

## Highlights

### **Balancing Accuracy versus Precision: Enhancing the Usability of Sub-Seasonal Forecasts**

Etienne Dunn-Sigouin, Erik W. Kolstad, C. Ole Wulff, Douglas J. Parker, Richard J. Keane

- Forecasts often have finer resolution than the scales they can accurately predict
- Developed user-method to evaluate trade-off between forecast accuracy and precision
- Reducing forecast precision improves accuracy and extends predictable lead-times
- Post-processing forecasts to accurate scales yields more actionable information

# Balancing Accuracy versus Precision: Enhancing the Usability of Sub-Seasonal Forecasts

Etienne Dunn-Sigouin<sup>a,b</sup>, Erik W. Kolstad<sup>a,b</sup>, C. Ole Wulff<sup>a,b</sup>, Douglas J. Parker<sup>a,b,c,d</sup>, Richard J. Keane<sup>c,e</sup>

<sup>a</sup>*NORCE Norwegian Research Center AS, Bergen, Norway*

<sup>b</sup>*Bjerknes Center for Climate Research, Bergen, Norway*

<sup>c</sup>*School of Earth and Environment, University of Leeds, Leeds, UK*

<sup>d</sup>*NCAS National Centre for Atmospheric Science, University of Leeds, Leeds, UK*

<sup>e</sup>*Met Office, Exeter, UK*

---

## Abstract

Forecasts are essential for climate adaptation and preparedness, such as in early warning systems and impact models. A key limitation to their practical use is often their coarse spatial grid spacing. However, another less frequently discussed but crucial limitation is that forecasts are often more precise than they are accurate when their grid spacing is finer than the scales they can accurately predict. Here, we adapt the fractions skill score, a metric conventionally used to quantify spatial forecast accuracy by the meteorological community, to help users navigate the trade-off between forecast accuracy versus precision. We demonstrate how this trade-off can be visualized for daily European precipitation, focusing on deterministic predictions of anomalies and probabilistic predictions of extremes, derived from three years of sub-seasonal forecasts from the European Centre for Medium-Range Weather Forecasts (ECMWF). Our results show that decreasing precision through spatial aggregation increases forecast accuracy, extends predictable lead times, and enhances the maximum possible accuracy relative to the grid scale, while increased precision diminishes these benefits. Notably, spatial aggregation benefits daily-accumulated forecasts more than weekly-accumulated ones, per unit lead-time. We demonstrate the practical value of our approach in three examples: communicating early warnings, managing hydropower capacity, and commercial aviation planning—each characterized by distinct

---

*Email address:* `etdu@norceresearch.no` (Etienne Dunn-Sigouin)

user constraints on accuracy, spatial scale, or lead-time. The results suggest a different approach for using forecasts; post-processing forecasts to focus on the most accurate scales rather than the default grid scale, thus offering users more actionable information.

*Keywords:* sub-seasonal forecasts, usability gap, forecast skill horizon, climate adaptation, ECMWF, fractions skill score

---

## 1 **Practical Implications**

2     Our results address a critical need for more accurate and actionable sub-  
3 seasonal forecasts, especially at longer lead times when crucial decisions are  
4 made. By adapting the fractions skill score, we illustrate how users can  
5 visualize and optimize the trade-off between a forecast’s spatial precision  
6 and its accuracy. Rather than relying on default high-resolution grids, we  
7 show that spatially aggregating forecasts can extend predictability and offer  
8 a clearer perspective on potential weather hazards. This approach not only  
9 complements existing forecast systems but also provides insights into when  
10 and where coarser-scale information is more dependable than finer scales.  
11 Ultimately, practitioners gain a practical tool that highlights where and how  
12 forecast aggregation pays dividends for planning at longer lead times. We  
13 demonstrate the value of the tool using three real-world examples.

14     In the first example, we show how early warning systems can benefit from  
15 the extended lead times offered by spatial aggregation. We demonstrate this  
16 using the 2023 Storm Hans case in Norway, which required timely alerts to  
17 protect lives and infrastructure. Even with the inherent uncertainty of pre-  
18 cipitation forecasts, aggregating them across broader areas yields more robust  
19 indications of impending extreme rainfall. This method could allow forecast-  
20 ers to issue warnings earlier, while policy makers and emergency managers  
21 stand to gain crucial time to mobilize resources.

22     The second example focuses on hydropower operations, where decisions  
23 are driven by localized hydrological processes but still benefit from a strategic  
24 view of precipitation patterns. Because releasing water from reservoirs too  
25 early can be costly, operators need maximum confidence in imminent rainfall  
26 forecasts. By matching the spatial aggregation scale to the watersheds of  
27 interest, hydropower managers can zero in on the most relevant signals. Our  
28 analyses highlight how post-processing forecasts at a watershed scale sharp-  
29 ens the focus on potential inflows, thereby supporting economically and en-

30 vironmentally sound reservoir management. Such tailored forecasting helps  
31 optimize water releases, reducing both flood risk, infrastructure damage and  
32 lost revenue opportunities.

33 In the third example, commercial aviation stands to benefit from spatial  
34 aggregation of forecasts when flights must be scheduled or canceled days in  
35 advance to minimize disruption. Spatial aggregation enables major carriers  
36 to detect the broader “footprint” of storms like Storm Hans well ahead of  
37 time, increasing confidence in decisions regarding flight cancellations, route  
38 changes, and resource allocation. Although individual airports lose fine-scale  
39 information, airlines can avoid the larger financial losses and passenger in-  
40 convenience that arise from last-minute adjustments over the entire network.  
41 This balance between lead time and accuracy can make flight networks more  
42 robust under uncertain weather conditions, ultimately improving safety and  
43 travel reliability while minimizing costs.

44 Taken together, these examples reveal how our adapted fractions skill  
45 score framework opens the door to “scalable” forecasts that users can cus-  
46 tomize to their unique spatial constraints. By offering a method to sys-  
47 tematically aggregate forecasts, decision-makers can glean earlier and more  
48 trustworthy signals of potential high-impact events, optimizing their inter-  
49 ventions in sectors ranging from disaster management to energy production  
50 and transportation. Our method does not eliminate all forecast uncertain-  
51 ties but provides a structured way to capitalize on known forecast strengths.  
52 In doing so, it encourages a shift from assuming that higher-resolution fore-  
53 casts are always better, to aligning forecast precision with the scales that can  
54 actually be predicted and those most relevant to users.

## 55 1. Introduction

56 In an era increasingly defined by climate change, the importance of weather  
57 and climate forecasts for society has surged, encompassing predictions from  
58 days to a decade ahead (Merryfield et al., 2020; White et al., 2022; O’Kane  
59 et al., 2023). This shift reflects a broader understanding of the critical role  
60 these forecasts play in managing the variable and often extreme environmen-  
61 tal conditions caused by climate variability and change, and their integration  
62 into society reflects a stronger push for adaptation and preparedness (God-  
63 dard, 2016; Trenberth et al., 2016; Coughlan de Perez et al., 2022). Forecasts  
64 support a number of international adaptation efforts such as the World Me-  
65 teorological Organisation’s Global Framework for Climate Services (Hewitt

et al., 2012), the United Nations Early Warnings for All initiative (EW4ALL, WMO, 2022), and the European Union’s financial sustainability taxonomy (European-Commission, 2020), and play a critical role in weather and climate services within the private sector (Cusick, 2019; Lam et al., 2023; Price et al., 2024). Forecasts are also used to predict impacts that are societally important but not directly modeled by their systems (Merz et al., 2020). These impact models vary widely in design and what they predict, such as floods, droughts, shipping routes, insurance risk, disease spread, agricultural cycles and renewable energy production (e.g., Torralba et al., 2017; Rösli et al., 2021; Graham et al., 2022; Haupt et al., 2018, 2019a,b).

There is, however, a well-known usability gap between the production of weather and climate information and its use (Lemos et al., 2012; Van den Hurk et al., 2018; Findlater et al., 2021). This gap is often attributed to the limited spatial resolution of forecasts, which often fails to meet the fine-scale precision required by users due to the prohibitive cost of high-resolution modeling.

Another less recognized yet crucial limitation to the practical use of forecasts, which is the focus of this paper, is the rapid decline in forecast skill at finer spatial scales as predictions extend into the future. This degradation in accuracy stems from faster error growth at smaller scales, where predictability is inherently linked to spatial size (Lorenz, 1969; Toth and Buizza, 2019). For instance, while slower large-scale cyclones spanning thousands of kilometers may be predictable over several days, smaller-scale thunderstorms operate on much shorter timescales, with predictability limited to a few hours. Spatial and temporal aggregation can be used to counteract this small-scale error growth by effectively filtering out high-frequency, small-scale noise, thereby enhancing the predictable signal from the lower-frequency, larger-scale circulation. Thus, aggregation helps to extend the limit of predictability, known as the “forecast skill horizon” (Buizza et al., 2015; Buizza and Leutbecher, 2015), but at the cost of spatial or temporal precision.

While temporal aggregation is routine, it is rare to find weather and climate information – either from forecasts or projections – that has been aggregated spatially. For example, sub-seasonal forecasts (often referred to as extended-range or monthly forecasts) are usually presented in the form of weekly aggregated information, such as in the online charts catalogue of the European Centre for Medium-Range Weather Forecasts (ECMWF, <https://charts.ecmwf.int/>). However, in this case no spatial aggregation is done;

104 the forecast charts are displayed at the model’s original grid spacing. Even  
105 commonly used daily-aggregated weather forecasts are shown on the default  
106 grid spacing despite a well known decline in forecast skill at smaller spatial  
107 scales over just a few days. This can lead to a critical mismatch between the  
108 apparent precision in forecast products and the underlying accuracy in the  
109 data. Such misrepresentation risks undermining effective use of weather and  
110 climate information (Nissan et al., 2019; Fiedler et al., 2021).

111 Two ways of making forecasts better fit user needs are improving their  
112 predictive skill at finer spatial scales or helping users more effectively uti-  
113 lize existing skill. While improving forecasts remains a formidable challenge  
114 (Bauer et al., 2015; Benjamin et al., 2019), recent advances in machine-  
115 learning-based models have made significant strides, offering performance  
116 that now rivals traditional dynamical forecast models (Lam et al., 2023;  
117 Ben Bouallègue et al., 2024; Price et al., 2024). Despite these breakthroughs,  
118 it is widely acknowledged that there are likely intrinsic limits to the forecast  
119 skill horizon, which no amount of model improvement can overcome (Lorenz,  
120 1969; Palmer et al., 2014). Thus, a pragmatic strategy involves helping users  
121 navigate the inherent trade-off between spatial accuracy and precision, opti-  
122 mizing existing forecasts for their needs.

123 Various strategies using spatio-temporal aggregation have been suggested  
124 (e.g., Gong et al., 2003; Gilleland et al., 2009; Jung and Leutbecher, 2008;  
125 Buizza and Leutbecher, 2015; Gehne et al., 2016; Toth and Buizza, 2019;  
126 Van Straaten et al., 2020; Young et al., 2020; Rivoire et al., 2023). These  
127 focus mostly on quantifying the forecast skill horizon where predictability is  
128 small and difficult to exploit in practice. An alternative, more user-oriented  
129 method is to use the fractions skill score (Roberts and Lean 2008), which  
130 quantifies where forecast predictability is high and usable. This approach,  
131 which aggregates forecasts over an increasing number of neighboring grid  
132 points, quantifies the trade-off between accuracy versus precision, and can  
133 be used to post-process the forecast according to the user’s preferred balance.  
134 The fractions skill score and a number of other closely related methods stand  
135 out for their intuitiveness and practical applicability, yet their use has been  
136 largely confined to the meteorological community (e.g., Gilleland et al., 2009;  
137 Jolliffe and Stephenson, 2012; Keane et al., 2016; Zhao and Zhang, 2018;  
138 Schwartz, 2019; Cafaro et al., 2021).

139 In this study, we propose a novel methodology, based on the fractions  
140 skill score, that realigns forecast capabilities with end-user requirements,  
141 thereby enhancing their practical application. The innovative aspect of this

142 method lies in shifting its focus from verifying forecasts for meteorologists to  
143 optimizing them for users. The method is applied to sub-seasonal forecasts  
144 with lead times ranging from 1 day to several weeks—an essential time-frame  
145 for many decision-making processes (Merz et al., 2020; White et al., 2022).

146 This work is part of the Climate Futures collaboration, an interdis-  
147 plinary and intersectoral initiative that, since 2020, has brought together  
148 public and private organizations in Norway to co-produce weather and cli-  
149 mate prediction-based tools and services. A key case study involved working  
150 with Tryg Forsikring, a private insurance company, to incorporate forecasts  
151 into their decision-making in order to comply with sustainable finance regu-  
152 lations. The collaboration served two purposes: to inform insurance profes-  
153 sionals on the practical limitations of using forecasts, and to maximize the  
154 utility of forecasts for insurance impact modeling. Our experience with the  
155 insurance sector suggests our approach could be broadly applied to support a  
156 wide range of climate adaptation efforts across industries. The intended users  
157 are meteorologists and climate scientists who supply forecasts, or industry  
158 professionals who can use forecast data but lack forecast expertise.

159 In the next section we start by introducing the forecast dataset and outlin-  
160 ing the methodology, which amounts to spatially aggregating forecasts before  
161 calculating commonly used metrics of forecast accuracy. In the results sec-  
162 tion, we apply the new method to evaluate the trade-off between precision  
163 and accuracy in European precipitation forecasts, derived from three years of  
164 sub-seasonal forecasts from the ECMWF. We explore how this method can  
165 aid users in interpreting deterministic predictions of precipitation anoma-  
166 lies and probabilistic predictions of extremes. To illustrate this, we use the  
167 example of Storm Hans, which struck Scandinavia, Northern Europe, and  
168 the Baltics in August 2023, with Norway bearing the brunt of its impact.  
169 Unusually approaching from the east rather than the west, the storm shat-  
170 tered century-old rainfall records in eastern Norway (Granerød et al., 2023).  
171 The resulting extreme rainfall triggered widespread flooding and landslides,  
172 severely damaging homes, roads, railways, and bridges, with estimated costs  
173 reaching 4 billion Norwegian Krone or 350 million euro (Ekroll, 2023). Over  
174 10,000 insurance claims were filed, and approximately 2,400 people were evac-  
175 uated—the largest such evacuation in Norway since World War II (KLP,  
176 2023). With extreme rainfall events expected to become more frequent due  
177 to climate change (Hanssen-Bauer et al., 2009), storm Hans exemplifies the  
178 growing challenges in climate adaptation.

179 Building on our findings and those of Roberts and Lean (2008), we end

the paper by discussing how the method could be applied to forecasts in three different contexts—communicating early warnings, managing hydropower capacity, and commercial aviation planning—each characterized by distinct user-constraints on accuracy, spatial scale, or lead-time. In each case, we enhance forecast utility by post-processing forecasts to focus on the most accurate spatial scales, rather than the default grid scale precision.

## 2. Data

We use three years (2020–2022) of sub-seasonal forecasts from the ECMWF (Buizza et al., 2018) downloaded from the MARS archive (ECMWF, 2024a). We use bi-weekly initializations on Mondays and Thursdays, for a total of 313 forecasts, each comprising 51 ensemble members running 46 days in the future. The initial 15 lead-time days are higher resolution ( $0.25^\circ \times 0.25^\circ$  grid spacing) than the last 31 days ( $0.5^\circ \times 0.5^\circ$ ), corresponding to approximately  $28 \text{ km}^2$  and  $56 \text{ km}^2$  at the equator, respectively. Accompanying each individual forecast is a set of retrospective forecasts. These were initialized on the same calendar day as the forecast over the previous 20 years and consist of 11 ensemble members. Such “hindcasts” provide an estimate of the climatological distribution accompanying each forecast.

The forecast-hindcast pairs correspond to different model versions over time (CY46R1, CY47R1, CY47R2, CY47R3) because the model is updated on the fly and our analysis spans multiple years. Changes in model cycles can influence model biases due to evolving model physics and data assimilation. While our approach (discussed in the next section) focuses on deviations from the model’s climatology, which helps mitigate systematic differences across model cycles, residual biases in the forecasts may remain. However, we do not expect these to qualitatively impact our results.

We focus our analysis on Europe ( $33^\circ\text{N}$  to  $73.5^\circ\text{N}$  and  $27^\circ\text{W}$  to  $35^\circ\text{E}$ ) and on predictions of daily and weekly-accumulated precipitation. A corresponding analysis of daily and weekly-mean 2-meter temperature forecasts for two years (2020–2021) is included in the supplementary materials. Forecast skill was verified relative to ERA5 reanalysis (Hersbach et al., 2023) for the same grid, domain, and time period as the forecast. Although ERA5 exhibits known biases, such as a tendency for excessive drizzle (Lavers et al., 2022), it remains a convenient benchmark for verification because its resolution matches that of the forecast. We note, however, that other observational datasets may be used for verification, and that we do not expect this choice to



216 impact our qualitative results. We also extended this analysis to storm Hans  
 217 in 2023, incorporating additional forecasts and hindcasts initialized between  
 218 3 and 7 August alongside ERA5 data.

219 Finally, to illustrate the spatial scale of the data, we convert its spa-  
 220 tial precision from gridpoint units to square kilometers, shown in the y-axis  
 221 labels of Fig. 1. Specifically, we simply rescale the nominal  $28 \text{ km}^2$  area  
 222 represented by one gridpoint<sup>2</sup> at the equator by the mean cosine of latitude  
 223 within the domain, consistent with the spherical geometry of Earth’s sur-  
 224 face. Consequently, one gridpoint<sup>2</sup> within the European domain corresponds  
 225 to approximately  $15 \text{ km}^2$ .

### 226 3. Methodology

227 We assess forecast accuracy as a function of precision and lead time using  
 228 modified versions of the Fractions Skill Score (FSS, Roberts and Lean, 2008).  
 229 Here, *accuracy* refers to the skill of the forecast quantified using a skill score,  
 230 and *precision* refers to the level of spatial aggregation of the forecast. We  
 231 begin by summarizing the original FSS developed for the meteorological com-  
 232 munity in section 3a, followed by our adaptations for end users in sections  
 233 3b,c,d. Next, we introduce a modified version of the Extreme Forecast Index  
 234 (EFI, Lalaurette, 2003), which we use to demonstrate the value of optimizing  
 235 forecast accuracy during Storm Hans in section 5.

236 To facilitate the computation of scores and indices in the following sec-  
 237 tions, it is useful to first convert the reanalysis verification into the same  
 238 format as the forecasts and hindcasts. Table 1 summarizes the variables and  
 239 their dimensions defined in section 3. Specifically, forecasts  $f(m, e, t, i, j)$  are  
 240 characterized by dimensions of forecast initialization ( $m$ ), ensemble mem-  
 241 ber ( $e$ ), lead time ( $t$ ), latitude ( $i$ ) and longitude ( $j$ ). These correspond to  
 242 a verification  $v_f(m, t, i, j)$  from ERA5 reanalysis, where  $e = 1$  and  $t = 1$   
 243 represents the 24-hour period after the forecast initialization date  $m$ . Simi-  
 244 larly, hindcasts  $h(m, y, e, t, i, j)$  which include a hindcast year dimension ( $y$ ),  
 245 correspond to a verification  $v_h(m, y, t, i, j)$  with  $e = 1$  that spans the past  
 246 twenty years for each calendar date of forecast initialization  $m$ .

#### 247 3.1. Fractions Skill Score

248 The FSS uses binary forecast and verification data to assess the skill of the  
 249 forecast at different levels of spatial aggregation. Roberts and Lean (2008)  
 250 developed their method using deterministic forecasts of precipitation, i.e.,

251 with only one ensemble-member. First, they converted the forecast  $f$  and  
 252 verification  $v_f$  to binary values based on a predefined absolute threshold (e.g.,  
 253 4 mm). If the precipitation amount exceeded this threshold, the value was  
 254 set to 1; otherwise, it was set to 0. Next, for each grid point, they averaged  
 255 surrounding points within a square of length  $n$  (this process is referred to  
 256 hereafter as *aggregation*), yielding an aggregated forecast  $F$  and verification  
 257  $V_F$  (see Equations 1 and 2). These aggregations are not binary, but have  
 258 fractional values between 0 and 1.

$$F(n, m, t, i, j) = \frac{1}{n^2} \sum_{k=1}^n \sum_{l=1}^n f \left[ m, t, i+k-1 - \frac{(n-2)}{2}, j+l-1 - \frac{(n-1)}{2} \right] \quad (1)$$

$$V_F(n, m, t, i, j) = \frac{1}{n^2} \sum_{k=1}^n \sum_{l=1}^n v_f \left[ m, t, i+k-1 - \frac{(n-2)}{2}, j+l-1 - \frac{(n-1)}{2} \right] \quad (2)$$

259 Grid points within the square of length  $n$  but outside the domain defined  
 260 in section 2 (i.e., Europe) were set to zero. By comparing the mean square  
 261 error of the aggregated forecast over the domain (equation 3) with that cal-  
 262 culated from an aggregated reference forecast for each  $n$ , they obtained the  
 263 FSS (equation 4).

$$MSE_F(n, m, t) = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J [F(n, m, t, i, j) - V_F(n, m, t, i, j)]^2 \quad (3)$$

$$FSS(n, m, t) = 1 - \frac{MSE_F(m, n, t)}{MSE_{REF}(n, t)} \quad (4)$$

264 An FSS value of 1 signifies perfect forecast accuracy relative to the ver-  
 265 ification data, while an FSS value of 0 or less indicates the forecasts are no  
 266 better or worse than a reference forecast. The choice of reference forecast is  
 267 up to the user (e.g., random forecast, climatology or something else).

### 268 3.2. Modified fractions skill scores

269 Next, we modify the original FSS to better adapt it to end-users. The  
 270 details of the modified scores are described in the next sections and the  
 271 primary steps can be summarized as:

- 272 1. Aggregate the raw forecast and verification fields spatially following  
273 equations 1 and 2.
- 274 2. Compute a standard grid point-wise score such as Mean Square Error  
275 or Brier Score.
- 276 3. Compute a skill score by comparing the score to a suitably aggregated  
277 reference forecast and average over all spatial grid-points and forecasts.

### 278 3.2.1. *Fractions Mean-Square Error Skill Score*

279 The Fractions Mean-Square Error Skill Score (FMSESS) quantifies the ac-  
280 curacy of ensemble-mean forecast anomalies, averaged across all grid points  
281 and forecasts, over varying spatial aggregation scales. Unlike the original  
282 FSS, which uses binary threshold-based values, the FMSESS incorporates  
283 anomalies relative to climatology. This modification offers several benefits to  
284 users: 1) it generalizes the widely used mean-square error skill score across  
285 multiple spatial scales (Jolliffe and Stephenson, 2012), enabling better com-  
286 parisons with past studies; 2) it simplifies the interpretation by providing a  
287 measure of accuracy independent of threshold; 3) it provides a better esti-  
288 mate of forecast skill by incorporating a mean-bias correction.

289 Forecast anomalies  $\tilde{f}$  are computed by taking the ensemble-mean of the  
290 difference between the forecast and the hindcast climatology (equation 5),  
291 while verification anomalies  $\tilde{v}_f$  are calculated by subtracting the verification  
292 climatology from each verification  $v_f$  (equation 6).

$$\tilde{f}(m, t, i, j) = \frac{1}{E} \sum_{e=1}^E \left[ f(m, e, t, i, j) - \frac{1}{Y} \sum_{y=1}^Y h(m, y, e, t, i, j) \right] \quad (5)$$

$$\tilde{v}_f(m, t, i, j) = v_f(m, t, i, j) - \frac{1}{Y} \sum_{y=1}^Y v_h(m, y, t, i, j) \quad (6)$$

293 We use a single date  $m$  to define the climatologies for each forecast  
294 and verification for simplicity. A more robust estimate of the climatology  
295 could be achieved by incorporating additional dates centered around the fore-  
296 cast/verification date, as demonstrated by ECMWF’s M-climate (ECMWF,  
297 2024b). However, we do not expect this choice to qualitatively affect our  
298 main results.

299 Aggregated forecast and verification anomalies  $\tilde{F}$  and  $\tilde{V}_F$  are then used  
300 to calculate the FMSESS similar to the original FSS (equations 7 and 9).

301 The aggregated version of the verification climatology (second term on the  
 302 right-hand side of equation 6) is used as the reference forecast to calculate  
 303 the reference mean-square error in the FMSESS (equation 8).

$$MSE_{\tilde{F}}(n, m, t) = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J [\tilde{F}(n, m, t, i, j) - \tilde{V}_F(n, m, t, i, j)]^2 \quad (7)$$

$$MSE_{R\tilde{E}F}(n, m, t) = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J [\tilde{F}(n, m, t, i, j) - \frac{1}{Y} \sum_{y=1}^Y V_H(n, m, y, t, i, j)]^2 \quad (8)$$

$$FMSESS(n, t) = 1 - \sum_{m=1}^M \frac{MSE_{\tilde{F}}(m, n, t)}{MSE_{R\tilde{E}F}(m, n, t)} \quad (9)$$

### 304 3.2.2. *Fractions Brier Skill Score*

305 The Fractions Brier Skill Score (FBSS) quantifies the accuracy of forecast  
 306 extremes, averaged over all grid points and forecasts, across varying spatial  
 307 aggregation scales. It introduces two key modifications to the original FSS.  
 308 First, it provides a probabilistic assessment of skill by utilizing an ensemble of  
 309 forecasts instead of a single deterministic forecast. Second, it uses a threshold  
 310 value from a predefined quantile based on the hindcast climatology, rather  
 311 than an absolute threshold. In this study, we demonstrate the method using  
 312 the 0.1 and 0.9 quantiles, corresponding to dry and wet extremes. Similar to  
 313 the FMSESS, the FBSS offers distinct advantages to users: 1) it generalizes  
 314 the widely-used Brier Skill Score across multiple spatial scales (Jolliffe and  
 315 Stephenson, 2012), enabling better comparisons with past studies; 2) it im-  
 316 proves the evaluation of extremes via probabilistic scoring and quantile-based  
 317 bias correction.

318 We start by defining a threshold value for extremes for a given quantile  
 319  $q$ . First, we calculate aggregated hindcast  $H(n, m, e, y, t, i, j)$  and verifica-  
 320 tion  $V_H(n, m, y, t, i, j)$  for each  $n$  following equations 1 and 2. The forecast  
 321 threshold value  $F_q(n, m, t, i, j)$  is then computed for the quantile  $q$  from a  
 322 sample of  $e$  ensemble members and  $y$  hindcast years in hindcast  $H$ . Corre-  
 323 spondingly, the verification threshold value  $V_q(n, m, t, i, j)$  is computed for  
 324 the quantile  $q$  from a sample of  $y$  hindcast years in verification hindcast  $V_H$ .

325 Next, we compute the Brier Score  $BS_q$  for a given quantile  $q$ . First, we cal-  
 326 culate the aggregated forecast  $F(n, m, e, t, i, j)$  and verification  $V_f(n, m, t, i, j)$   
 327 for each  $n$ . Then, we compute the forecast probability  $P_{F_q}(n, m, t, i, j)$  by  
 328 determining the fraction of ensemble members  $e$  in forecast  $F(n, m, e, t, i, j)$   
 329 that exceed the threshold value  $F_q(n, m, t, i, j)$ . Similarly, we compute the bi-  
 330 nary verification  $P_{V_q}(n, m, t, i, j)$  based on whether the verification  $V_f(n, m, t, i, j)$   
 331 crosses the threshold value  $V_q(n, m, t, i, j)$ , assigning 1 to values above the  
 332 threshold and 0 below. The squared difference between the forecast probabil-  
 333 ity and the binary verification is then averaged over all values in the domain  
 334 (equation 10).

$$BS_q(n, m, t) = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J [P_{F_q}(n, m, t, i, j) - P_{V_q}(n, m, t, i, j)]^2 \quad (10)$$

335 Finally, the FBSS is computed by comparing the Brier Score of the fore-  
 336 cast with one calculated from a reference forecast and averaging over all fore-  
 337 casts (equation 11). The reference forecast used in the reference Brier Score  
 338 is simply the quantile  $q$  used to define the threshold ( $q = 0.9$  or  $q = 0.1$ ).

$$FBSS_q(n, t) = 1 - \sum_{m=1}^M \frac{BS_q(n, m, t)}{BS_{REF_q}(n, m, t)} \quad (11)$$

### 339 3.2.3. Statistical Significance of Skill Scores

340 Statistical significance of the FMSESS and FBSS (equations 9 and 11) is  
 341 evaluated using bootstrapping. We generate a distribution of scores for each  
 342 spatial scale  $n$  and lead-time  $t$  by resampling the forecasts  $m$  10,000 times  
 343 with replacement. The null hypothesis is that the score is zero or negative,  
 344 i.e., less than the climatological reference forecast. A score is considered  
 345 significantly more skillful (at the 5% level) than the reference if 95% of the  
 346 resampled distribution is greater than zero.

### 347 3.3. Fractions Extreme Forecast Index

348 The skill scores introduced in the previous sections quantify the accuracy  
 349 of past forecasts. However, for these insights to inform real-time decision-  
 350 making, users require a method to apply them operationally. Roberts and  
 351 Lean (2008) proposed post-processing forecasts via spatial aggregation, en-  
 352 abling users to optimize the balance between accuracy and precision based

on a past assessment of forecast skill. A key limitation of spatial aggregation is that it reduces the amplitude of raw forecast values, making them less intuitive for users accustomed to working with unprocessed data. To address this, we propose normalizing the aggregated forecast by an aggregated reference, similar to the procedure to define the Extreme Forecast Index (EFI, Lalaurette, 2003). The EFI measures how extreme a probabilistic forecast is relative to its climatology by comparing the cumulative distributions of the forecast and its corresponding hindcast, and is operationally employed by the ECMWF. In this subsection, we outline the computation of the EFI and then detail our modifications to enhance its applicability for end-users.

To compute the EFI, threshold values are first defined for each quantile  $q$  varying from  $0 < q < 1$  in steps of  $\Delta q$ . The hindcast threshold  $h_{thresh}(q, m, t, i, j)$  is determined as the value of the  $q^{th}$  quantile from a sample of  $e$  ensemble members and  $y$  hindcast years in hindcast  $h(m, e, y, t, i, j)$ . Then, for each forecast  $f(m, e, t, i, j)$ , the fraction of ensemble members below the hindcast threshold  $h_{thresh}(q, m, t, i, j)$  is computed, called  $fr_f(q, m, t, i, j)$ . The EFI is derived by summing the difference between the quantile  $q$  and the corresponding forecast fractions  $fr_f(q, m, t, i, j)$  across all quantiles, normalized by  $q(1 - q)$ , and multiplied by the quantile step  $\Delta q$  (equation 12).

$$EFI(m, t, i, j) = \frac{2}{\pi} \sum_{q=0}^1 \frac{q - fr_f(q, m, t, i, j)}{q(1 - q)} \Delta q \quad (12)$$

EFI values range from -1 to 1, where -1 indicates that the entire cumulative forecast distribution is below the cumulative hindcast distribution, and +1 indicates it is entirely above. An  $|EFI| > 0.8$  typically signifies an extreme event (ECMWF, 2024c).

For comparison with the forecast, we define an analogous Extreme Verification Index (EVI) using the verification  $v(m, t, i, j)$ , verification hindcast  $v_H(m, y, t, i, j)$ , verification hindcast threshold  $v_{thresh}(q, m, t, i, j)$  and verification fraction  $fr_v(q, m, t, i, j)$  in equation 13. The verification hindcast threshold is determined as the value of the  $q^{th}$  quantile from a sample of  $y$  hindcast years in the verification hindcast, and the verification fraction is set to 1 if the verification lies above the threshold and to 0 if it lies below.

$$EVI(m, t, i, j) = \frac{2}{\pi} \sum_{q=0}^1 \frac{q - fr_v(q, m, t, i, j)}{q(1 - q)} \Delta q \quad (13)$$

383 The Fractions Extreme Forecast Index (FEFI) is computed in the same  
 384 way as the original EFI except it utilizes the aggregated forecasts  $F(n, m, e, t, i, j)$ ,  
 385 hindcasts  $H(n, m, e, y, t, i, j)$ , hindcast threshold  $H_{thresh}(n, q, m, t, i, j)$  and  
 386 forecast fraction  $FR_F(n, q, m, t, i, j)$  with the additional aggregation dimension  
 387  $n$  (equation 14).

$$FEFI(n, m, t, i, j) = \frac{2}{\pi} \sum_{q=0}^1 \frac{q - FR_F(n, q, m, t, i, j)}{q(1 - q)} \Delta q \quad (14)$$

388 Thus, the FEFI quantifies how extreme the forecast is relative its its  
 389 climatology across different spatial scales.

#### 390 4. Quantifying Accuracy Versus Precision in European Precipita- 391 tion Forecasts

392 In this section, we evaluate the impact of spatial aggregation on sub-  
 393 seasonal precipitation forecasts. We calculate FMSESS and FBSS for daily  
 394 and weekly-accumulated precipitation over the entire European domain, ex-  
 395 amine their regional variations, and compare the results with those for 2-  
 396 meter temperature forecasts. By quantifying the trade-off between spatial  
 397 accuracy and precision, we offer a clearer understanding of how spatial ag-  
 398 gregation influences forecast performance.

399 Spatial aggregation improves the accuracy of daily-accumulated precip-  
 400 itation forecasts. Figure 1a shows the FMSESS for daily-accumulated pre-  
 401 cipitation anomalies across various lead times and spatial scales. At the  
 402 grid-scale, accuracy is high initially ( $> 0.8$ ) but decreases with lead-time  
 403 ( $< 0.1$ ), with forecasts remaining skillful for up to 10 days (in agreement  
 404 with Rivoire et al., 2023). Spatial aggregation not only increases accuracy  
 405 for a given lead time but also extends the forecast skill horizon (as pointed out  
 406 by Buizza et al., 2015; Buizza and Leutbecher, 2015), indicated by the right-  
 407 slanted skill contours and significance hatching. This approach also improves  
 408 accuracy for extreme precipitation, measured by the FBSS for the 0.9 and  
 409 0.1 quantiles, although forecasting extremes is generally less accurate than  
 410 forecasting anomalies (compare Fig. 1a and 1c,e). Notably, low precipitation  
 411 extremes can show reduced accuracy with spatial aggregation (left-slanted  
 412 contours close to the grid-scale for lead times greater than 8 days in Fig. 1e).  
 413 Near the grid-scale, forecasts are more accurate for low precipitation thresh-  
 414 olds due to the model’s tendency to predict no precipitation. Greater spatial

415 aggregation raises the likelihood of non-zero precipitation thresholds, making  
416 the predictions more challenging and reducing accuracy.

417 Reducing precision can extend predictable lead-times of daily-accumulated  
418 precipitation by a few days. The green shading on the right-hand side panels  
419 of Fig. 1 illustrates the lead-time gained by spatial aggregation. For exam-  
420 ple, a forecast with a precision of 1 grid point<sup>2</sup> and an accuracy of 0.5 that  
421 is spatially aggregated to a precision of 33 grid points<sup>2</sup> represents an accu-  
422 racy gain of 0.25 (black arrow). This increase in accuracy is equivalent to  
423 gain in 2 lead-time days (blue arrow) because the aggregated forecast drops  
424 to the same level of accuracy as the grid-scale forecast two days later (red  
425 arrow). For low precipitation extremes, reducing precision can also lead to  
426 a loss of lead-time relative to the grid-scale (pink shading, Fig. 1f), though  
427 this mainly occurs where forecast accuracy is low ( $< 0.2$ ).

428 The highest levels of forecast accuracy are only achievable with spatial  
429 aggregation. Cross-hatching in Fig. 1b,d,e shows where forecast accuracy  
430 exceeds the maximum achievable accuracy at the grid scale, i.e., lead-time  
431  $t = 1$ . Gains from spatial aggregation are substantial for anomalies (0.8  
432 to 0.95) and high precipitation extremes (0.5 to 0.7), and even greater for  
433 low precipitation extremes (0.3 to 0.7). Overall similar results are found for  
434 winter and summer only forecasts, where winter generally exhibits higher  
435 accuracy (Figs. S1 and S2).

436 Spatial aggregation also improves regional forecast accuracy of daily-  
437 accumulated precipitation. Figure 2 shows latitude-longitude maps of fore-  
438 cast accuracy for daily-accumulated precipitation at lead-day 5, comparing  
439 two spatial precision levels: the grid scale (left) and 33 grid points<sup>2</sup> (right).  
440 The FMSESS and FBSS are calculated regionally at each latitude and longi-  
441 tude by omitting the domain average in equations 7 and 10. It is important  
442 to note that the spatially aggregated forecasts are displayed with the same  
443 grid spacing as the raw forecasts (e.g.,  $0.25^\circ \times 0.25^\circ$ ), but their *effective*  
444 spatial resolution is reduced since each grid point represents the aggregate  
445 of its neighboring grid points. At the grid scale, forecast accuracy varies  
446 regionally, with higher accuracy over mountainous regions like western Nor-  
447 way and the Alps, and is higher for anomalies than extremes (Fig. 2 left).  
448 Spatial aggregation increases overall forecast accuracy and extends the fore-  
449 cast horizon, as indicated by the darker shading and reduced hatching in the  
450 right-hand versus left-hand panels of Fig. 2. This suggests that the European  
451 domain-averaged results in Fig. 1 generally hold regionally. Similar patterns  
452 are observed for different lead times and seasons (not shown).



453 When measured per unit lead-time, spatial aggregation benefits weekly-  
 454 accumulated precipitation forecasts less than daily-accumulated ones. Figure  
 455 3 shows forecast accuracy for weekly-accumulated precipitation anomalies  
 456 and extremes across various lead-times and spatial scales. At the grid-scale,  
 457 forecast accuracy is high in the first week, low in the second, and as skill-  
 458 ful as climatology in the third or fourth, consistent with daily-accumulated  
 459 precipitation forecasts (compare Figs. 3 and 1, left). However, the improve-  
 460 ments from spatial aggregation for weekly-accumulated forecasts are marginal  
 461 compared to those for daily-accumulated forecasts when assessed per unit  
 462 lead-time: 0.2-0.4 weeks versus 2-3 days (compare the slanted grey accuracy  
 463 contours in Figs. 3 and 1, right). Regionally, spatial aggregation modestly  
 464 improves accuracy and extends the forecast skill horizon for anomalies, but  
 465 less so for extremes (compare left and right-hand panels of Fig. 4).

466 Surface temperature forecasts, while generally more accurate than pre-  
 467 cipitation forecasts, benefit less from spatial aggregation. Figure S3 shows  
 468 forecast accuracy for daily-mean temperature anomalies and extremes across  
 469 various lead times and spatial scales. Forecast accuracy remains similar with  
 470 spatial aggregation, in contrast to precipitation, as indicated by the more  
 471 vertical skill contours and shorter lead-time gains relative to the grid-scale  
 472 (shading, compare Fig. S3 with Fig. 1). Weekly averaged temperature  
 473 anomalies and extremes display similar traits, with a forecast skill horizon  
 474 of 3 weeks across all spatial scales (Fig. S4). More spatially homogeneous  
 475 temperature fields compared to precipitation fields result in more accurate  
 476 forecasts at smaller scales, diminishing the benefits of spatial aggregation.

477 In summary, our results highlight and quantify the fundamental trade-off  
 478 between accuracy and precision in sub-seasonal precipitation forecasts. Spa-  
 479 tial aggregation, which reduces precision, increases forecast accuracy, extends  
 480 predictable lead times, and enhances maximum possible accuracy compared  
 481 to the grid scale. Conversely, increased precision tends to diminish these  
 482 benefits. This trade-off is more important at higher temporal precision (e.g.,  
 483 daily versus weekly aggregation), and for spatially inhomogeneous variables  
 484 (e.g., precipitation versus temperature). It is important to note that our  
 485 results are based on averages over hundreds of forecasts, but there are win-  
 486 dows of opportunity where individual forecasts can predict more accurately  
 487 and further ahead (Mariotti et al., 2020). While these findings are well-  
 488 recognized within the meteorological community (Buizza and Leutbecher,  
 489 2015; Toth and Buizza, 2019), they are often underappreciated by users who  
 490 could benefit from them, even potentially leading to the misuse of forecasts

491 (Nissan et al., 2019; Fiedler et al., 2021). In the next section, we demonstrate  
 492 how users can leverage these findings with three practical examples.

## 493 5. Use-Cases

494 In this section, we illustrate how spatial aggregation can make forecasts  
 495 more usable. Roberts and Lean (2008) proposed post-processing real-time  
 496 forecasts via spatial aggregation, optimizing spatial precision based on the  
 497 accuracy of historical forecasts. They envisioned a coarse forecast at longer  
 498 lead-times, becoming finer as the forecast horizon shortens and predictabil-  
 499 ity increases at smaller scales. However, they did not provide a practical  
 500 demonstration.

501 Here, we take their idea further and present three different use-cases for  
 502 spatial aggregation, schematically illustrated in Fig. 5: optimized accuracy  
 503 (blue arrow), fixed spatial precision (red arrow), and fixed lead time (black  
 504 arrow). These use cases are tailored to a smaller Scandinavian domain (see  
 505 Fig. 6a) and illustrate how the approach can be employed strategically to  
 506 optimize forecasts given specific user-defined constraints. Rather than re-  
 507 placing existing practices using grid-scale forecasts, our approach offers a  
 508 complementary perspective, and suggests avenues for further investigation in  
 509 each of the three examples presented.

### 510 5.1. Optimized Accuracy

511 Forecast accuracy often constrains how far ahead decisions can be made.  
 512 This is particularly true for forecasters, who need to issue early-warnings at  
 513 extended lead times. Using Storm Hans as an example, which first struck  
 514 Norway on August 7<sup>th</sup> 2023, we illustrate how spatial forecast aggregation  
 515 can give forecasters an earlier indication of extreme precipitation and thus  
 516 help them issue more timely early warnings.

517 Figure 6a–f shows forecasts of ensemble and daily-mean precipitation on  
 518 August 7<sup>th</sup> at various lead times, comparing grid-scale precision (left) and  
 519 spatially aggregated forecasts (right), while Fig. 6g shows the corresponding  
 520 observed precipitation at the grid-scale. Spatial aggregation is progressively  
 521 increased with lead time to increase forecast accuracy (blue arrow in Fig. 5).  
 522 Spatial aggregation could be used to maintain a constant forecast accuracy  
 523 by following along a specific contour (e.g., 0.7), but since these do not span  
 524 several lead-time days, here we ‘optimize’ accuracy by aggregating diagonally  
 525 across contours of constant accuracy. Regions with FEFI and EVI values

526 above 0.8, marked by red stippling, highlight areas of extreme precipitation  
527 in the forecasts and observations, respectively.

528 At the grid scale, forecasts capture the general pattern of precipitation  
529 over eastern Norway and Sweden, with the FEFI signaling extreme precip-  
530 itation reasonably well up to lead day 3 (see Fig. 6a,c,e with g). Spatially  
531 aggregated forecasts at lead days 1, 3 and 5 achieve accuracy comparable to  
532 grid-scale forecasts at lead day 1, as demonstrated by FMSESS and FBSS<sub>0.9</sub>  
533 values (contrast Fig. 6a,c,e with b,d,f). Most importantly, it is difficult to  
534 infer from the grid-scale forecast alone that localized extreme rainfall at lead  
535 day 5 would evolve into a widespread event across Scandinavia (compare red  
536 stippling Fig. 6a and 6g). In contrast, the spatially aggregated forecast at  
537 lead day 5 offers a clearer and earlier indication of the approaching large-  
538 scale extreme event (compare red stippling Fig. 6b and 6g). This is because  
539 spatial aggregation filters out small-scale noise and amplifies the larger-scale,  
540 predictable signal at longer lead times. Similar results are found up to lead  
541 day 7 (Fig. S5).

542 The Norwegian Meteorological Office issued a red alert for extreme rainfall  
543 on August 6<sup>th</sup>, just one day before Storm Hans struck Norway (Granerød  
544 et al., 2023). The decision to issue such alerts is guided by stringent internal  
545 procedures, balancing the need for timely warnings against the risk of false  
546 alarms. Our findings suggest that incorporating spatial aggregation into  
547 existing forecasting workflows could extend the lead time of early warnings  
548 for high-impact events like Storm Hans.

## 549 5.2. Fixed Spatial Precision

550 Forecast users are often constrained by the spatial scales at which they  
551 operate. In hydropower, for instance, the size of the watershed dictates  
552 both the volume of incoming precipitation and the timing of downstream  
553 impacts. When reservoirs near full capacity, Norwegian operators are often  
554 forced to discharge water in anticipation of rainfall events (NRK, 2020; TV2,  
555 2024) to avert flooding and infrastructure damage. However, such preemptive  
556 measures can be costly since the water could be used to generate higher-  
557 priced electricity at a different time. Leveraging the watershed’s size to  
558 refine precipitation forecasts could help operators optimize decision-making.

559 By spatially aggregating forecasts to match the watershed’s spatial ex-  
560 tent (red arrow in Fig. 5), operators could enhance forecast accuracy at the  
561 scale that matters most for their decisions. This refinement could enable hy-  
562 dropower managers to decide sooner, and with increased certainty, when and

563 how much to discharge water. Crucially, this approach goes beyond simply  
 564 spatially averaging the raw forecast over the watershed, since it maintains  
 565 the default grid-spacing but each grid cell aggregates data from its surround-  
 566 ing neighbors (e.g., Fig. 2). As a result, the method incorporates the spatial  
 567 uncertainty of precipitation that could occur near but not necessarily over  
 568 the watershed at longer lead times (e.g., as illustrated by Storm Hans in Fig.  
 569 6), thereby extracting a stronger predictable signal from the forecast. This  
 570 added granularity yields a more nuanced perspective of potential rainfall that  
 571 could ultimately feed into the basin.

### 572 5.3. *Fixed Lead-Time*

573 Fixed operational lead times can constrain how weather forecasts are ap-  
 574 plied. In commercial aviation, for instance, large-scale mid-latitude storms,  
 575 such as Storm Hans (Fig. 6), often disrupt flight schedules across multiple  
 576 airports and regions over extended periods. Because major carriers maintain  
 577 extensive route networks on tight schedules, they often decide two to three  
 578 days in advance whether to cancel, reroute or maintain flights given incoming  
 579 severe weather (NBC, 2011; nytimes, 2017). These early go/no-go choices are  
 580 critical for allocating resources and enabling passengers to arrange alterna-  
 581 tive travel plans. Delaying such arrangements significantly heightens the risk  
 582 of disruptions for passengers and additional operating costs.

583 Aggregating forecasts at broader spatial scales offers a practical means of  
 584 increasing their accuracy at these fixed lead times (black arrow in Fig. 5).  
 585 As demonstrated by the forecast maps in Fig. 6, airlines could synthesize  
 586 forecasts across larger geographic areas to gain a more robust sense of whether  
 587 a major winter system will develop at the time when critical decisions need  
 588 to be made. Although this aggregated approach drops fine-grained precision  
 589 at individual airports, it can strengthen confidence in the storm’s overall  
 590 footprint. This would facilitate timely and decisive operational adjustments  
 591 that span the geographically extensive networks of major carriers.

## 592 6. Conclusions and Discussion

593 Despite an abundance of available forecast data, much of it remains un-  
 594 derutilized, pointing to a critical usability gap (Lemos et al., 2012; Van den  
 595 Hurk et al., 2018; Findlater et al., 2021). In this study, we highlight and  
 596 quantify a crucial yet often overlooked challenge to their practical use: fore-  
 597 casts are often more precise than they are accurate when they are presented

598 on a denser grid spacing than the scales they can accurately predict. This  
599 mismatch stems from the loss of accuracy at smaller spatial scales as fore-  
600 casts extend further out in time (Lorenz, 1969; Toth and Buizza, 2019): a  
601 constraint known as the forecast skill horizon (Buizza et al., 2015; Buizza  
602 and Leutbecher, 2015). While many meteorologists recognize this trade-off,  
603 non-expert users may remain unaware, risking both the missed potential of  
604 predictability at larger scales and over-interpretation at finer scales. As fore-  
605 casts become increasingly important to support climate adaptation and pre-  
606 paredness, users and providers can benefit from recognizing and accounting  
607 for the trade-off between forecast accuracy and precision.

608 It is informative to consider this limitation through the lens of Murphy’s  
609 (1993) classic framework for “good” forecasts, which comprises three key  
610 measures: correspondence between the forecaster’s judgment and the de-  
611 livered forecast (consistency), how well the forecast corresponds to observed  
612 conditions (quality), and the practical benefit of the forecast to users (value).  
613 Although much attention focuses on how forecast quality impacts value, we  
614 argue that issuing forecasts at the default grid spacing, even when forecasters  
615 recognize diminished accuracy at those scales, reduces consistency. This lack  
616 of consistency, in turn, diminishes value: users unaware of the limitation may  
617 make suboptimal decisions, whereas those who are aware must contend with  
618 greater complexity in using the forecast. Indeed, Murphy (1993) anticipated  
619 this problem, noting that “forecasts and judgments may be inconsistent . . .  
620 in terms of their spatial and/or temporal specificity,” and called for practical  
621 solutions to address such mismatches.

622 Here, we modified the original fractions skill score to help users balance  
623 the trade-off between forecast precision and accuracy, transforming this met-  
624 ric from its traditional role in verifying spatial forecast accuracy for meteo-  
625 rologists into a tool for optimizing forecasts for end-users. We applied this  
626 approach to daily European precipitation forecasts, quantifying the balance  
627 for both deterministic predictions of anomalies and probabilistic predictions  
628 of extremes, using three years of sub-seasonal data from the European Cen-  
629 tre for Medium-Range Weather Forecasts (ECMWF). Our results show that  
630 decreasing precision through spatial aggregation increases forecast accuracy,  
631 extends predictable lead times, and enhances the maximum possible accuracy  
632 relative to the grid scale, while increased precision diminishes these benefits.

633 We hope that users will employ our approach to optimize forecasts for  
634 their specific application. Implementing this involves: 1) identifying the  
635 geographic region of interest, 2) verifying past forecasts with a chosen met-

636 ric, and 3) aggregating real-time forecasts accordingly. We demonstrated  
637 the practical value of our approach in three contexts: communicating early  
638 warnings, managing hydropower capacity, and commercial aviation plan-  
639 ning—each characterized by distinct user-constraints on accuracy, spatial  
640 scale, or lead-time. These use-cases showed that focusing on the scales where  
641 forecasts are most accurate, rather than the default grid-scale, can offer users  
642 more actionable information. Instead of replacing existing practices, our ap-  
643 proach offers a complementary perspective, and highlights multiple avenues  
644 for further investigation in each of the three examples.

645 Aggregating forecasts is a well-established practice, yet its implementa-  
646 tion is often done by forecast providers rather than end-users. For example,  
647 the “ready-set-go” framework (Goddard et al., 2014) links forecasts across  
648 timescales to general preparedness levels, from monthly seasonal predictions  
649 (ready), to weekly sub-seasonal forecasts (set), and finally to daily weather  
650 forecasts (go), with each system having finer spatial and temporal precision  
651 as the forecast window shortens. Another approach is to filter forecasts into  
652 a few distinct large-scale patterns that are more predictable, called weather  
653 regimes (Michelangeli et al., 1995), as done operationally by ECMWF. Both  
654 these approaches implicitly involve aggregation, but the scales are set by ei-  
655 ther the modeling system or the regime classification, not the users. Our  
656 approach builds on these approaches but takes a step further, giving the user  
657 the ability to tailor forecasts according to their desired spatial scales.

658 Machine learning models hold considerable promise to improve both fore-  
659 cast accuracy and precision, potentially narrowing the usability gap (Eyring  
660 et al., 2024). However, these new approaches are likely not a panacea, and  
661 are subject to similar physical constraints which limit conventional models  
662 (Ben Bouallègue et al., 2024). It is widely acknowledged that there are fun-  
663 damental limits to the forecast skill horizon (Lorenz, 1969; Palmer et al.,  
664 2014), and emerging evidence suggests that machine learning-based forecasts  
665 are not exempt from this constraint (Keane et al. 2025, in review). So, even  
666 if machine learning models produce ever finer-scale forecasts, users are still  
667 likely to face the challenge posed by the forecast skill horizon, and to require  
668 methods to deal with it, as discussed here.

669 Our findings show that spatial aggregation enhances the accuracy of  
670 daily-accumulated precipitation forecasts to a greater extent than weekly-  
671 accumulated ones, per unit lead time. This suggests that temporal ag-  
672 gregation compensates to some extent for spatial aggregation, which has  
673 implications for how forecasts are used and communicated. For example,

674 sub-seasonal forecasts are often presented as weekly aggregates ([https://](https://charts.ecmwf.int/)  
675 [charts.ecmwf.int/](https://charts.ecmwf.int/)), leveraging time averaging to enhance accuracy at ex-  
676 tended lead times. Correspondingly, intuition based on the forecast skill  
677 horizon suggests that these predictions should be interpreted on spatial scales  
678 larger than the grid-scale. However, our results for both precipitation and  
679 temperature show weekly-aggregated forecasts have similar accuracy across  
680 spatial scales (Fig. 3 and S4). Thus, the decision not to spatially aggregate  
681 these forecasts is sound; grid-scale forecasts appear to be more usable than  
682 we expected, provided temporal aggregation is applied.

683 Further refinements to our approach could enhance its usability and  
684 broaden its relevance. For instance, employing alternative fractions skill  
685 scores, such as those proposed by Necker et al. (2024), may better assess  
686 probabilistic spatial forecast accuracy. Accounting for geographic variations  
687 in forecast skill, as shown in Fig. 2, would further improve adaptability across  
688 regions and applications. While refining the domain to a smaller area of in-  
689 terest is a straightforward solution, other approaches are possible, such as ap-  
690 plying different levels of spatial aggregation in different regions of the domain  
691 as suggested by Roberts and Lean (2008). Beyond accuracy, incorporating  
692 other metrics of forecast quality, such as reliability and discrimination, could  
693 be more relevant to users (Murphy, 1993; Weisheimer and Palmer, 2014). In  
694 practice, it is possible to quantify these metrics on the same axes as Fig. 5,  
695 swapping out forecast accuracy for alternatives. Finally, the fractions ex-  
696 treme forecast index could be further modified to weight the most extreme  
697 quantiles, changing its sensitivity to severe events.

698 Our results suggest that spatially aggregating weather forecasts could en-  
699 hance the accuracy of downstream impact models, which are increasingly  
700 important for climate adaptation and preparedness (Merz et al., 2020). Op-  
701 timizing their inputs might be simpler and more effective than improving  
702 their basic design. Implementing this approach, however, poses technical  
703 challenges. Physics-based impact models, such as those producing hydro-  
704 logical forecasts, require physical consistency between input variables (e.g.,  
705 rainfall and temperature) which is disrupted by aggregation. On the other-  
706 hand, data-driven impact models, which are not constrained by physical laws  
707 and are being used in a variety of sectors including insurance, agriculture,  
708 and even hydrology, could be trained to capture the relationship between  
709 aggregated weather forecasts and impacts. Exploring the feasibility of this  
710 approach could be an interesting avenue for further research.

## 711 Acknowledgments

712 E.D.S, E.K. and O.W. were supported by the Research Council of Norway  
713 grant 309562 Climate Futures, and E.K. was also supported by the Eu-  
714 ropean Union grant 101137847 ACACIA. R.K. was supported by the Centre  
715 for Environmental Modelling and Computation (CEMAK) at the Univer-  
716 sity of Leeds. D.P. was supported by the Norad-funded project RAF:23/006  
717 Improving Smallholder Resilience through Customised Climate Services and  
718 the WISER Early Warnings for Southern Africa (WISER-EWSA) project,  
719 funded by the Met Office as part of the Weather and Climate Information  
720 Services (WISER) programme on behalf UK government’s Foreign, Common-  
721 wealth and Development Office (FCDO). The computations were performed  
722 on resources provided by the Sigma2 NS9873K project, as part of the Na-  
723 tional Infrastructure for High-Performance Computing and Data Storage in  
724 Norway. E.D.S acknowledges Håkon Otneim, Sondre Holleland and Bjørnar  
725 Jensen for feedback on an early version of the manuscript and Thordis Tho-  
726 rarinsdottir for early discussions which contributed to the main idea for the  
727 manuscript.

## 728 Data statement

729 Sub-seasonal forecasts from ECWMF (ECMWF, 2024a) and ERA5 re-  
730 analysis data (Hersbach et al., 2023) are freely available from the MARS  
731 archive (<https://apps.ecmwf.int/archive-catalogue/>) and the Coperni-  
732 cus Climate Data Store ([https://cds.climate.copernicus.eu/datasets/](https://cds.climate.copernicus.eu/datasets/reanalysis-era5-single-levels?tab=overview)  
733 [reanalysis-era5-single-levels?tab=overview](https://cds.climate.copernicus.eu/datasets/reanalysis-era5-single-levels?tab=overview)), respectively. Materials  
734 to reproduce the figures in this paper are provided on Github ([https://](https://github.com/edunnsigouin/accuracyvsprecision)  
735 [github.com/edunnsigouin/accuracyvsprecision](https://github.com/edunnsigouin/accuracyvsprecision)).

## 736 Declaration of generative AI and AI-assisted technologies in the 737 writing process

738 During the preparation of this work the author used ChatGPT to assist  
739 with the writing of the manuscript. After using this tool, the author reviewed  
740 and edited the content as needed and takes full responsibility for the content  
741 of the published article.



## References

- Bauer, P., Thorpe, A., Brunet, G., 2015. The quiet revolution of numerical weather prediction. *Nature* 525, 47–55. doi:10.1038/nature14956.
- Ben Bouallègue, Z., Clare, M.C., Magnusson, L., Gascon, E., Maier-Gerber, M., Janoušek, M., Rodwell, M., Pinault, F., Dramsch, J.S., Lang, S.T., et al., 2024. The rise of data-driven weather forecasting: A first statistical assessment of machine learning-based weather forecasts in an operational-like context. *Bulletin of the American Meteorological Society* 105, E864–E883. doi:10.1175/BAMS-D-23-0162.1.
- Benjamin, S.G., Brown, J.M., Brunet, G., Lynch, P., Saito, K., Schlatter, T.W., 2019. 100 years of progress in forecasting and nwp applications. *Meteorological Monographs* 59, 13–1. doi:10.1175/AMSMONOGRAPHS-D-18-0020.1.
- Buizza, R., Balmaseda, M.A., Brown, A., English, S., Forbes, R., Geer, A., Haiden, T., Leutbecher, M., Magnusson, L., Rodwell, M., et al., 2018. The development and evaluation process followed at ECMWF to upgrade the Integrated Forecasting System (IFS). *European Centre for Medium Range Weather Forecasts*. doi:10.21957/xzopnhty9.
- Buizza, R., Leutbecher, M., 2015. The forecast skill horizon. *Quarterly Journal of the Royal Meteorological Society* 141, 3366–3382. doi:10.1002/qj.2619.
- Buizza, R., Leutbecher, M., Thorpe, A., 2015. Living with the butterfly effect: A seamless view of predictability. *ECMWF newsletter* 145, 18–23. doi:10.21957/x4h3e8w3.
- Cafaro, C., Woodhams, B.J., Stein, T.H., Birch, C.E., Webster, S., Bain, C.L., Hartley, A., Clarke, S., Ferrett, S., Hill, P., 2021. Do convection-permitting ensembles lead to more skillful short-range probabilistic rainfall forecasts over tropical east africa? *Weather and Forecasting* 36, 697–716. doi:10.1175/WAF-D-20-0172.1.
- Cusick, D., 2019. Tech offers a virtual window into future climate change risk. URL: <https://www.scientificamerican.com/article/tech-offers-a-virtual-window-into-future-climate-change-risk/>. (Accessed: 13.08.2024).

775 ECMWF, 2024a. Meteorological archival and retrieval system (mars).  
 776 URL: [https://www.ecmwf.int/en/forecasts/access-forecasts/](https://www.ecmwf.int/en/forecasts/access-forecasts/access-archive-datasets)  
 777 [access-archive-datasets](https://www.ecmwf.int/en/forecasts/access-forecasts/access-archive-datasets). (Accessed: 13.08.2024).

778 ECMWF, 2024b. Section 5.3.1 m-climate, the ens model climate.  
 779 URL: [https://confluence.ecmwf.int/display/FUG/Section+5.3.1+](https://confluence.ecmwf.int/display/FUG/Section+5.3.1+M-climate%2C+the+ENS+Model+Climate)  
 780 [M-climate%2C+the+ENS+Model+Climate](https://confluence.ecmwf.int/display/FUG/Section+5.3.1+M-climate%2C+the+ENS+Model+Climate). (Accessed: 16.08.2024).

781 ECMWF, 2024c. Section 8.1.9.2 extreme forecast index - efi.  
 782 URL: [https://confluence.ecmwf.int/display/FUG/Section+8.1.9.](https://confluence.ecmwf.int/display/FUG/Section+8.1.9.2+Extreme+Forecast+Index+-+EFI)  
 783 [2+Extreme+Forecast+Index+-+EFI](https://confluence.ecmwf.int/display/FUG/Section+8.1.9.2+Extreme+Forecast+Index+-+EFI). (Accessed: 16.08.2024).

784 Ekroll, H.C., 2023. Ekstremværet hans vil trolig koste langt mer enn stor-  
 785 flommen i 1995. URL: [https://www.aftenposten.no/norge/i/76w3oW/](https://www.aftenposten.no/norge/i/76w3oW/ekstremvaeret-hans-vil-trolig-koste-langt-mer-enn-storflommen-i-1995)  
 786 [ekstremvaeret-hans-vil-trolig-koste-langt-mer-enn-storflommen-i-1995](https://www.aftenposten.no/norge/i/76w3oW/ekstremvaeret-hans-vil-trolig-koste-langt-mer-enn-storflommen-i-1995).  
 787 (Accessed: 13.08.2024).

788 European-Commission, 2020. Eu taxonomy for sustainable activi-  
 789 ties. URL: [https://finance.ec.europa.eu/sustainable-finance/](https://finance.ec.europa.eu/sustainable-finance/tools-and-standards/eu-taxonomy-sustainable-activities_en)  
 790 [tools-and-standards/eu-taxonomy-sustainable-activities\\_en](https://finance.ec.europa.eu/sustainable-finance/tools-and-standards/eu-taxonomy-sustainable-activities_en).  
 791 (Accessed: 13.08.2024).

792 Eyring, V., Gentine, P., Camps-Valls, G., Lawrence, D.M., Reichstein,  
 793 M., 2024. Ai-empowered next-generation multiscale climate modelling  
 794 for mitigation and adaptation. *Nature Geoscience* , 1–9doi:10.1038/  
 795 [s41561-024-01527-w](https://doi.org/10.1038/s41561-024-01527-w).

796 Fiedler, T., Pitman, A.J., Mackenzie, K., Wood, N., Jakob, C., Perkins-  
 797 Kirkpatrick, S.E., 2021. Business risk and the emergence of cli-  
 798 mate analytics. *Nature Climate Change* 11, 87–94. doi:10.1038/  
 799 [s41558-020-00984-6](https://doi.org/10.1038/s41558-020-00984-6).

800 Findlater, K., Webber, S., Kandlikar, M., Donner, S., 2021. Climate services  
 801 promise better decisions but mainly focus on better data. *Nature Climate*  
 802 *Change* 11, 731–737. doi:10.1038/s41558-021-01125-3.

803 Gehne, M., Hamill, T.M., Kiladis, G.N., Trenberth, K.E., 2016. Comparison  
 804 of global precipitation estimates across a range of temporal and spatial  
 805 scales. *Journal of Climate* 29, 7773–7795. doi:10.1175/JCLI-D-15-0618.  
 806 1.

- 807 Gilleland, E., Ahijevych, D., Brown, B.G., Casati, B., Ebert, E.E., 2009. In-  
808 tercomparison of spatial forecast verification methods. *Weather and fore-*  
809 *casting* 24, 1416–1430. doi:10.1175/2009WAF2222269.1.
- 810 Goddard, L., 2016. From science to service. *Science* 353, 1366–1367. doi:10.  
811 1126/science.aag3087.
- 812 Goddard, L., Baethgen, W.E., Bhojwani, H., Robertson, A.W., 2014. The  
813 international research institute for climate & society: why, what and how.  
814 *Earth Perspectives* 1, 1–14. doi:10.1186/2194-6434-1-10.
- 815 Gong, X., Barnston, A.G., Ward, M.N., 2003. The effect of spatial aggrega-  
816 tion on the skill of seasonal precipitation forecasts. *Journal of Climate* 16,  
817 3059–3071. doi:10.1175/1520-0442(2003)016<3059:TEOSA0>2.0.CO;2.
- 818 Graham, R.M., Browell, J., Bertram, D., White, C.J., 2022. The application  
819 of sub-seasonal to seasonal (s2s) predictions for hydropower forecasting.  
820 *Meteorological Applications* 29, e2047. doi:10.1002/met.2047.
- 821 Granerød, M., Stabell, D., Mjelstad, H., Tajet, H., 2023. Ek-  
822 stremværet’hans’, ekstremt mye nedbør i deler av sør-norge 07.-09. august  
823 2023. The Norwegian Meteorological Institute (ed) METinfo 26.
- 824 Hanssen-Bauer, I., Drange, H., Førland, E., Roald, L., Børsheim, K., His-  
825 dal, H., Lawrence, D., Nesje, A., Sandven, S., Sorteberg, A., et al.,  
826 2009. Climate in norway 2100. Background information to NOU Cli-  
827 mate adaptation (In Norwegian: Klima i Norge 2100. Bakgrunnsmate-  
828 riale til NOU Klimatilplassing), Oslo: Norsk klimasenter URL: [https:](https://klimaservicesenter.no/kss/rapporter/kin2100)  
829 [//klimaservicesenter.no/kss/rapporter/kin2100](https://klimaservicesenter.no/kss/rapporter/kin2100).
- 830 Haupt, S.E., Hanna, S., Askelson, M., Shepherd, M., Fragomeni, M.A., Deb-  
831 bage, N., Johnson, B., 2019a. 100 years of progress in applied meteorology.  
832 part ii: Applications that address growing populations. *Meteorological*  
833 *Monographs* 59, 23–1. doi:10.1175/AMSMONOGRAPHS-D-18-0007.1.
- 834 Haupt, S.E., Kosović, B., McIntosh, S.W., Chen, F., Miller, K., Shepherd,  
835 M., Williams, M., Drobot, S., 2019b. 100 years of progress in applied  
836 meteorology. part iii: Additional applications. *Meteorological Monographs*  
837 59, 24–1. doi:10.1175/AMSMONOGRAPHS-D-18-0012.1.

- 838 Haupt, S.E., Rauber, R.M., Carmichael, B., Knievel, J.C., Cogan, J.L.,  
839 2018. 100 years of progress in applied meteorology. part i: Ba-  
840 sic applications. *Meteorological Monographs* 59, 22–1. doi:10.1175/  
841 AMSMONOGRAPHS-D-18-0004.1.
- 842 Hersbach, H., Bell, B., Berrisford, P., G., B., Horányi, A., Muñoz Sabater, J.,  
843 Nicolas, J., Peubey, C., Radu, R., Rozum, I., Schepers, D., Simmons, A.,  
844 Soci, C., Dee, D., Thépaut, J.N., 2023. ERA5 hourly data on single levels  
845 from 1940 to present. Copernicus Climate Change Service (C3S) Climate  
846 Data Store (CDS) doi:10.24381/cds.adbb2d47.
- 847 Hewitt, C., Mason, S., Walland, D., 2012. The global framework for climate  
848 services. *Nature Climate Change* 2, 831–832. doi:10.1038/nclimate1745.
- 849 Van den Hurk, B., Hewitt, C., Jacob, D., Bessembinder, J., Doblas-Reyes,  
850 F., Döscher, R., 2018. The match between climate services demands and  
851 earth system models supplies. *Climate Services* 12, 59–63. doi:10.1016/  
852 j.cliser.2018.11.002.
- 853 Jolliffe, I.T., Stephenson, D.B., 2012. Forecast verification: a practi-  
854 tioner’s guide in atmospheric science. John Wiley & Sons. doi:10.1002/  
855 9781119960003.
- 856 Jung, T., Leutbecher, M., 2008. Scale-dependent verification of ensemble  
857 forecasts. *Quarterly Journal of the Royal Meteorological Society* 134, 973–  
858 984. doi:https://doi.org/10.1002/qj.255.
- 859 Keane, R.J., Plant, R.S., Tennant, W.J., 2016. Evaluation of the plant–craig  
860 stochastic convection scheme (v2. 0) in the ensemble forecasting system  
861 mogreps-r (24 km) based on the unified model (v7. 3). *Geoscientific Model*  
862 *Development* 9, 1921–1935. doi:10.5194/gmd-9-1921-2016.
- 863 KLP, K.L.G.F., 2023. Dette lærte kommunene av ek-  
864 stremværet hans. URL: [https://www.klp.no/artikler/  
865 dette-laerte-kommunene-av-ekstremvaeret-hans](https://www.klp.no/artikler/dette-laerte-kommunene-av-ekstremvaeret-hans). (Accessed:  
866 13.08.2024).
- 867 Lalaurette, F., 2003. Early detection of abnormal weather conditions using  
868 a probabilistic extreme forecast index. *Quarterly Journal of the Royal*  
869 *Meteorological Society* 129, 3037–3057. doi:10.1256/qj.02.152.

- 870 Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirsberger, P., Fortunato, M.,  
871 Alet, F., Ravuri, S., Ewalds, T., Eaton-Rosen, Z., Hu, W., et al., 2023.  
872 Learning skillful medium-range global weather forecasting. *Science* 382,  
873 1416–1421. doi:10.1126/science.adi233.
- 874 Lavers, D.A., Simmons, A., Vamborg, F., Rodwell, M.J., 2022. An evaluation  
875 of ERA5 precipitation for climate monitoring. *Quarterly Journal of the*  
876 *Royal Meteorological Society* 148, 3152–3165. doi:10.1002/qj.4351.
- 877 Lemos, M.C., Kirchhoff, C.J., Ramprasad, V., 2012. Narrowing the climate  
878 information usability gap. *Nature climate change* 2, 789–794. doi:10.  
879 1038/nclimate1614.
- 880 Lorenz, E.N., 1969. The predictability of a flow which possesses many scales  
881 of motion. *Tellus* 21, 289–307. doi:10.1111/j.2153-3490.1969.tb00444.  
882 x.
- 883 Mariotti, A., Baggett, C., Barnes, E.A., Becker, E., Butler, A., Collins,  
884 D.C., Dirmeyer, P.A., Ferranti, L., Johnson, N.C., Jones, J., et al., 2020.  
885 Windows of opportunity for skillful forecasts subseasonal to seasonal and  
886 beyond. *Bulletin of the American Meteorological Society* 101, E608–E625.  
887 doi:10.1175/BAMS-D-18-0326.1.
- 888 Merryfield, W.J., Baehr, J., Batté, L., Becker, E.J., Butler, A.H., Coelho,  
889 C.A., Danabasoglu, G., Dirmeyer, P.A., Doblas-Reyes, F.J., Domeisen,  
890 D.I., et al., 2020. Current and emerging developments in subseasonal to  
891 decadal prediction. *Bulletin of the American Meteorological Society* 101,  
892 E869–E896. doi:10.1175/BAMS-D-19-0037.1.
- 893 Merz, B., Kuhlicke, C., Kunz, M., Pittore, M., Babeyko, A., Bresch, D.N.,  
894 Domeisen, D.I., Feser, F., Koszalka, I., Kreibich, H., et al., 2020. Impact  
895 forecasting to support emergency management of natural hazards. *Reviews*  
896 *of Geophysics* 58, e2020RG000704. doi:10.1029/2020RG000704.
- 897 Michelangeli, P.A., Vautard, R., Legras, B., 1995. Weather regimes: Re-  
898 currence and quasi stationarity. *Journal of the atmospheric sciences* 52,  
899 1237–1256. doi:10.1175/1520-0469(1995)052<1237:WRRQ>2.0.CO;2.
- 900 Murphy, A.H., 1993. What is a good forecast? an essay on the nature  
901 of goodness in weather forecasting. *Weather and forecasting* 8, 281–293.  
902 doi:10.1175/1520-0434(1993)008<0281:WIAGFA>2.0.CO;2.

903 NBC, 2011. Airlines cancel flights before snow flies. URL: <https://www.nbcnews.com/id/wbna41502385>. (Accessed: 09.01.2025).  
904

905 Necker, T., Wolfgruber, L., Kugler, L., Weissmann, M., Dorninger, M., Serafin, S., 2024. The fractions skill score for ensemble forecast verification. Quarterly Journal of the Royal Meteorological Society doi:10.1002/qj.4824.  
906  
907  
908

909 Nissan, H., Goddard, L., de Perez, E.C., Furlow, J., Baethgen, W., Thomson, M.C., Mason, S.J., 2019. On the use and misuse of climate change projections in international development. Wiley Interdisciplinary Reviews: Climate Change 10, e579. doi:10.1002/wcc.579.  
910  
911  
912

913 NRK, 2020. Overfylte vannmagasin gir rekordlave strømpriser: – får betalt for å bruke strøm. URL: [https://www.nrk.no/vestland/overfylte-vannmagasin-gir-rekordlave-strompriser\\_-\\_far-betalt-for-a-bruke-strom-1.15110896](https://www.nrk.no/vestland/overfylte-vannmagasin-gir-rekordlave-strompriser_-_far-betalt-for-a-bruke-strom-1.15110896). (Accessed: 08.01.2025).  
914  
915  
916

917 nytimes, 2017. Airlines take a proactive approach to potential weather woes. URL: <https://www.nytimes.com/2017/03/14/business/airline-cancellations-delays-snowstorm.html>. (Accessed: 09.01.2025).  
918  
919  
920

921 O’Kane, T.J., Scaife, A.A., Kushnir, Y., Brookshaw, A., Buontempo, C., Carlin, D., Connell, R.K., Doblas-Reyes, F., Dunstone, N., Förster, K., et al., 2023. Recent applications and potential of near-term (interannual to decadal) climate predictions. Frontiers in Climate 5, 1121626. doi:10.3389/fclim.2023.1121626.  
922  
923  
924  
925

926 Palmer, T., Döring, A., Seregin, G., 2014. The real butterfly effect. Nonlinearity 27, R123. doi:10.1088/0951-7715/27/9/R123.  
927

928 Coughlan de Perez, E., Harrison, L., Berse, K., Easton-Calabria, E., Marunye, J., Marake, M., Murshed, S.B., Zauisomue, E.H., et al., 2022. Adapting to climate change through anticipatory action: The potential use of weather-based early warnings. Weather and Climate Extremes 38, 100508. doi:10.1016/j.wace.2022.100508.  
929  
930  
931  
932

933 Price, I., Sanchez-Gonzalez, A., Alet, F., Andersson, T.R., El-Kadi, A., Masters, D., Ewalds, T., Stott, J., Mohamed, S., Battaglia, P., et al., 2024.  
934

- 935 Probabilistic weather forecasting with machine learning. *Nature* 637, 84–  
936 90. doi:10.1038/s41586-024-08252-9.
- 937 Rivoire, P., Martius, O., Naveau, P., Tuel, A., 2023. Assessment of  
938 subseasonal-to-seasonal (s2s) ensemble extreme precipitation forecast skill  
939 over europe. *Natural Hazards and Earth System Sciences* 23, 2857–2871.  
940 doi:10.5194/nhess-23-2857-2023.
- 941 Roberts, N.M., Lean, H.W., 2008. Scale-selective verification of rainfall ac-  
942 cumulations from high-resolution forecasts of convective events. *Monthly*  
943 *Weather Review* 136, 78–97. doi:10.1175/2007MWR2123.1.
- 944 Röösl, T., Appenzeller, C., Bresch, D.N., 2021. Towards operational impact  
945 forecasting of building damage from winter windstorms in switzerland.  
946 *Meteorological applications* 28, e2035. doi:10.1002/met.2035.
- 947 Schwartz, C.S., 2019. Medium-range convection-allowing ensemble forecasts  
948 with a variable-resolution global model. *Monthly Weather Review* 147,  
949 2997–3023. doi:10.1175/MWR-D-18-0452.1.
- 950 Torralba, V., Doblas-Reyes, F.J., MacLeod, D., Christel, I., Davis, M., 2017.  
951 Seasonal climate prediction: a new source of information for the man-  
952 agement of wind energy resources. *Journal of Applied Meteorology and*  
953 *Climatology* 56, 1231–1247. doi:10.1175/JAMC-D-16-0204.1.
- 954 Toth, Z., Buizza, R., 2019. Weather forecasting: What sets the forecast skill  
955 horizon?, in: *Sub-Seasonal to Seasonal Prediction*. Elsevier, pp. 17–45.  
956 doi:10.1016/B978-0-12-811714-9.00002-4.
- 957 Trenberth, K.E., Marquis, M., Zebiak, S., 2016. The vital need for a climate  
958 information system. *Nature Climate Change* 6, 1057–1059. doi:10.1038/  
959 nclimate3170.
- 960 TV2, 2024. Ikke opplevd maken på 50 år: – ek-  
961 stremt. URL: [https://www.tv2.no/nyheter/innenriks/](https://www.tv2.no/nyheter/innenriks/ikke-opplevd-maken-pa-50-ar-ekstremt/17173589/)  
962 [ikke-opplevd-maken-pa-50-ar-ekstremt/17173589/](https://www.tv2.no/nyheter/innenriks/ikke-opplevd-maken-pa-50-ar-ekstremt/17173589/). (Accessed:  
963 08.01.2025).
- 964 Van Straaten, C., Whan, K., Coumou, D., van den Hurk, B., Schmeits, M.,  
965 2020. The influence of aggregation and statistical post-processing on the

- 966 subseasonal predictability of european temperatures. Quarterly Journal of  
967 the Royal Meteorological Society 146, 2654–2670. doi:10.1002/qj.3810.
- 968 Weisheimer, A., Palmer, T., 2014. On the reliability of seasonal climate  
969 forecasts. Journal of the Royal Society Interface 11, 20131162. doi:10.  
970 1098/rsif.2013.1162.
- 971 White, C.J., Domeisen, D.I., Acharya, N., Adefisan, E.A., Anderson, M.L.,  
972 Aura, S., Balogun, A.A., Bertram, D., Bluhm, S., Brayshaw, D.J., et al.,  
973 2022. Advances in the application and utility of subseasonal-to-seasonal  
974 predictions. Bulletin of the American Meteorological Society 103, E1448–  
975 E1472. doi:10.1175/BAMS-D-20-0224.1.
- 976 WMO, 2022. Wmo and the early warnings for all initiative.  
977 URL: [https://wmo.int/activities/early-warnings-all/  
978 wmo-and-early-warnings-all-initiative](https://wmo.int/activities/early-warnings-all/wmo-and-early-warnings-all-initiative). (Accessed: 16.08.2024).
- 979 Young, M., Heinrich, V., Black, E., Asfaw, D., 2020. Optimal spatial scales  
980 for seasonal forecasts over africa. Environmental Research Letters 15,  
981 094023. doi:10.1088/1748-9326/ab94e9.
- 982 Zhao, B., Zhang, B., 2018. Assessing hourly precipitation forecast skill with  
983 the fractions skill score. Journal of Meteorological Research 32, 135–145.  
984 doi:10.1007/s13351-018-7058-1.



Table 1: Overview of variables and their dimensions defined in section 3

Dimension/Variable	Description	section
$m$	forecast initialization date	3
$e$	ensemble member	3
$t$	lead-time day	3
$i$	latitude	3
$j$	longitude	3
$y$	hindcast year	3
$n$	spatial aggregation level	3
$q$	quantile	3
$f(m, e, t, i, j)$	forecast	3
$h(m, y, e, t, i, j)$	hindcast	3
$v_f(m, t, i, j)$	forecast verification	3
$v_h(m, y, t, i, j)$	hindcast verification	3
$F(n, m, t, i, j)$	aggregated forecast	3.1
$V_F(n, m, t, i, j)$	aggregated forecast verification	3.1
$MSE_F(n, m, t)$	mean-square error of aggregated forecast	3.1
$MSE_{REF}(n, t)$	mean-square error of aggregated forecast relative to reference forecast	3.1
$FSS(n, m, t)$	fractions skill score	3.1
$\tilde{f}(m, t, i, j)$	forecast anomaly	3.2.1
$\tilde{v}_f(m, t, i, j)$	forecast verification anomaly	3.2.1
$\tilde{F}(m, t, i, j)$	aggregated forecast anomaly	3.2.1
$\tilde{V}_F(m, t, i, j)$	aggregated forecast verification anomaly	3.2.1
$MSE_{\tilde{F}}(n, m, t)$	mean-square error of aggregated forecast anomalies	3.2.1
$MSE_{\tilde{REF}}(n, m, t)$	mean-square error of aggregated forecast anomalies relative to reference forecast	3.2.1
$FMSESS(n, t)$	fractions mean-square error skill score	3.2.1
$H(n, m, t, i, j)$	aggregated hindcast	3.2.2
$V_H(n, m, t, i, j)$	aggregated hindcast verification	3.2.2
$F_q(n, m, t, i, j)$	aggregated forecast threshold value for quantile q	3.2.2
$V_q(n, m, t, i, j)$	aggregated verification threshold value for quantile q	3.2.2
$P_{F_q}(n, m, t, i, j)$	aggregated forecast probability for quantile q	3.2.2
$P_{V_q}(n, m, t, i, j)$	aggregated binary verification for quantile q	3.2.2
$BS_q(n, m, t)$	brier-score of aggregated forecasts for quantile q	3.2.2
$BS_{REF_q}(n, m, t)$	brier-score of aggregated forecasts for quantile q relative to reference forecast	3.2.2
$FBSS_q(n, t)$	fractions brier skill score for quantile q	3.2.2
$h_{thresh}(q, m, t, i, j)$	hindcast threshold value	3.3
$fr_f(q, m, t, i, j)$	fraction of forecast ensemble members over threshold	3.3
$EFI(m, t, i, j)$	extreme forecast index	3.3
$v_{h_{thresh}}(q, m, t, i, j)$	verification hindcast threshold value	3.3
$fr_v(q, m, t, i, j)$	binary verification over threshold	3.3
$EVI(m, t, i, j)$	extreme verification index	3.3
$H_{thresh}(n, q, m, t, i, j)$	aggregated hindcast threshold value	3.3
$FR_F(n, q, m, t, i, j)$	fraction of aggregated forecast ensemble members over threshold	3.3
$FEFI(n, m, t, i, j)$	fractions extreme forecast index	3.3

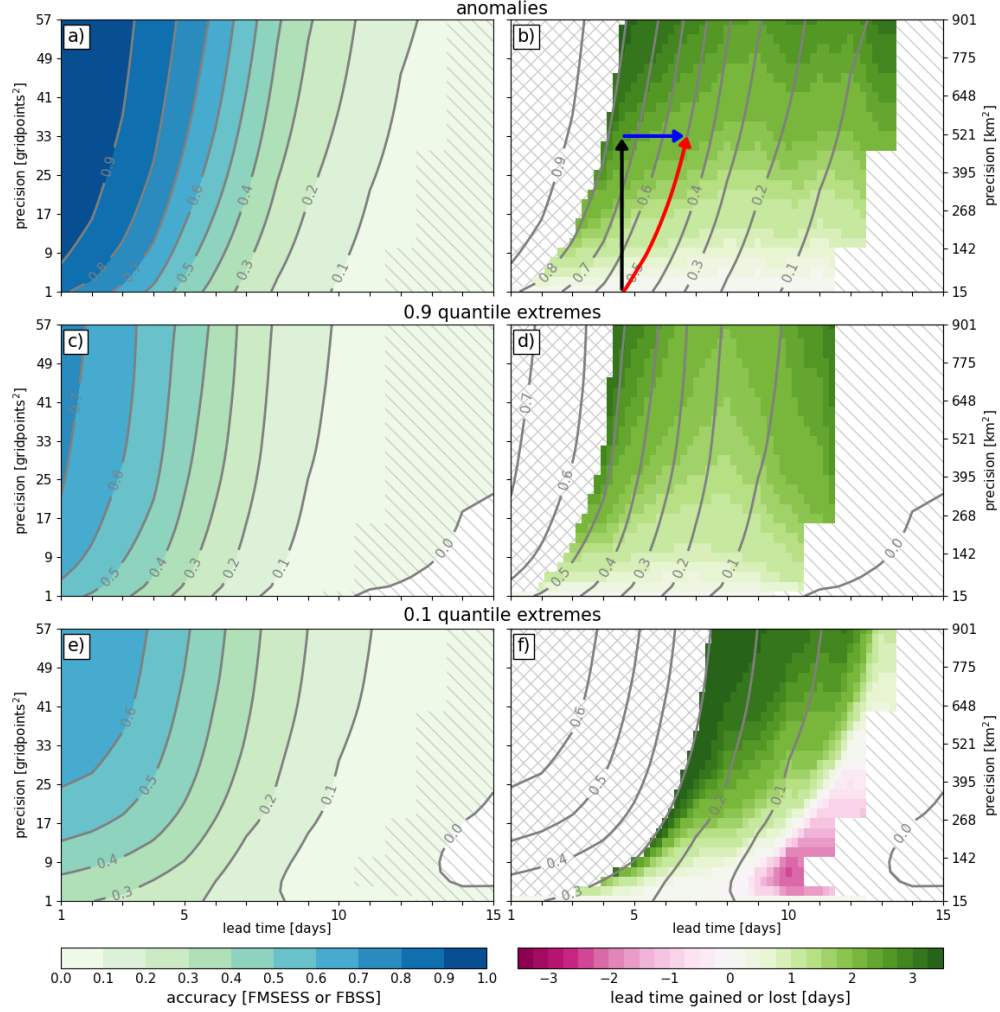


Figure 1: Forecast accuracy and lead-time gained for daily-accumulated precipitation over Europe as a function of lead time and spatial precision. Left panels show the Fractions Mean-Square Error Skill Score (FMSESS) for anomalies (a) and the Fractions Brier Skill Score (FBSS) for 0.9 and 0.1 quantile extremes in (c) and (e), respectively. Hatching indicates skill scores that are not statistically significant at the 5% level assessed via bootstrapping. Cross-hatching denotes scores surpassing the highest accuracy obtained at the grid scale. Right panels follow the same conventions as the left ones, except that shading represents the lead time gained or lost by spatially aggregating the forecasts. Black, blue and red arrows in (b) illustrate the lead-time gained when spatially aggregating an example forecast at the grid-scale. See Section 4 for further details.

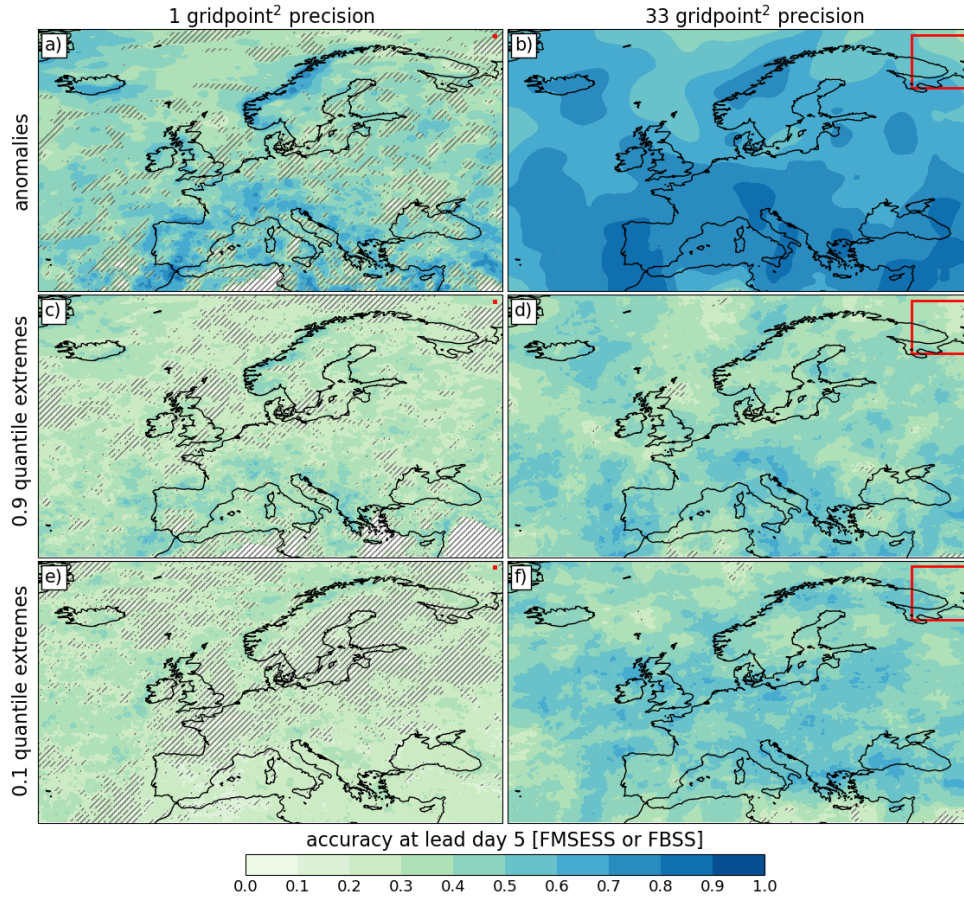


Figure 2: Forecast accuracy at lead-day 5 for daily-accumulated precipitation over Europe, evaluated at two levels of spatial precision: the grid-scale (left) and 33 gridpoints<sup>2</sup> (right). Shading denotes the Fractions Mean-Square Error Skill Score (FMSESS) for anomalies in (a,b) and the Fractions Brier Skill Score (FBSS) for the 0.9 and 0.1 quantile extremes in (c,d) and (e,f), respectively. Hatching indicates skill scores that are not statistically significant at the 5% level assessed using bootstrapping. Red squares in the upper-right corner of the panels denote the spatial scale of forecast precision.

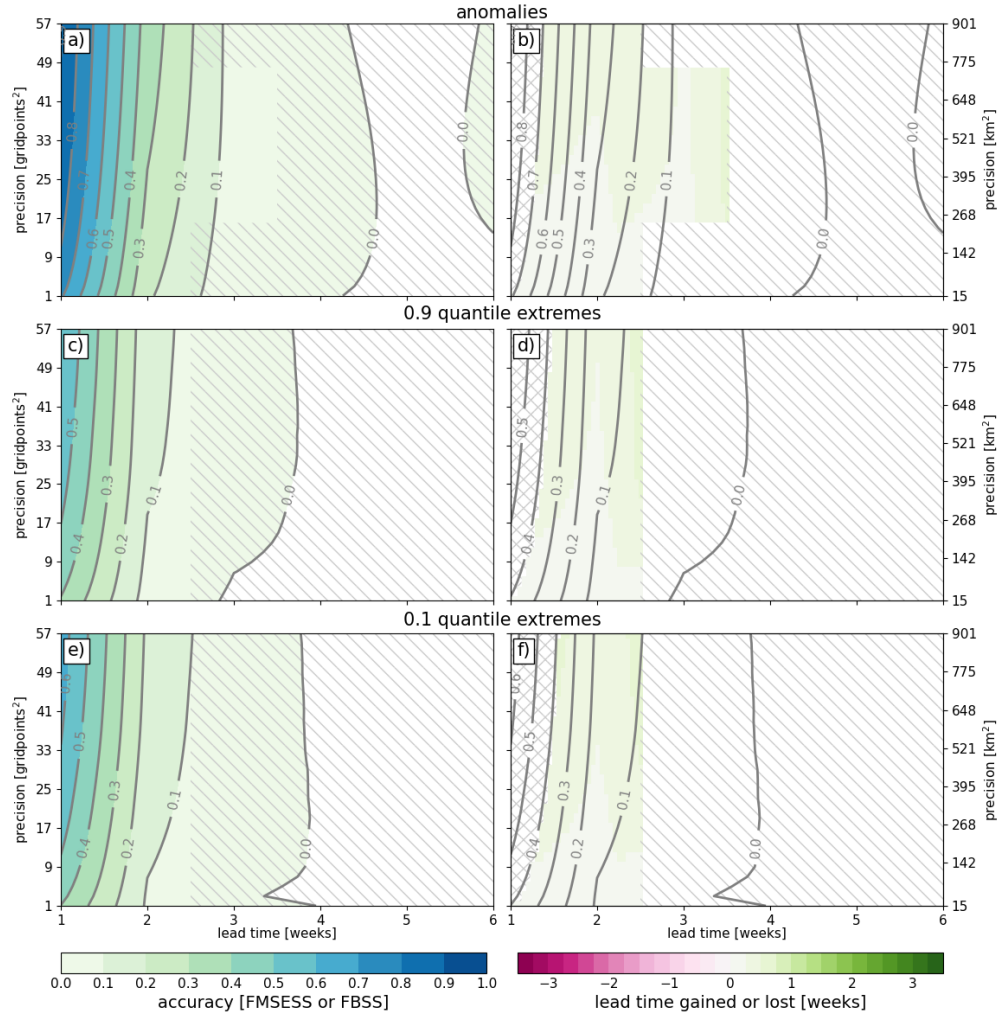


Figure 3: As in Fig. 1 except for weekly-accumulated precipitation forecasts.

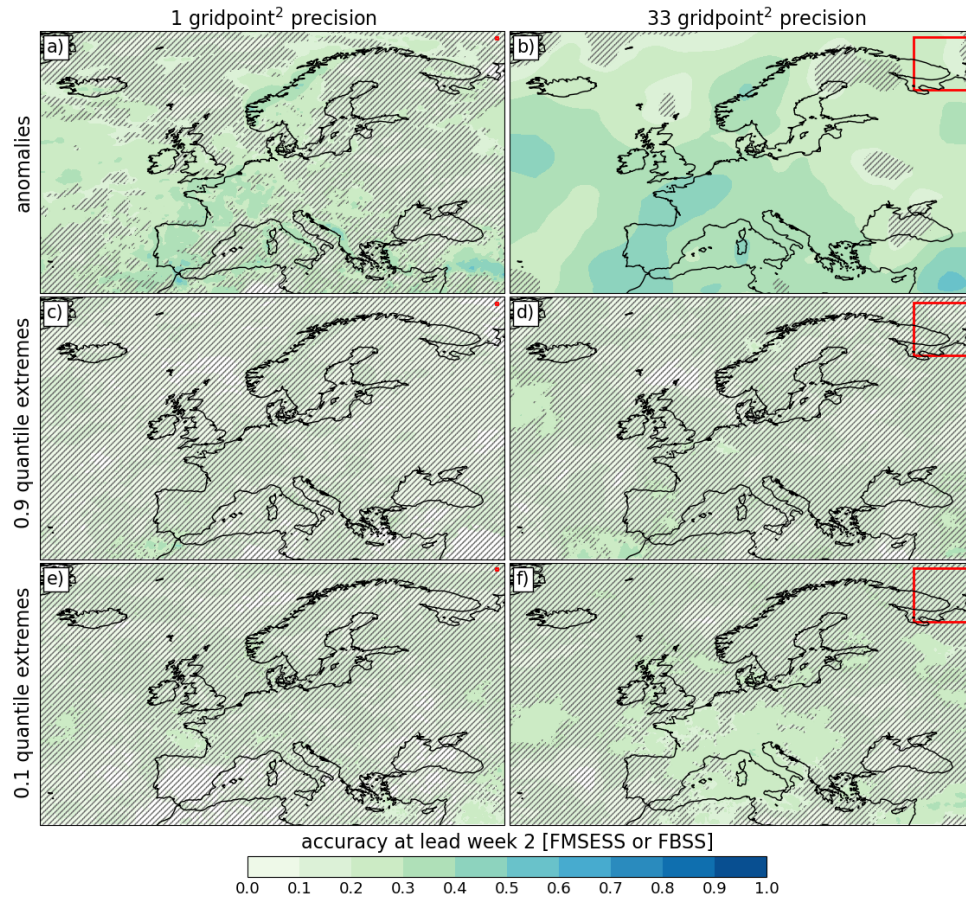


Figure 4: As in Fig. 2 except for weekly-accumulated precipitation forecasts at lead-week 2.



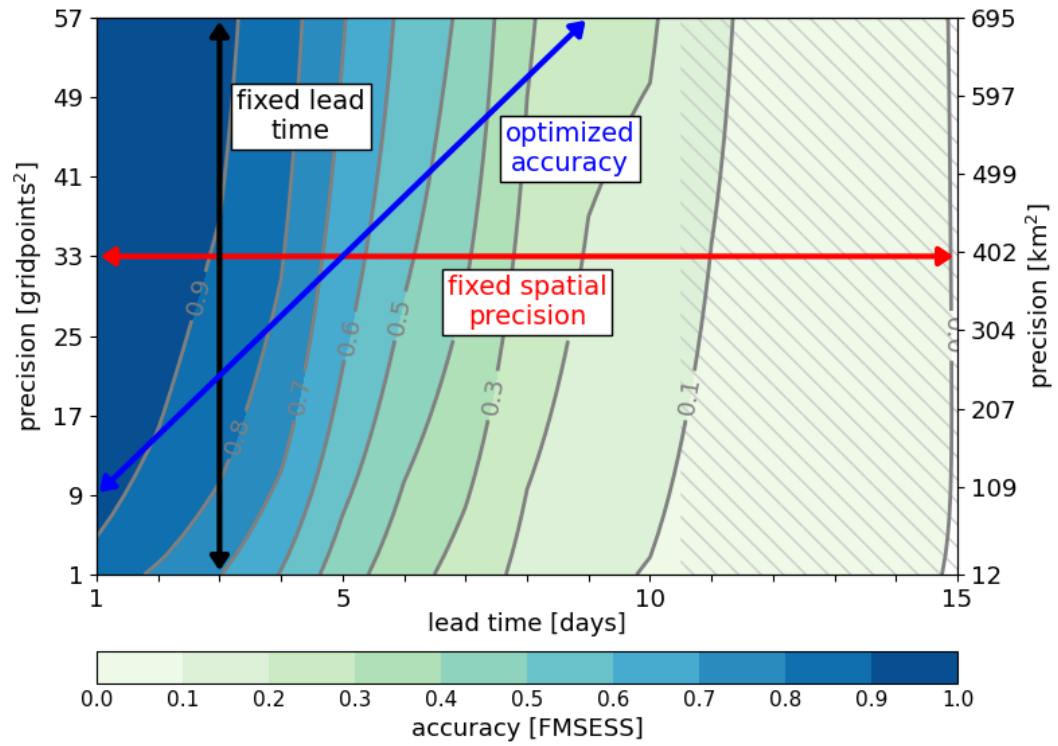


Figure 5: As in Fig. 1a except for the Fractions Mean-Square Error Skill Score (FMSESS) over the Scandinavian domain shown in Fig. 6. Black, blue and red arrows denote the use-cases described in section 5.

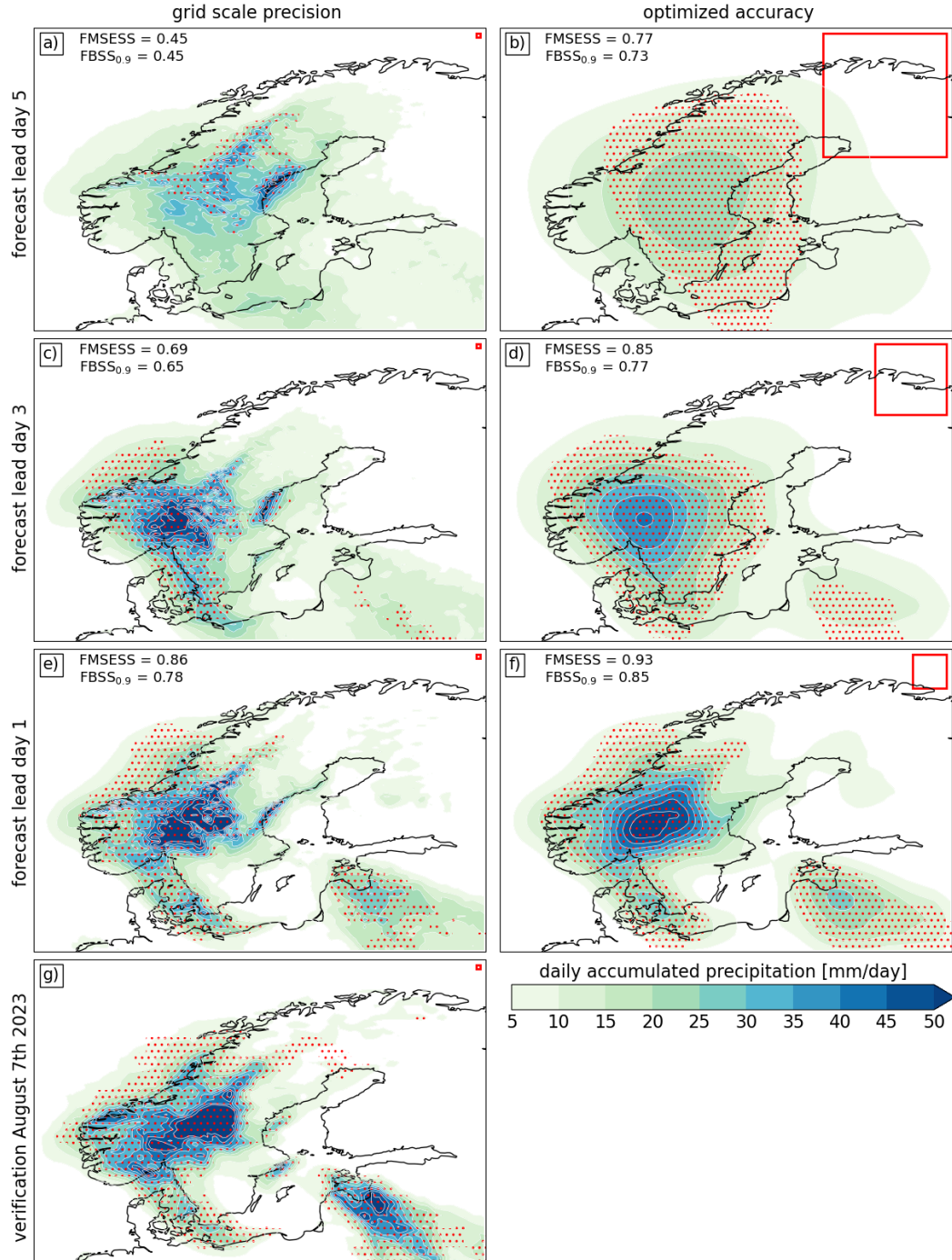


Figure 6: (a-f) Forecasts of storm Hans on August 7<sup>th</sup> 2023. Shading denotes the ensemble-mean daily-accumulated precipitation for lead days 5 (a,b), 3 (c,d) and 1 (e,f), shown at the grid-scale (left) and with progressively increasing levels of spatial precision (right). Red squares in the upper-right corner of the panels denote the spatial scale of forecast precision. The FMSESS and FBSS for 0.9 quantile extremes are displayed in the top left hand corner of each panel. (g) Daily-accumulated precipitation during Storm Hans on August 7<sup>th</sup> 2023 according to data from the ERA5 reanalysis. Red stippling denotes Fractions Extreme Forecast Index (FEFI, panels a-f) and Extreme Verification Index (EVI, panel g) values greater than 0.8, signaling an extreme event.