}, & 0 < x \leq u \\ \sin\left(\frac{\pi}{2} \cdot \frac{u}{x}\right), & u \leq x \end{cases}, \quad (6)$$

288 Similarly, for different magnitudes of forecasted precipitation ( $u = 10, 25, 50$  and  $100$  mm) and  
289 observed precipitation ranging from  $0$  to  $300$  mm, the corresponding forecast scores are shown in Fig.  
290 2. The forecast scores also vary with the degree of proximity between forecasts and observations.  
291 Figs. 1 and 2 exhibit similar trends but are not identical, illustrating that the PAS possesses moderate  
292 symmetry.

## 293 **2.2 PAS verification for precipitation forecasts**

294 From the properties of the core verification function of the PAS, it is noted that when the  
295 observed precipitation  $u < 10$  mm, there is a large gradient in the PAS curve. A slight change in the  
296 forecasted value ( $x$ ) can result in a large fluctuation in the PAS. To account for this characteristic,  
297 based on a comprehensive analysis in combination with the sensitivity of forecasters and the public  
298 to small-scale precipitation, a smoothing optimization scheme is applied to the PAS curve for  
299 accumulated precipitation below  $10$  mm. Similarly, the IPS, EPS, IEPS and PASC curves are  
300 appropriately smoothed and optimized according to their respective definitions.

301 Let forecasted precipitation  $x = 0$  mm,  $1 > \text{PAS} > 0$  for observed precipitation, and  $0 < u < 10$   
302 mm.

303 Assumption:

304 (1)  $\text{PAS} = 0.6\text{PAS}|_{u \rightarrow 0}$  when  $u = 0$  mm, and  $x \neq 0$  mm;

305  $\text{PAS}|_{u \rightarrow 0}$  denotes the PAS for the case of observed precipitation  $0 < u \leq 0.1$  mm;



306 (2)  $PAS = 0.6PAS|_{u \rightarrow 0}$  when  $x = 0$  mm, and  $u < 10$  mm;

307  $PAS|_{x \rightarrow 0}$  denotes the PAS for the case of forecasted precipitation  $0 < x \leq 0.1$  mm.

308 1. When the observed precipitation  $u = 0$  mm (Fig. 3a), let  $PAS = 0.6PAS|_{u \rightarrow 0}$ . Then,

$$309 \quad PAS = 0.6e^{-\left(\frac{x}{10}\right)^2} \quad x > 0 \quad (7)$$

310 2. When the forecasted precipitation  $x = 0$  mm and the observed precipitation  $u < 10$  mm

311 (Fig. 3b), let  $PAS = 0.6PAS|_{x \rightarrow 0}$ . Then,

$$312 \quad PAS = 0.6 \sin\left(\frac{\pi}{2} \cdot \frac{10-u}{10}\right), \quad 0 < u < 10 \quad (8)$$

313 3. When the observed precipitation  $0 < u < 10$  mm and the forecasted precipitation  $x \neq 0$  (Fig.

314 3c). Then,

$$315 \quad PAS = \begin{cases} \sin\left(\frac{\pi}{2} \cdot \frac{x-u+10}{10}\right), & 0 \leq x < u, \quad 0 \leq u < 10 \\ e^{-\left(\frac{x-u}{10}\right)^2}, & u \leq x, \quad 0 \leq u < 10 \end{cases} \quad (9)$$

316 Note that  $x$  and  $u$  are not both equal to 0 at the same time.

317 4. When the observed precipitation  $u \geq 10$  mm (Fig. 3d), then

$$318 \quad PAS = \begin{cases} \sin\left(\frac{\pi}{2} \cdot \frac{x}{u}\right), & 0 \leq x < u, \quad u \geq 10 \\ e^{-\left(\frac{x-u}{u}\right)^2}, & u \leq x, \quad u \geq 10 \end{cases} \quad (10)$$

319

320 To compare with the traditional scoring method, the new scoring method for precipitation

321 forecasting adopts the “classification before verification, no classification during verification”

322 approach. Scoring for precipitation processes over different accumulated periods is based on but not

323 limited to the commonly used precipitation classification approaches in practical operations, as

324 shown in Tables 3–5.

### 325 **3 Ideal experimental validation of the new verification method**



### 326 **3.1 Validation of forecast scoring results for general precipitation**

327 General precipitation refers to precipitation ranging from light rain to heavy rain, i.e., 24-hour  
328 accumulated precipitation between [0.1 mm, 50 mm). Figure 4 shows the schematic diagram of PAS  
329 scores for general precipitation. The forecasted amounts are compared under conditions when the  
330 24-hour accumulated precipitation is 10 mm, 25 mm and 45 mm and the PAS scores are 0.8, 0.7, 0.5  
331 and 0.3 (Table 6). When the observed precipitation is 10 mm, the forecasted amounts of 5.9 mm and  
332 14.7 mm both have a PAS score of 0.8, with differences from the perfect forecast value (10 mm) of  
333 4.1 mm and 4.7 mm, respectively; the forecasted amounts with a PAS score of 0.3 are 1.9 and 21.0  
334 mm, differing by 8.1 mm and 11.0 mm from the perfect forecast value (10 mm), respectively. When  
335 the observed precipitation is 25 mm, the forecasted amounts with a PAS score of 0.8 are 14.7 mm  
336 and 36.8 mm, with differences from the perfect forecast value (25 mm) of 10.3 mm and 11.8 mm,  
337 respectively; the forecasts with a PAS score of 0.5 are 8.3 mm and 45.8 mm, differing by 16.7 mm  
338 and 20.8 mm from the perfect forecast value (25 mm), respectively.

339 For forecasts with the same observed precipitation and the same scores, the absolute errors of an  
340 insufficient forecast and observation are smaller than those of an excessive forecast and observation,  
341 and the higher the scores are, the closer the absolute errors of the forecasts. When the observed  
342 precipitation is 50 mm, only the insufficient precipitation forecast is scored since a precipitation  
343 forecast exceeding 50 mm is not considered within the scope of general precipitation evaluation. The  
344 scoring experimental results align with expectations.

### 345 **3.2 Validation of forecast scoring results for precipitation at the level of torrential** 346 **rain and above**

347 Figure 5 shows a schematic diagram of the PASs when the amount of precipitation exceeds the





348 storm magnitude. The predicted precipitation is compared when the 24-hour cumulative observed  
349 precipitation is 25 mm, 50 mm, and 100 mm with PAS scores of 0.877, 0.7, 0.5, 0.3, and 0.1 (Table  
350 7). When the observed precipitation is 25 mm, forecasts  $\geq 50$  mm are only involved in the rating,  
351 with PASs of 0.3 and 0.1 for forecasts of 52.4 and 62.9 mm, respectively.

352 When the PAS is 0.877 and the observed precipitation is 50 mm, the predicted values are 34.1  
353 and 68.1 mm, respectively; when the observed precipitation is 100 mm, the predicted values are 68.1  
354 and 136.2 mm, respectively. When the observed precipitation is 50 or 100 mm, the prediction is 68.1  
355 mm, with a score of 0.877. The absolute error is 18.1 mm for the excessive precipitation forecast and  
356 31.9 mm for the insufficient precipitation forecast. This result indicates that the scoring tolerance  
357 increases as the grade of observed precipitation increases and is gradually expanded through  
358 continuous changes, avoiding discontinuous increases caused by magnitude changes.

359 When the observed precipitation is 50 mm and the PAS is 0.3, the insufficient forecast is 9.7  
360 mm and the excessive forecast is 104.9 mm. When the observed precipitation is 100 mm, the  
361 predictions for a PAS of 0.3 are 19.4 and 209.7 mm, respectively. When the observed precipitation is  
362 50 mm, the insufficient forecast with a PAS of 0.1 is 3.2 mm, and the excessive forecast is 125.9 mm.  
363 When the observed precipitation is 100 mm, the predictions with a PAS of 0.1 are 6.1 and 251.7 mm,  
364 respectively.

365 Under constant observed precipitation conditions, for forecasts with the same score, the absolute  
366 error between the insufficient forecast and the observed precipitation is smaller than that between the  
367 excessive forecast and the observed precipitation. The higher the score is, the smaller the absolute  
368 error between the forecast and the observation. Moreover, the scoring tolerance increases with  
369 increasing observed precipitation. The scoring experimental results conform to expectations.



## 370 **4 Example-based comparative experiments for the new verification method**

### 371 **4.1 Introduction of typical cases**

372 Comparative experiments of traditional TS and PAS are conducted for 12-hour accumulated  
373 precipitation of two typical cases. One case pertains to the precipitation weather process occurring at  
374 12:00 UTC on July 16, 2019 (referred to as “Case 1”), which is dominated by a weak weather system.  
375 The other case relates to the precipitation weather process occurring at 12:00 UTC on June 13, 2020  
376 (referred to as “Case 2”), which is predominantly associated with a strong weather system.

377 Both of these precipitation cases are associated with precipitation during the Meiyu period. Case  
378 1, which occurred during the Meiyu period of 2019 and was featured by scattered precipitation under  
379 weak synoptic-scale forcing. The low intensity shear line system is located south of the Yangtze  
380 River. There are two precipitation concentration areas, one at the intersection of Hunan Province and  
381 Jiangxi Province and the other covering the majority of Zhejiang Province. The precipitation process  
382 in Case 2 (July 11–13) represents the first round of widespread rainstorms during the Meiyu period  
383 of 2020, including heavy precipitation affected by a low-level vortex shear system. The western  
384 section of the low-level vortex shear is relatively stable, while the eastern section slightly presses  
385 southwards. Southwesterly airflow developed and pushed northwards, and a strong wind speed belt  
386 persisted for a long time in the Jianghuai region. Moreover, the Jiangnan–Jianghuai region  
387 maintained a high-energy and high-moisture state, resulting in persistent heavy rainfall.

388 A subjective analysis of these two weather processes reveals that for the event on July 16, 2019  
389 (Fig. 6), the forecasted precipitation intensity and rainfall areas are relatively consistent with the  
390 observations. There are two distinct heavy rainfall areas in the east and south parts of the Yangtze  
391 River, with particularly high accuracy in forecasting scattered rainstorms in Zhejiang Province



392 located in the eastern section. In contrast, for the precipitation weather process on June 13, 2020 (Fig.  
393 7), it is evident that there is an overestimation of the precipitation forecast.

#### 394 **4.2 Data and methods**

395 The observed precipitation data are provided by the China Meteorological Administration  
396 multisource merged precipitation analysis system (CMPAS), developed by the National  
397 Meteorological Information Center of China. The CMPAS integrates hourly precipitation data from  
398 nearly 40,000 automatic meteorological stations in China and provides radar-based quantitative  
399 precipitation estimation and satellite-retrieved precipitation products with a spatial resolution of  $0.05^\circ$   
400  $\times 0.05^\circ$ . The predicted precipitation data with 3 km resolution are from the Precision Weather  
401 Analysis and Forecasting System (PWAFS) model, a regional refined forecast model, developed by  
402 the Jiangsu Provincial Meteorological Bureau. These data are output once per hour.

403 The specific methods are as follows.

404 (1) Determine the verification domain and verification points. The verification domain covers  
405 the Huang–Huai region of China ( $28^\circ\text{N}$ – $38^\circ\text{N}$ ,  $111^\circ\text{E}$ – $123^\circ\text{E}$ ). The verification points are defined  
406 based on the grid points of the observed precipitation data, their spatial resolution is  $0.05^\circ \times 0.05^\circ$ ,  
407 and the total number of verification grid points is 48,000 ( $200 \times 240$ ).

408 (2) Prepare the observed and forecasted precipitation data and interpolate the forecasted  
409 precipitation data onto the observed grid points. The observed 12-hour accumulated precipitation  
410 data are derived by accumulating the hourly precipitation data from the CMPAS. The forecasted  
411 12-hour accumulated precipitation data are obtained by subtracting the zero-field data from the  
412 12-hour forecast field data. Since the grid points of the observed and forecasted precipitation data do  
413 not coincide and the grid spacing is small, the nearest neighbour method is used in this study to



414 match the forecasted data to the grid points of the observed precipitation. Specifically, the forecasted  
415 data on the grid point nearest to the observed grid point are used as the forecasted value for this  
416 observed grid point.

417 (3) Analyse the relationship between forecasted precipitation and observed precipitation. The  
418 scores for each verification grid point and the overall scores for each verification area are calculated  
419 based on the scoring formula for each index in the GCEM system. Then, the verification result file is  
420 generated in NetCDF format. On this basis, the distribution maps for the scores of various indices in  
421 the GCEM system are produced. Additionally, the total TS and clear/rainy TS for different  
422 precipitation magnitudes within the verification area are calculated based on the TS and clear/rainy  
423 TS formulas.

#### 424 **4.3 Analysis of the comparative experiment results**

425 For the precipitation process on July 16, 2019, traditional TSs for different rainfall categories,  
426 such as clear/rainy and 12-hour accumulated precipitation  $\geq 0.1$  mm,  $\geq 10$  mm,  $\geq 25$  mm and  $\geq 50$   
427 mm, are all lower than the traditional scores for the weather process on June 13, 2020. For example,  
428 the TS is 0.381 for 12-hour accumulated precipitation  $\geq 0.1$  mm at 12:00 UTC on July 16, 2019  
429 (Table 8), while this score is 0.625 for that at 12:00 UTC on June 13, 2020 (Table 9), which differs  
430 from the subjective judgement.

431 For the precipitation process at 12:00 UTC on July 16, 2019, the PASs for clear/rainy and  
432 12-hour accumulated precipitation  $\geq 0.1$  mm,  $\geq 10$  mm and  $\geq 25$  mm are all higher than those for the  
433 precipitation process on June 13, 2020. For instance, the overall PAS is 0.617 for 12-hour  
434 accumulated precipitation  $\geq 0.1$  mm at 12:00 UTC on July 16, 2019. This PAS is higher than the PAS  
435 of 0.457 for the precipitation process at 12:00 UTC on June 13, 2020, which aligns with the



436 subjective judgement.

437 For the precipitation process at 12:00 UTC on July 16, 2019, the PAS for each magnitude is  
438 higher than the corresponding TS, addressing the issue of TSs being lower. For the precipitation  
439 process at 12:00 UTC on June 13, 2020, the PASs for clear/rainy and the magnitudes of  $\geq 0.1$  mm  
440 and  $\geq 10$  mm are lower than the corresponding TSs, whereas the PASs for the magnitudes of  $\geq 25$   
441 mm and  $\geq 50$  mm are higher than the corresponding TSs. This result indicates that the PAS is  
442 different from the magnitude-improved TS and the neighbourhood spatial verification method, both  
443 of which blindly increase the tolerance, leading to a monotonous increase in scores. This result also  
444 demonstrates that the PAS has good discrimination ability for extreme events. The PAS assigns  
445 objective scores based on the proximity of the forecast to the observation, making it more reliable for  
446 precipitation evaluation than the TS.

#### 447 **4.4 Analysis of the indices in the new verification method**

448 The GCEM includes not only the core element of the PAS index but also the IPS, EPS, IEPS  
449 and PASC indices.

450 Regarding the issue of analysing the sources of errors from the verification results, objectively  
451 tracing these errors back from a single score can only determine whether an error was “insufficient  
452 (missed alarm)” or “excessive (false alarm)”. However, the advantage of the GCEM lies in its ability  
453 to decompose the score to each verification point and examine the forecast performance at each point,  
454 which is different from the dichotomous evaluation with only 0 and 1. These indices not only provide  
455 overall scores for individual cases similar to the TS but also offer two-dimensional score distribution  
456 plots, which can comprehensively reflect the performance and characteristics of precipitation  
457 forecasts.



458 Figure 8 shows the distribution for the PASC scores of 12-hour accumulated precipitation. In  
459 these two cases, due to the high accuracy of non-precipitation forecasts, the overall PASC scores are  
460 relatively high. However, for Case 1, the scores in Zhejiang are lower and scattered within a small  
461 area. In contrast, for Case 2, there is a large area occupying most of the Jianghuai region with low  
462 scores. Therefore, the PASC score of Case 1 (0.808) is higher than that of Case 2 (0.734).

463 Figure 9 shows the PAS distributions of 12-hour accumulated precipitation with magnitudes of  
464  $\geq 0.1$  mm,  $\geq 10$  mm and  $\geq 25$  mm. The blank points in the figure are the points that are excluded in  
465 the scoring, following the scoring principle of “classification before verification, no classification  
466 during verification” described in Section 2. From the PAS distributions of different magnitudes, for  
467 Case 1, the scoring areas in Zhejiang exhibit alternatively distributed high and low scores. In contrast,  
468 for Case 2, the scoring areas in the Jianghuai region have a larger area of low scores than high scores.  
469 Therefore, Case 1 has higher PASs for the three categories ( $\geq 0.1$  mm,  $\geq 10$  mm and  $\geq 25$  mm) than  
470 Case 2, and the distributions also allow for distinguishing the areas with better and poorer forecasting  
471 performance.

472 Figure 10 shows the IPS, EPS and IEPS distributions of 12-hour accumulated precipitation. In  
473 terms of the IPS, for Case 1, the large-value IPS areas are located at the intersection of Anhui,  
474 Zhejiang and Jiangxi in the Hunan–Jiangxi region, as well as in the southern part of Hebei. For Case  
475 2, the large-value IPS areas are situated along the Yangtze River in Anhui and Jiangxi, as well as at  
476 the intersection of Henan and Shanxi. The IPSs for Case 1 and Case 2 are  $-0.376$  and  $-0.400$ ,  
477 respectively, indicating that Case 2 shows a slightly higher level of insufficient forecasts (Table 10).  
478 In terms of the EPS, for Case 1, the large-value EPS areas are in Zhejiang and Jiangxi. In contrast,  
479 for Case 2, the large-value EPS areas are located in most of Hunan, Hubei, Anhui and Jiangsu,



480 exhibiting a wide southwest–northeast orientation with a large area and depth. The EPS for Case 2 is  
481 larger than that for Case 1. The IEPS score is a comprehensive reflection of under- and over-  
482 precipitation, and its value reflects the degree of insufficient and excessive precipitation forecasts.  
483 From the distributions of insufficient and excessive precipitation forecasts in Case 1, it is evident that  
484 the insufficient and excessive forecasts are roughly equivalent, each with an IEPS of 0.057. However,  
485 for Case 2, the distribution of the excessive forecast is obviously larger than that of the insufficient  
486 forecast, with an IEPS of 0.325. This result indicates that Case 2 has poorer forecasting performance,  
487 with larger excessive forecasts being an important factor.

488 Consequently, analysing the locations of insufficient and excessive precipitation forecasts from  
489 the figures in conjunction with the characteristics of the forecasting process can provide useful  
490 insights for improving forecasts.

## 491 **5 Discussion and conclusion**

492 By analysing the advantages and disadvantages of the traditional TS, magnitude-improved TS  
493 and neighbourhood spatial verification method, a new precipitation verification method GCEM was  
494 designed and constructed from the perspective of the proximity of the forecast to the observation.  
495 This method consists of the core element of the PAS, as well as multiple elements such as IPS, EPS,  
496 IEPS and PASC.

497 The PAS index consists of sine and e-exponential functions. Additionally, considering the  
498 characteristics of large fluctuations in the function curves when observed precipitation is less than 10  
499 mm, the formula has been smoothed for optimization. The PAS method adopts the principle of  
500 “classification before verification, no classification during verification”, which can serve as an  
501 alternative to skill scores such as the TS and ETS for verifying quantitative precipitation forecasts.



502 This method is characterized by objective and transparent rules and easy generalization. Moreover,  
503 this possesses attributes of an ideal precipitation scoring method, such as fairness, boundedness and  
504 moderate symmetry. Therefore, it can be used to calculate the accuracy of numerical models or  
505 quantitative precipitation forecasts, as well as evaluate the comprehensive forecasting capabilities of  
506 various refined quantitative precipitation forecast products. The GCEM can also evaluate the  
507 performance of numerical forecasts on clear/rain forecasts, as well as insufficient precipitation  
508 forecasts, excessive precipitation forecasts and precipitation forecast biases. In addition to the overall  
509 score, two-dimensional score distribution maps can be generated for each index in the GCEM system.  
510 These maps offer a comprehensive reflection of the precipitation forecasting performance of the  
511 numerical models and serve as a reference for improving model forecasts.

512 This new verification method is validated based on the forecast scoring results for general  
513 precipitation and precipitation at the level of torrential rain and above, and the verification results  
514 align with expectations. Comparative experiments are also conducted on two cases using the new  
515 verification method. For Case 1, the subjective judgement is relatively good, but the TS is lower.  
516 Conversely, for Case 2, the subjective judgement is poorer, yet the TS is higher. Verification using  
517 the PAS reveals that forecasts with better subjective judgement receive higher scores, and forecasts  
518 with poorer subjective judgement receive lower scores. Therefore, the PAS aligns with public  
519 expectations.

520 In addition, the National Meteorological Center of China conducted long-term series large-scale  
521 sample testing on this method in 2023. Based on the ECMWF model's 24-hour and 48-hour  
522 precipitation forecasts from March 2022 to February 2023, the assessment results show that  
523 compared with the TS, the PAS is less affected by the randomness of the sample, and the relative size





524 relationship of different time forecast scores is more stable.

525 From the construction of the GCEM to ideal experiments and case analyses, it is evident that  
526 this evaluation system, especially the PAS method, is a suitable method for quantitative precipitation  
527 evaluation. However, the PAS still has subjective flaws, such as the determination of coefficients in  
528 the PAS expression [0.6 in Eqs. (7) and (8)] when the observed or forecasted precipitation is 0 mm.  
529 Once these coefficients are determined, they apply to all precipitation scoring, thus becoming an  
530 objective component in practice.

531 **Code and data availability.** The data are provided at <https://www.doi.org/10.5281/zenodo.10251028>.

532 The source code of this work can be found at <https://www.doi.org/10.5281/zenodo.10251028>.

533 **Author contributions.** BZ designed the evaluation method, completed the experiments, and wrote  
534 the paper. MZ provided advice on the planning and application of the evaluation method. AH  
535 provided suggestions for the evaluation method and contributed to paper revisions, ZQ contributed to  
536 paper revisions, and CL provided long-term series large-scale sample comparison test results for the  
537 evaluation method. All authors discussed the results and commented on the paper.

538 **Competing interests.** The contact author has declared that none of the authors has any competing  
539 interests.

540 **Disclaimer.** Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional  
541 claims made in the text, published maps, institutional affiliations, or any other geographical  
542 representation in this paper. While Copernicus Publications makes every effort to include appropriate  
543 place names, the final responsibility lies with the authors.

544 **Acknowledgements.** This study was supported by grid precipitation analysis data, forecasted  
545 precipitation data and numerical computing capabilities provided by the Jiangsu Provincial



546 Meteorological Bureau of China, as well as long-term series large-scale sample testing conducted at  
547 the National Meteorological Center of China.

548 **Financial support.** Financial support. This research has been supported by the National Key  
549 Research and Development Program of China (Grant no. 2021YFC3000904).

550

## 551 **References**

552

553 Ahijevych D, Gilleland E, Barbara G B, et al.: Application of spatial verification methods to idealized and  
554 NWP-gridded precipitation forecasts, *Wea Forecasting*, 24:1485-1497, 2009.

555 Bi B G, Dai K, Wang Y, et al.: Advances in techniques of quantitative precipitation forecast, *J Appl Meteor*  
556 *Sci*,27(5):534-549 (in Chinese) doi: 10.11898/1001-7313.20160503, 2016.

557 Casati D B, Wilson L J, Stephenson D B, et al.: Forecast verification: current status and future directions,  
558 *Meteorological Applications*, 15(1):3-18, 2008.

559 Chen F., Chen J., Wei Q., Li J., Liu C., Yang D., Zhao B., Zhang Z.: A new verification method for heavy rainfall  
560 forecast based on predictability II: Verification method and test[J], *Acta Meteorologica Sinica*, 77(1): 28-42.  
561 doi: 10.11676/qxxb2019.003, 2019.

562 Chen H M, Li P X, and Zhao Y.: A review and outlook of verification and evaluation of precipitation forecast at  
563 convection-permitting resolution, *Adv Meteor Sci Technol*,11(3):155-164 (in Chinese) doi:  
564 10.3969/j.issn.2095-1973.2021.03.018, 2021.

565 Donaldson R J, Dyer R M, Krauss M J.: An objective evaluator of techniques for predicting severe weather events,  
566 9th Conf Severe Local Storms, Norman, Oklahoma, Amer Meteor Soc, 321-326, 1975.

567 Doswell CA III, Davies-Jones R, and Keller DL.: On summary measures of skill in rare event forecasting based on  
568 contingency tables, *Weather and Forecasting* 5: 576–585, 1990.

569 Doolittle M H.: Association ratios, *Bull Philos Soc Washington*, 10:83-87;94-96, 1888.

570 Doolittle M H.: The verification of predictions, *Amer Meteor J*, 2:327-329, 1885.

571 Ebert E E.: Fuzzy verification of high-resolution gridded forecasts: a review and proposed framework[J], *Meteor*  
572 *Appl*,15(1):51-64, 2008.

573 Finley J P. :Tomado predictions, *Amer Meteor J*, 1: 85-88, 1884.



- 574 Flueck J A.: A study of some measures of forecast verification, 10th Conf Probability and Statistics in Atmospheric  
575 Science, Edmonton, Alberta, Amer Meteor Soc, 69-73, 1987.
- 576 Gandin L S, Murphy A H.: Equitable scores for categorical forecasts, Monthly Weather Review 120: 361–370,  
577 1992.
- 578 Gilbert G K.: Finley's tornado predictions, Amer Meteor J, 1: 166-172, 1884.
- 579 Gofa F, Boucouvala D, Louka P, et al.: Spatial verification approaches as a tool to evaluate the performance of high  
580 resolution precipitation forecasts[J]. Atmos Res, 208:78-87, 2018.
- 581 Han F., Tang W., Zhou C., Sheng J., and Zhang X.: Improving a precipitation nowcasting algorithm based on the  
582 SWAN system and related application assessment. Acta Meteorologica Sinica , 81(2):304-315 doi:  
583 10.11676/qxxb2023.20220066, 2023.
- 584 Hanssen A W, Kuipers W J A.: On the relationship between the frequency of rain and various meteorological  
585 parameters, Mededeelingen en Verhandelingen, 81:2-15, 1965.
- 586 Hao C., Yu B., Dai Y., Zhi X., And Zhang Y.: Statistical downscaling research on spatio-temporal distributions of  
587 summer precipitation across the beijing region[J], Meteor Mon, 49(7):843-854, 2023.
- 588 Heidke P.: Berechnung des Erfolges und der Güte der Windstärkevorhersagen im Sturmwarnungsdienst (Calculation  
589 of the success and goodness of strong wind forecasts in the storm warning service), Geogr Ann Stockholm, 8:  
590 301-349, 1926.
- 591 Liu C., Lin J., Dai K., Cao Y., Wei Q.: An evaluation method suitable for precipitation forecasts and services[J],  
592 Torrential Rain and Disasters, 41(6): 712-719. doi: 10.12406/byzh.2021-203, 2022.
- 593 Liu C, Dai K, Lin J, et al.: Design and implementation of whole process evaluation program library of weather  
594 forecast[J], Meteor Mon, 49(3):351-364(in Chinese), 2023.
- 595 Liu J., Ren C., Zhao Z., Chen C., Wang Y., And Cai K.: Comparative analysis on verification of heavy rainfall  
596 forecasts in different regional models[J]. Meteor Mon, 48(10):1292-1302, 2022.
- 597 Marzban C.: Scalar measures of performance in rare-event situations. Weather and Forecasting 13: 753–763, 1998.
- 598 Mason IB.: Binary events. in forecast verification: a practitioner's guide in atmospheric science, Jolliffe IT,  
599 Stephenson DB (eds). John Wiley and Sons: Chichester, UK; 37–76, 2003.
- 600 Mason IB.: Dependence of the critical success index on sample climate and threshold probability, Australian  
601 Meteorological Magazine 37: 75–81, 1989.
- 602 Murphy A H.: The finley affair: a signal event in the history of forecast verification, Wea Forecasting, 11:



- 603 3-20,1996.
- 604 Murphy A H, Daan H.: Forecast evaluation // Murphy A H, Katz R W. Probability, statistics, and decision making in  
605 the atmospheric sciences. Westview : Westview Press, 379-437, 1985.
- 606 Palmer W C, Allen R A.: Note on the accuracy of forecasts concerning the rain problem, US Weather  
607 Bureau,1-4,1949.
- 608 Peirce C S.:The numerical measure of the success of prediction, Science, 4:453-454, 1884.
- 609 Roberts N M, Lean H W.: Scale-selective verification of rainfall accumulations from high-resolution forecasts of  
610 convective events[J], Mon Wea Rev, 136:78-97, 2008.
- 611 Schaefer JT.: The critical success index as an indicator of warning skill, Weather and Forecasting 5: 570–575, 1990.
- 612 Stephenson D B, Casati B, Ferro C A T, et al.: The extreme dependency score: A non-vanishing measure for  
613 forecasts of rare events. Meteor Appl, 15(1): 41-50, 2008.
- 614 Tang W., Zhou Q., Liu X., et al.: Analysis on verification of national severe convective weather categorical  
615 forecasts, Meteorological Monthly, 43(1): 67-76, DOI: 10.7519/j.issn.1000-0526.2017.01.007, 2017.
- 616 Wei Q., Li W., Peng S., et al.: Development and application of national verification system in CMA, J Appl Meteor  
617 Sci, 30(2): 245-256. DOI: 10.11898/1001-7313.20190211, 2019.
- 618 Weygandt S S, Loughe A F, Benjamin S G ,et al.: Scale sensitivities in model precipitation skill scores during  
619 IHOP[C]//22nd Conf Severe Local Storms, Amer Met Soc, Hyannis, MA, 4-8 October 2004.
- 620 Wilks, D. S.: Statistical methods in the atmospheric sciences. 2nd ed. Elsevier, 627 pp, 2006.
- 621 Yang Y., Yin J., Wang D., Liu Y., Lu Y., Zhang W., and Xu S.: ABM-based emergency evacuation modelling during  
622 urban pluvial floods: A “7.20” pluvial flood event study in Zhengzhou, Henan Province. Science China Earth  
623 Sciences, 66(2): 282–291, <https://doi.org/10.1007/s11430-022-1015-6>, 2023.
- 624 Yang D., Gao X., and Zhang W.: Research and improvement of a new rainstorm forecast accuracy verification  
625 scheme [C]//Proceedings of the 34th Annual Meeting of the Chinese Meteorological Society S1 Disaster  
626 Weather Monitoring, Analysis and Forecast, 549-559,2017.
- 627 Yates E., S Anquetin, V Ducrocq, et al.: Point and areal validation of forecast precipitation fields[J]. Meteorol Appl,  
628 13: 1-20, 2006.
- 629 Zepeda-Arce J, E Foufoula-Georgiou, Droegemeier K K.: Space-time rainfall organization and its role in validating  
630 quantitative precipitation forecasts[J], J Geophys Res, 105(D8), 10 129-10 146,2000.
- 631 Zhang, B., Zeng, M., Huang, A., Qin, Z., Liu, C., Shi, W., Li, X., Zhu, K., Gu, C., and Zhou, J.: A general



632 comprehensive evaluation method for cross-scale precipitation forecasts, Zenodo [code and data set],  
633 <https://www.doi.org/10.5281/zenodo.10251028>, 2023.

634 Zhang Y., Yu H., Zhang M., Yang Y., and Meng Z.: Uncertainties and error growth in forecasting the  
635 record-breaking rainfall in Zhengzhou, Henan on 19–20 July 2021. *Science China Earth Sciences*, 65(10):  
636 1903–1920, <https://doi.org/10.1007/s11430-022-9991-4>, 2022.

637 Zhao B., Zhang B.: Application of neighborhood spatial verification method on precipitation evaluation[J],  
638 *Torrential Rain and Disasters*,37(1):1-7, 2018.

639 Zhong Q, Sun Z, Chen H, Li J, and Shen L.: Multi model forecast biases of the diurnal variations of intense rainfall  
640 in the Beijing-Tianjin-Hebei region. *Science China Earth Sciences*, 65(8):1490–1509,  
641 <https://doi.org/10.1007/s11430-021-9905-4>, 2022.

642  
643  
644  
645  
646  
647  
648  
649  
650  
651

652 **Table 1.** Gradient decrease scoring table for station-by-station (time) rainstorm forecasts. The values are  
653 normalized, i.e., score = original data/100.

Observation (mm)	Forecast (mm)			
	25-49.9	50.0-99.9	100.0-249.9	≥250
<25.0	--	0	0	0
25.0-39.9	--	0.4	0	0
40.0-49.9	--	0.7	0.4	0
50.0-99.9	0.4	1	0.8	0.4
100.0-249.9	0	0.8	1	0.9
≥250.0	0	0.4	0.8	1

654



655

**Table 2.** Examples of station-specific rainstorm precipitation scoring

	Observation	Forecast	Forecast	Correct, Reasonable
		A	B	or False
Precipitation	50 mm	48 mm	98 mm	--
Forecast effect	--	Good	Bad	Correct
Classic TS	--	0	1	False
Improved TS	--	0.4	1	False
PAS	--	0.998	0.398	Reasonable

656

657

658

659

660

661

662

663

664

665

666

667

**Table 3.** Classification of PAS for short-term heavy rainfall.

Scoring name	Notes on the scoring application
$PAS_{ ux10}$	PAS score for 1 h observed precipitation $u \geq 10$ mm or forecasted precipitation $x \geq 10$ mm
$PAS_{ ux20}$	PAS score for 1 h observed precipitation $u \geq 20$ mm or forecasted precipitation $x \geq 20$ mm



668

669

**Table 4.** Classification of PAS for 12 h accumulated precipitation.

Scoring name	Notes on the scoring application
PASC	12 h PAS clear and precipitation forecast accuracy score 12 h PAS overall precipitation prediction verification score
PAS <sub> ux0.1</sub>	i.e., PAS score for observed precipitation $u \geq 0.1$ mm or forecasted precipitation $x \geq 0.1$ mm
PAS <sub> ux10</sub>	PAS score for 12 h observed precipitation $u \geq 10$ mm or forecasted precipitation $x \geq 10$ mm
PAS <sub> ux25</sub>	PAS score for 12 h observed precipitation $u \geq 25$ mm or forecasted precipitation $x \geq 25$ mm
PAS <sub> ux50</sub>	PAS score for 12 h observed precipitation $u \geq 50$ mm or forecasted precipitation $x \geq 50$ mm
PAS <sub> ux100</sub>	PAS score for 12 h observed precipitation $u \geq 100$ mm or forecasted precipitation $x \geq 100$ mm

670

671

672

673

674

675

676

677

678

679

680



681

**Table 5.** Classification of PAS for 24 h accumulated precipitation.

Scoring name	Notes on the scoring application
PASC	24 h PAS clear and precipitation forecast accuracy score 24 h PAS overall precipitation prediction verification score
PAS <sub> ux0.1</sub>	i.e., PAS score for observed precipitation $u \geq 0.1$ mm or forecasted precipitation $x \geq 0.1$ mm
PAS <sub> ux10</sub>	PAS score for 24 h observed precipitation $u \geq 10$ mm or forecasted precipitation $x \geq 10$ mm
PAS <sub> ux25</sub>	PAS score for 24 h observed precipitation $u \geq 25$ mm or forecasted precipitation $x \geq 25$ mm
PAS <sub> ux50</sub>	PAS score for 24 h observed precipitation $u \geq 50$ mm or forecasted precipitation $x \geq 50$ mm
PAS <sub> ux100</sub>	PAS score for 24 h observed precipitation $u \geq 100$ mm or forecasted precipitation $x \geq 100$ mm

682

**Table 6.** Examples of forecast verification scores for general precipitation ( $u = 25, 50$  and  $100$  mm).

PAS value	Observation $u=10$ mm		Observation $u=25$ mm		Observation $u=45$ mm	Observation $u=50$ mm
	Insufficient	Excessive	Insufficient	Excessive	Insufficient	Insufficient
	forecast x	forecast x	forecast x	forecast x	forecast x	forecast x
PAS=0.8	5.9	14.7	14.7	36.8	26.6	29.5
PAS=0.7	4.9	16.0	12.3	39.9	22.2	24.7
PAS=0.5	3.3	18.3	8.3	45.8	15.0	16.7
PAS=0.3	1.9	21.0	4.8	--	8.7	9.7

684

685

686





687 **Table 7.** Same as Table 6, but for precipitation at the level of torrential rain and above ( $u = 25, 50$  and  $100$  mm).

PAS value	Observation $u=25$ mm		Observation $u=50$ mm		Observation $u=100$ mm	
	Excessive	Insufficient	Excessive	Insufficient	Excessive	Insufficient
	forecast x	forecast x	forecast x	forecast x	forecast x	forecast x
PAS=0.877	--	34.1	68.1	68.1	136.2	
PAS=0.7	--	24.7	79.9	49.4	159.7	
PAS=0.5	--	16.7	91.6	33.3	183.3	
PAS=0.3	52.4	9.7	104.9	19.4	209.7	
PAS=0.1	62.9	3.2	125.9	6.4	251.7	

688

689 **Table 8.** PAS and TS of 12 h accumulated precipitation at 12:00 UTC on July 16, 2019.

	Clear/rainy	$\geq 0.1$ mm	$\geq 10$ mm	$\geq 25$ mm	$\geq 50$ mm
PAS	0.808	0.617	0.256	0.200	0.104
TS	0.690	0.381	0.194	0.076	0.006

690

691

692

693

**Table 9.** Same as Table 8, but for 12:00 UTC on June 13, 2020.

	Clear/rainy	$\geq 0.1$ mm	$\geq 10$ mm	$\geq 25$ mm	$\geq 50$ mm
PAS	0.734	0.457	0.228	0.185	0.116
TS	0.816	0.625	0.338	0.149	0.036

694

695 **Table 10.** Accuracy scores of insufficient precipitation forecast (IPS), excessive precipitation forecast (EPS) and

696 insufficient and excessive precipitation forecast (IEPS) of 12 h accumulated precipitation for two precipitation

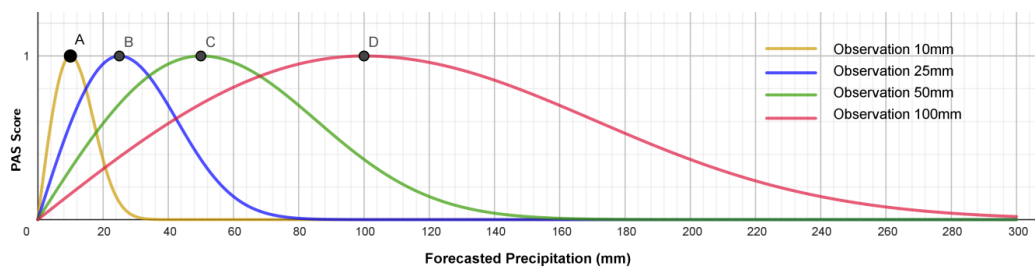
697

processes.

	IPS	EPS	IEPS
Case 1	-0.376	0.389	0.057
Case 2	-0.400	0.597	0.325

698

699



700

701

**Figure 1.** Schematic diagram of the precipitation forecast accuracy score (PAS) curves when the observed

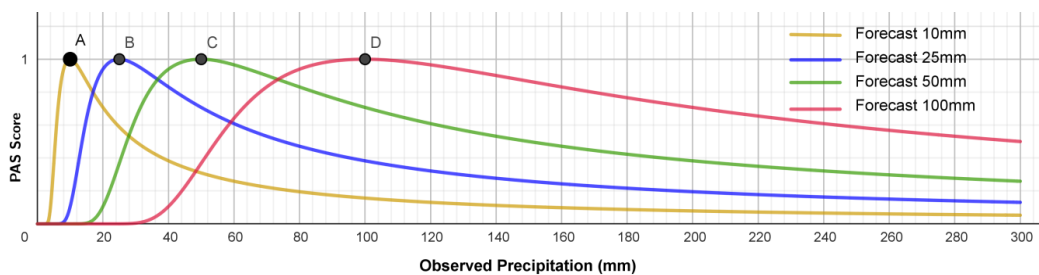
702

precipitation amounts are 10, 25, 50 and 100 mm.

703

704

705



706

707

**Figure 2.** PAS curves corresponding to different forecasted precipitation amounts ( $u = 10, 25, 50$  and  $100$

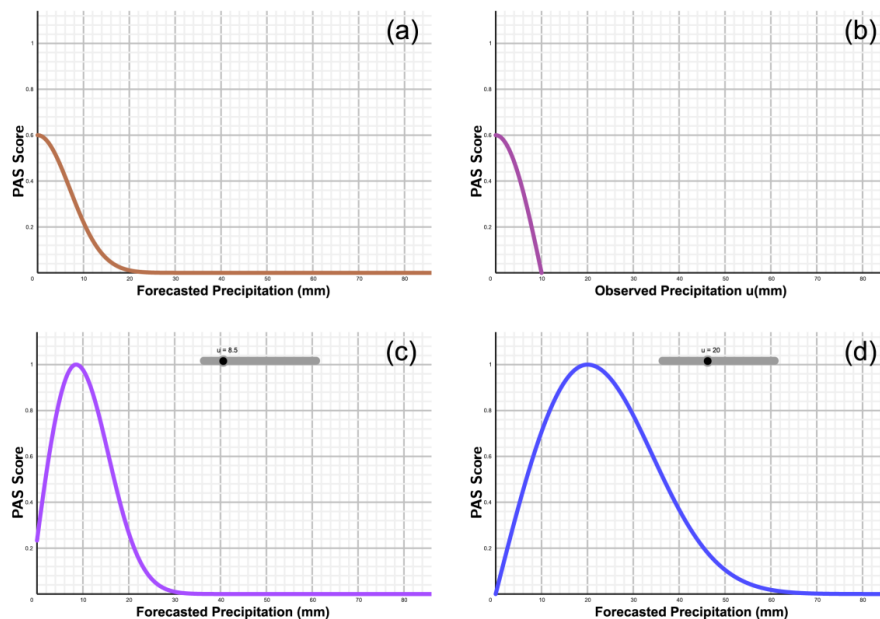
708

mm).

709

710

711



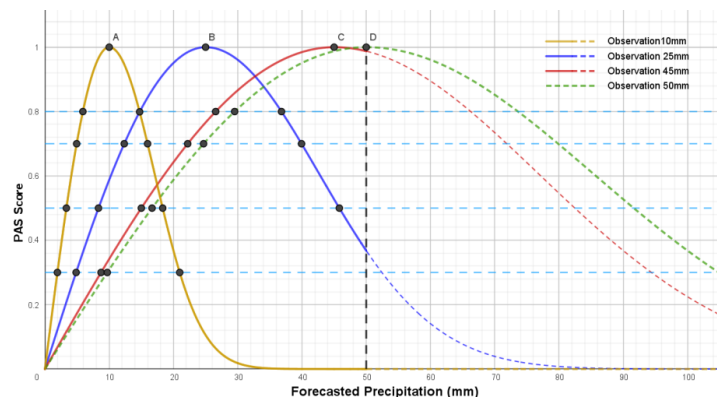
712

713 **Figure 3.** PAS curves of precipitation forecasts when (a) the observed precipitation  $u = 0$  mm, (b) the observed  
 714 precipitation  $u < 10$  mm and the forecasted precipitation  $x = 0$  mm (the horizontal coordinate denotes the  
 715 observed precipitation  $u$ ), (c) the observed precipitation  $0 \leq u < 10$  mm, and (d) the observed precipitation  $u \geq$   
 716 10 mm.

717

718

719



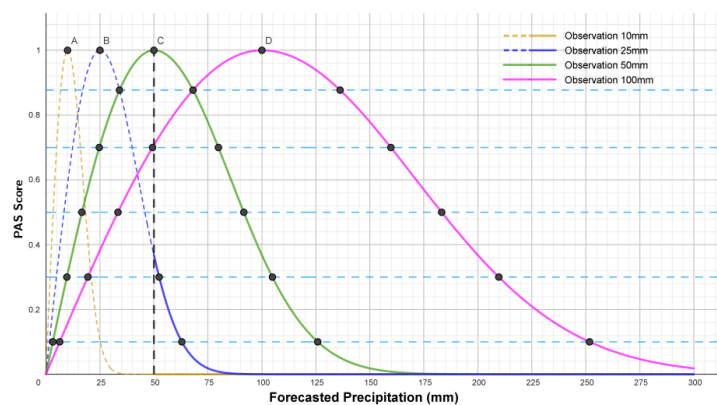
720

721

**Figure 4.** PAS curves of the forecasts under general precipitation conditions ( $u = 10, 25$  and  $45$  mm).



722

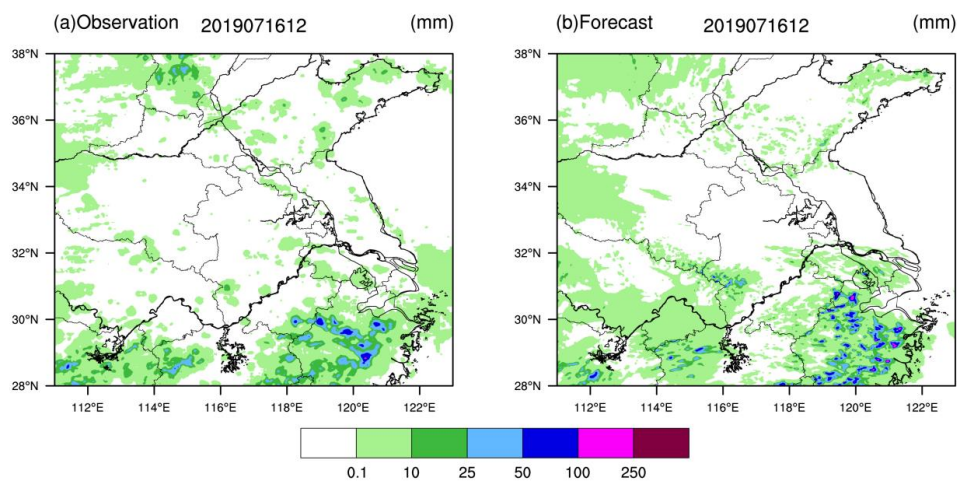


723

724 **Figure 5.** Same as Fig. 4, but for precipitation at the level of torrential rain and above ( $u = 25, 50$  and  $100$  mm).

725

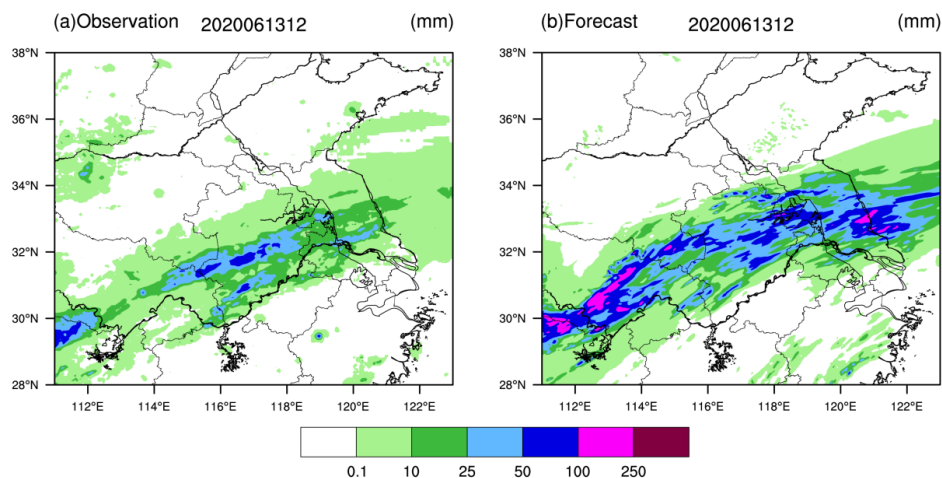
726



727

728 **Figure 6.** Accumulated precipitation (a) observed and (b) forecasted from 00:00 to 12:00 UTC on July 16, 2019.

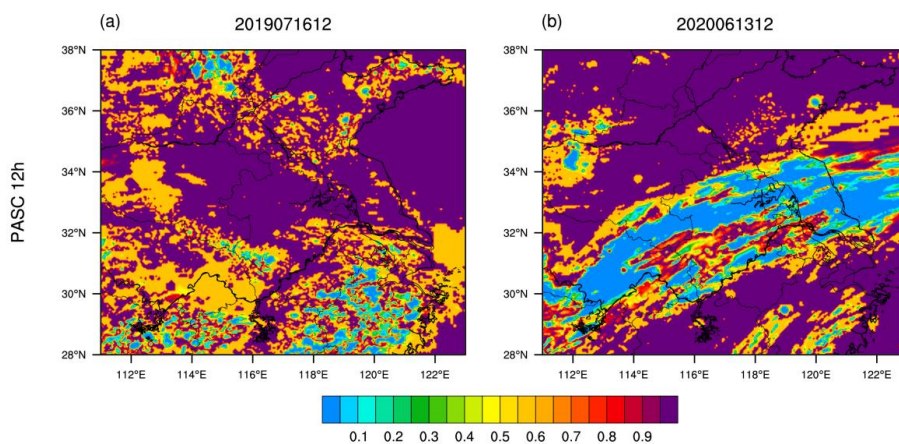
729



730

731 **Figure 7.** Accumulated precipitation (a) observed and (b) forecasted from 00:00 to 12:00 UTC on June 13, 2020.

732

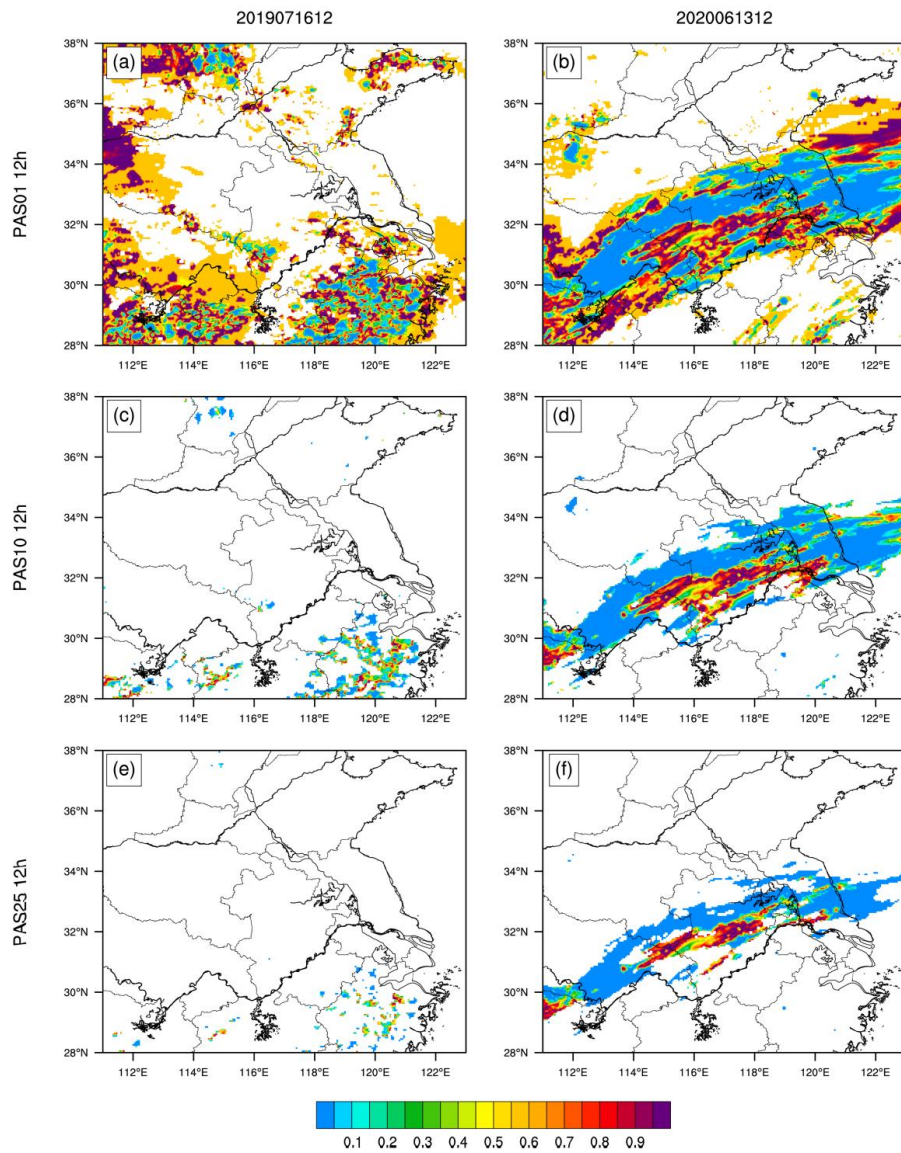


733

734 **Figure 8.** Distributions of the PAS clear/rainy forecast accuracy score (PASC) of 12 h accumulated precipitation for

735 (a) case 1 from 00:00–12:00 UTC on July 16, 2019, and (b) case 2 from 00:00–12:00 UTC on June 13, 2020.

736



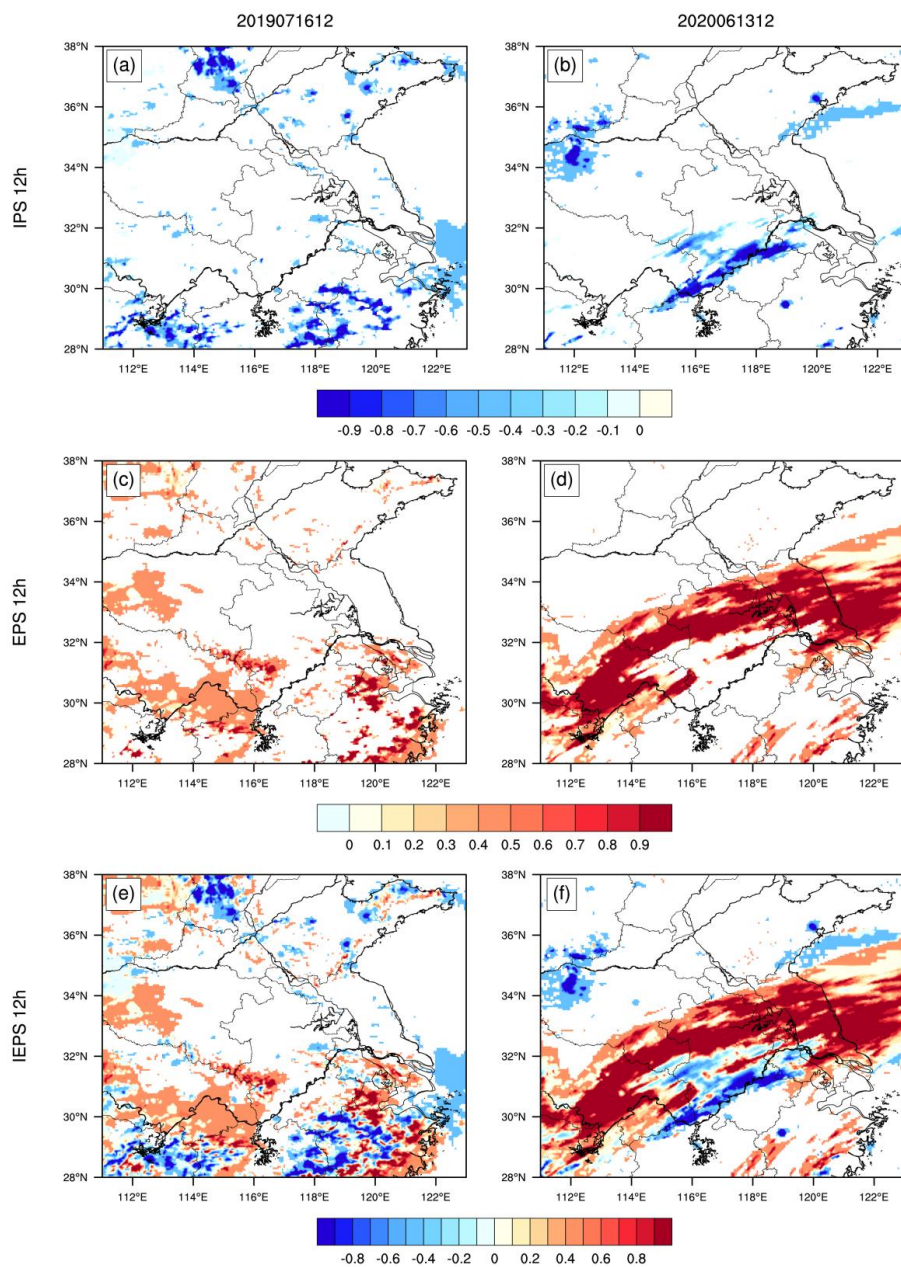
737

738

739  $\geq 0.1$  mm for (a) case 1 from 00:00–12:00 UTC on July 16, 2019 and (b) case 2 from 00:00–12:00 UTC on June 13,

740

2020,  $\geq 10$  mm for (c) case 1 and (d) case 2, and  $\geq 25$  mm for (e) case 1 and (f) case 2.



741

742 **Figure 10.** Distributions of IPS of 12 h accumulated precipitation for (a) case 1 from 00:00–12:00 UTC on July 16,

743 2019, and (b) case 2 from 00:00–12:00 UTC on June 13, 2020, EPS for (c) case 1 and (d) case 2, and IEPS for (e)

744

case 1 and (f) case 2.

745