# Accepted Manuscript

Evaluation of ensemble streamflow predictions in Europe

Lorenzo Alfieri, Florian Pappenberger, Fredrik Wetterhall, Thomas Haiden, David Richardson, Peter Salamon

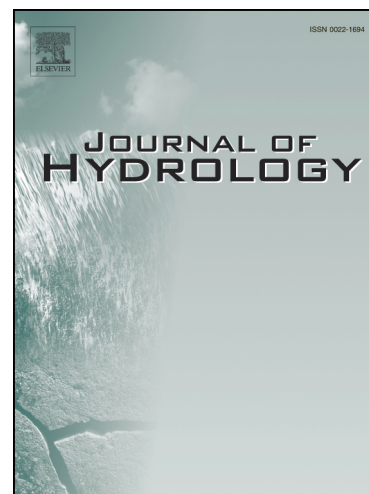Please cite this article as: Alfieri, L., Pappenberger, F., Wetterhall, F., Haiden, T., Richardson, D., Salamon, P., Evaluation of ensemble streamflow predictions in Europe, *Journal of Hydrology* (2014), doi: http://dx.doi.org/10.1016/j.jhydrol.2014.06.035

**Evaluation of ensemble streamflow predictions in Europe**

Lorenzo Alfieri[1,2], Florian Pappenberger[1], Fredrik Wetterhall[1], Thomas Haiden[1], David

Richardson[1], and Peter Salamon[2].

[1] European Centre for Medium-Range Weather Forecasts, Reading, UK

[2] European Commission - Joint Research Centre, Ispra, Italy

Correspondence to: Lorenzo Alfieri

European Commission - Joint Research Centre

TP 122, Via E. Fermi 2749

21027 Ispra (VA)

Italy

Tel: +39 0332 78 3835

Email: lorenzo.alfieri@jrc.ec.europa.eu

**Abstract**

In operational hydrological forecasting systems, improvements are directly related to the

continuous monitoring of the forecast performance. An efficient evaluation framework must

be able to spot issues and limitations and provide feedback to the system developers. In

regional systems, the expertise of analysts on duty is a major component of the daily

evaluation. On the other hand, large scale systems need to be complemented with semi-

automated tools to evaluate the quality of forecasts equitably in every part of their domain.

This article presents the current status of the monitoring and evaluation framework of the

European Flood Awareness System (EFAS). For each grid point of the European river

network, 10-day ensemble streamflow predictions are evaluated against a reference simulation which uses observed meteorological fields as input to a calibrated hydrological model. Performance scores are displayed over different regions, forecast lead times, basin sizes, as well as in time, considering average scores for moving 12-month windows of forecasts. Skilful predictions are found in medium to large rivers over the whole 10-day range. On average, performance drops significantly in river basins with upstream area smaller than 300 $km^2$, partly due to underestimation of the runoff in mountain areas. Model limitations and recommendations to improve the evaluation framework are discussed in the final section.

Keywords: flood early warning; ensemble streamflow predictions; CRPS; skill scores; distributed hydrological modelling.

# 1 Introduction

Operational hydrological forecasting systems play a key role in the water resources management and in the preparedness against extreme events. Assessing their performance is crucial for the error diagnostic and in the planning of development work to improve the system accuracy and extend the forecast lead time. A vast number of regional and national hydro-meteorological centres have flood forecasting and early warning systems in place based on weather predictions (see Alfieri et al., 2012 for a recent review of European systems). At the same time, the number of ensemble-based systems is increasing (Cloke and Pappenberger, 2009; Wetterhall et al., 2013), with the aim of describing part of the uncertainty embedded in the forecasts. The evaluation of the forecast accuracy is regularly performed in many operational systems, where verification scores need to be complemented by the local knowledge and experience of analysts on duty. Further, skill scores are rarely displayed publicly, to prevent misinterpretation of results and avoid the need for simplifying their information content for a wider recipient of users. Yet, reporting on past performance by means of verification scores is listed as one of the main priorities of users, to increase the trust in forecasting systems (Wetterhall et al., 2013).

Assessing the forecast performance over large domains raises the challenge of comparing river points with different upstream area and hydrological regimes. In these cases, a widespread approach to tackle the forecast verification is to compute scores based on the probability of thresholds exceedance (e.g., warning levels), that can be defined in a consistent way for every point. While this is a standard practice for early warning systems (e.g., Bartholmes et al., 2009; Gourley et al., 2012), it is also applied to the verification of categorical events for any set of thresholds (Thirel et al., 2008). If quantitative values are considered, the choice of performance scores becomes wider (Legates and McCabe, 1999;

25 Wilks, 2006), though only a relatively small subset is specifically dedicated to evaluate the

26 quality of ensemble forecasts (Brown et al., 2010). The comparison of forecast skill in several

27 river sections is often performed through benchmarking against simplified simulations

28 (Pappenberger et al., submitted), previous model versions (Arheimer et al., 2011), different

29 input data (e.g., Renner et al., 2009), or climatological values (Demargne et al., 2010;

30 Verkade et al., 2013; Wood et al., 2005). An alternative method consists in normalizing

31 forecasts and reference values before the evaluation (Pappenberger et al., 2010). Trinh et al.

32 (2013) used a similar concept to propose a modified Continuous Ranked Probability Score

33 (CRPS) which is suitable to compare forecast performance at different river sections. In

34 operational systems, the forecast performance must be monitored and updated continuously in

35 time. Hence, a skill assessment based on different scores and benchmarks (e.g., Alfieri et al.,

36 2013a; Randrianasolo et al., 2010) is often preferred in order to analyze different aspects of

37 the forecast performance at several locations and quickly detect trends over time or

38 weaknesses.

39 In 2012, after the transfer of the EFAS operational suite to the European Centre for Medium-

40 Range Weather Forecasts (ECMWF), a commitment was made to set up an evaluation

41 framework of the hydrological forecasts, in order to monitor their performance over time and

42 after major system updates. The idea was to implement an automated procedure to regularly

43 produce and update summary skill scores for the whole computation domain, able to spot a

44 variety of possible problems and address subsequent in-depth analysis. Among the main

45 challenges to face was the choice of appropriate skill scores, the handling of large data sets,

46 and the visualization of results through concise and intuitive graphs.

47 This article presents the current status of implementation of such an evaluation framework,

48 after one year of operational runs at ECMWF. Streamflow forecasts at every grid point of the

4

49    river network are verified against a reference simulation which uses observed meteorological

50    fields as input to a calibrated hydrological model.

51    **2    Data and Methods**

52    **2.1    Model framework**

53    The main components of the EFAS hydro-meteorological forecasting chain are: a) a

54    hydrological model, b) weather forecasts, and c) meteorological observations, to update the

55    initial model states and for verification purpose (see Figure 1). Each of these three

56    components has inherent uncertainty, which can be described in the modelling framework and

57    propagated to the output discharge. The current EFAS system is a multi-model ensemble

58    approach, in that it accounts for the uncertainty of input weather forecasts using model runs

59    from different meteorological centres in Europe. These include two deterministic forecasts,

60    from the ECMWF (ECMWF-HiRes, Miller et al., 2010) and from the German Weather

61    Service (DWD, see Majewski et al., 2002; Steppeler et al., 2003), and two ensemble forecasts,

62    from the COSMO Consortium (COSMO-LEPS, Marsigli et al., 2005) and from ECMWF

63    (ECMWF-ENS, Miller et al., 2010). The version of the evaluation framework presented here

64    is based on the performance of the ECMWF-ENS forecasts only, though it is foreseen to

65    extend it to include the other model simulations. The system setup and additional details on

66    how weather forecasts are handled in EFAS are documented in the published literature

67    (Bartholmes et al., 2009; Pappenberger et al., 2010; Thielen et al., 2009), therefore we refer

68    the reader to these articles for additional information not included in the present work, and

69    focus on the analysis of the evaluation framework.

70

71    *Figure 1: Schematic view of the EFAS hydro-meteorological forecasting system.*

## 2.2  Meteorological data

ECMWF-ENS is a 51-member ensemble forecast run twice per day, at 00 UTC and 12 UTC as part of the operational production suite of ECMWF Integrated Forecast System (IFS, see Bechtold et al., 2014; Miller et al., 2010). ENS forecasts are run globally at T639 spectral resolution, corresponding to about 32 km horizontal resolution, with forecast lead time (LT) up to 10 days. After day 10, the model run is extended up to day 15 (day 32 twice per week) at a coarser horizontal resolution of about 65 km. Currently, EFAS uses only the first 10 days of forecast as input to the hydrological model. For this work, ENS forecasts from January 2009 to the present were extracted and used in the hydrological simulations, considering those available at the time of the forecasts (i.e., no reforecast with more recent IFS versions was used). Meteorological forecast fields used are total precipitation, evaporation, and 2-metre temperature, which are regridded to the same spatial resolution of the hydrological model (see next section).

A database of observed meteorological fields for Europe was provided by the Joint Research Centre of the European Commission. It consists of maps of spatially interpolated point measurements of precipitation and temperature at the surface level. The database includes daily data from the 1990 to the present, and it is populated by an increasing number of reporting gauges over time, with the latest figures showing on average more than 6000 stations for precipitation and more than 4000 for temperature (see Figure 2 for a recent example of daily data). A subset of the same meteorological station network is used to generate interpolated potential evapotranspiration maps using the Penman-Monteith method.

*Figure 2: stations reporting observed precipitation (left) and average temperature (right) on the 1st October 2013.*

96

## 2.3   Hydrological modelling

98   In EFAS, hydrological simulations are performed with Lisflood, a hybrid between a

99   conceptual and a physical rainfall–runoff distributed model, designed to reproduce the main

100   hydrological processes of medium to large river basins (see van der Knijff et al., 2010). The

101   considered model setup for Europe was calibrated at 481 river gauges, using the observed

102   meteorological fields as input and up to 7 years of gauged discharge. A reference hydrological

103   simulation starting in 1990 was run for the European window with the calibrated Lisflood

104   model at 5x5 km resolution, using the observed meteorological fields as input. The

105   operational model is updated daily using the initial states of the previous day and the most

106   recent meteorological observations acquired with about 1 day lag. This simulation, hereafter

107   referred to as EFAS Water Balance (EFAS-WB), represents our best estimate of the

108   hydrological states in the European rivers. The EFAS-WB is used in EFAS with regard to

109   three main aspects (see Figure 1): *I*) deriving climatological features of the runoff in each

110   point of the river network (e.g., average conditions, extremes, alert thresholds, seasonality);

111   *II*) creating initial conditions for daily hydrological runs driven by the latest weather

112   predictions; *III*) providing a reference simulation which is as realistic as possible, to be used

113   as a proxy to evaluate streamflow forecasts in every grid point of the simulation domain.

114   Further details on the EFAS-WB are described by Alfieri et al. (2013b). The same calibrated

115   Lisflood setup is used to perform 10-day EFAS streamflow forecasts updated twice per day,

116   by forcing the hydrological model with initial conditions from the EFAS-WB and with

117   forecast weather fields (described in the previous section) with 1-day temporal resolution.

## 3    Evaluation strategy

EFAS forecasts are run at the ECMWF twice per day since October 2012, using weather

predictions initialized at 00 and 12 UTC. This operational dataset of hydrological forecasts

was complemented by running 4 years of daily hindcasts with the same model configuration,

starting on January 2009. To reduce the computing load, the hindcasts were run only once per

day, using forecast runs from 12 UTC. Ensemble streamflow predictions (ESP) are validated

against the EFAS-WB for each point of the modelled European river network, comprising

38452 grid points. Such an approach enables a quick spatial overview of skill scores on every

region of the computation domain, rather than just at stations where observed discharge is

provided. On the other hand it does not account for the potential mismatch between actual

river discharge and the simulated EFAS-WB used as reference.

Average scores are calculated over 1-year time windows. This choice proved to be effective

as it includes one full hydrological year and dampens the seasonal variability of skill scores.

In practice, the verification of dry months leads to higher scores than those of rainy months,

as the quantitative forecast of high precipitation amounts is more challenging than forecasting

days with zero precipitation. As a result, the evaluation framework was set up to select the

first day of each month and calculate the average skill scores of the previous 365 days,

starting on the 1st January 2010. The procedure was then semi-automated and skill scores are

now updated every month to include results of the latest forecasts.

Skill scores to evaluate the ESP were chosen so that grid points with different upstream area

and climatic regime could be compared together in the same graphs and in the same maps. To

this end, four different dimensionless skill scores were selected, able to stress different

aspects of the forecast performance. These are described in the following sub-sections and

summarized in Table 1.

### 3.1 Nash-Sutcliffe efficiency

The Nash-Sutcliffe efficiency (NS, Nash and Sutcliffe, 1970) applied to discharge forecasting can be defined as:

$$NS = 1 - \frac{\sum_{t=1}^{N}[q_{sim}(t) - q_{fc}(t)]^2}{\sum_{t=1}^{N}[q_{sim}(t) - \overline{q}_{sim}]^2} \, , \tag{1}$$

where $q_{sim}$ is the proxy discharge given by the EFAS-WB and $q_{fc}$ is the forecast discharge at the same time step. $t$ is a time index spanning all $N$ forecasts included in the evaluation window, that is $N$=730 in operational forecasts (when two forecasts per day are evaluated) and $N$=365 for hindcasts between 2009 and 2012. In the case of the considered ESP, $q_{fc}$ represents the mean of the 51-member ensemble. The NS values range from -∞ to 1, the latter corresponding to perfect forecasts. NS above 0 means that forecasts perform better than climatological values, in the form of their average discharge $\overline{q}_{sim}$. In the presented work, NS values are calculated for fixed forecast lead times between 1 and 10 days, and the average values over 1 year windows are shown, as described in the previous section.

### 3.2 Forecast bias

Monitoring the bias of ensemble streamflow predictions is of vital importance for a flood awareness system based on a threshold exceedance approach as in EFAS. Flood alerts are detected by comparing EFAS simulations driven by weather forecasts as input, against reference warning thresholds, derived from the EFAS-WB. If weather forecasts were persistently different from observed meteorological values, discharge forecasts would be consequently biased, which may result in statistically significant over- or under-prediction of flood alerts. The main potential source of bias in ESP is the quantitative forecast of precipitation, particularly for high flow events. However, biased forecast values of temperature may induce cyclical drifts of discharge predictions, particularly in hydrological

165 regimes where the snow accumulation and melting processes play a prominent role. In

166 addition, precipitation, temperature and evapo-transpiration are key drivers for the soil

167 moisture state, therefore consistent bias in their forecast values can affect the streamflow

168 potentially over long ranges (i.e., monthly to inter-annual time scales). In the presented

169 evaluation framework, the bias at each grid point is rescaled by the corresponding average

170 discharge for the same period, calculated from the EFAS-WB:

171 $$\mathbf{Pbias} = \frac{\frac{1}{N}\sum_{t=1}^{N}[q_{sim}(t) - q_{fc}(t)]}{\bar{q}_{sim}} \tag{2}$$

172 Being a linear operator, the sum of the percentage bias (Pbias) of all ensemble members is

173 equal to the percentage bias of the ensemble mean.

### 3.3 Coefficient of variation of the RMSE

175 The Root Mean Squared Error (RMSE) has long been used to assess the magnitude of the

176 error of deterministic forecasts. It has the advantage that it retains the units of the forecast

177 variable and it includes the effect of both bias and variance of estimation. In addition, the

178 RMSE depends on a quadratic function of the estimation residuals. This lead to some

179 peculiarities, among which: 1) it is highly affected by few large errors and 2) it is often used

180 as an error function to be minimized in a wide range of calibration and optimization

181 processes. On the other hand it is difficult to compare RMSE values among different river

182 stations, as their climatological discharge values may be substantially different. One option to

183 compare the RMSE at different locations is to rescale it by the corresponding average

184 discharge, as shown in Reed et al. (2007), so that resulting values become dimensionless:

185 $$CV = \frac{\sqrt{\frac{\sum_{t=1}^{N}[q_{sim}(t) - q_{fc}(t)]^2}{N}}}{\bar{q}_{sim}}, \tag{3}$$

10

186 The resulting score is commonly referred to as coefficient of variation (CV) of the RMSE

187 and, as for the RMSE, values close to zero are preferable. Also, when CV values are close to

188 1 it means that the RMSE of estimation is of the same order as the average discharge. Indeed,

189 it can be associated to an inverse of the signal-to-noise ratio. By definition the CV penalizes

190 river reaches with low average discharge compared to its variability, therefore higher CV

191 values are expected in small or flash-flood prone river basins, such as those along the

192 Mediterranean coast, where the predictability is indeed shorter than in large river basins.

193 **3.4 Continuous Ranked Probability Skill Score**

194 To fully exploit and assess the added value of probabilistic predictions, the Continuous

195 Ranked Probability Skill Score (CRPSS) is used to evaluate the quantitative skills of the ESP.

196 The CRPSS (e.g., Hersbach, 2000) is defined as:

197
$$CRPSS = \frac{\overline{CRPS_{ref}} - \overline{CRPS_{forecast}}}{\overline{CRPS_{ref}}},$$
(4)

198 where

199
$$CRPS = \int_{-\infty}^{\infty} [F(y) - F_0(y)]^2 dy$$
(5)

200 and

201
$$F_0(y) = \begin{cases} 0, & y < observed\ value \\ 1, & y \geq observed\ value \end{cases}$$
(6)

202 while $F(y)$ is the stepwise cumulative distribution function (cdf) of the ESP of each

203 considered forecast. The CRPSS is a dimensionless indicator of the skill of ensemble

204 predictions, measured by ($CRPS_{forecast}$), compared to that of a reference forecast ($CRPS_{ref}$).

205 The CRPSS ranges between 1 (for perfect predictions) to -∞, though ESP are valuable only

206 when CRPSS>0, i.e., when the forecasts perform better than the reference. In this work, we

207 compare and discuss the use of two different $CRPS_{ref}$ to evaluate the CRPSS, the first based

11

208    on the average climatological discharge $\bar{q}_{sim}$ ($CRPS_{ref,ad}$), and the second based on a

209    persistence forecast ($CRPS_{ref,pf}$), meaning a forecast given by assuming the same value used to

210    initialize the ESP. It is worth noting that both reference CRPS are based on deterministic

211    predictions, hence the $CRPS_{ref}$ reduces to the mean absolute error (Hersbach, 2000):

212    $$CRPS_{ref,ad} = \frac{1}{M}\sum_{t=1}^{M}|q_{sim}(t) - \bar{q}_{sim}| \qquad (7)$$

213    where $t$ is a daily time index going from 1/1/1990 to the present. On the other hand,

214    $$CRPS_{ref,pf}(LT) = \frac{1}{N}\sum_{t=1}^{N}|q_{sim}(t) - q_{sim}(t - LT)| \qquad (8)$$

215    where $N$ has the same meaning as in Eq.1.

216    Two significant differences between Eq.7 and 8 can be seen. The $CRPS_{ref,ad}$ is a constant

217    value and only depends on the location, though it needs climatological information to be

218    evaluated, in the form of a reference time series of observations or proxy simulations (i.e., the

219    EFAS-WB in this case). On the other hand, the $CRPS_{ref,pf}$ depends on the lead time of the

220    forecast (LT). It does not need any prior climatological information on the discharge regime

221    at the point but the discharge value used to initialize the forecast.

## 4    Results

223    Skill scores of the last available year are now routinely calculated on the 13th day of each

224    month, after all meteorological observations to update the EFAS-WB are received and the

225    hydrological model is run. Simulated proxy discharges need to be computed until the 11th of

226    the same month, so that 10-day ESP starting on the 1st can be evaluated. Scores described in

227    Sect. 3 are shown in **Figure 3**. NS, CV and Pbias are deterministic scores; hence they are

228    calculated on the ensemble mean, while the CRPSS take into account the whole ensemble. A

229    forecast lead time of 5 days is chosen for most figures in the article, being representative of

230    the general behaviour of the ESP and a frequent lead time of EFAS flood alerts. One can see

12

231    that, for LT=5 days, in the vast majority of grid points the ESP is more skilful than a

232    persistence forecast (i.e., $CRPSS_{pf}>0$). The NS and the CV suggest that higher performance is

233    achieved in large rivers of Central and Northern Europe. Excluding Iceland, lower skills are

234    mostly seen in Southern Europe and can be explained by a) resolution issues in small basins,

235    b) less skilful precipitation forecast in mountainous areas, c) a comparatively lower station

236    density to run the EFAS-WB, and d) the higher proportion of convective precipitation,

237    leading to higher space-time variability of rainfall rates and larger extremes over short (i.e., 1-

238    day or sub-daily) durations. Similarly, the Pbias (on gray background in **Figure 3**) shows a

239    widespread underestimation of discharge over the main mountain ranges (i.e., Pyrenees, Alps

240    and Balkans, among others), mostly in the range 10% to 50% of the corresponding average

241    flow. These findings are in line with previous works by Wittmann et al. (2010) and

242    Pappenberger et al. (2013), who showed increasing underestimation of precipitation and

243    streamflow forecast in the Alpine region during intense precipitation events. The apparently

244    poor performance over Iceland in **Figure 3** is actually imputable to an incorrect reference

245    streamflow. Indeed, the number of reporting stations for this region is very low (see an

246    example in **Figure 2**), particularly for precipitation, thus leading to a considerable under-

247    prediction of the streamflow. In other words, although EFAS streamflow forecasts over

248    Iceland may be skilful, the current availability of meteorological observations prevents from

249    simulating reliable reference discharge to perform forecast evaluation in this area. In the

250    following analyses, summary scores of grid points in Iceland are excluded from all figures,

251    which brings the dataset to a subset of 37588 points.

252

253    *Figure 3: $CRPSS_{pf}$, CV, NS and Pbias over Europe for 1 year of daily forecasts ending on the*

254    *1st October 2013 (5-day lead time).*

## 4.1 Performance versus forecast range

Skill scores as in **Figure 3** are shown in **Figure 4** for each forecast lead time between 1 and 10 days. A solid line indicates the mean value among all grid points, while grey shades denote the 5%-95% (light grey) and the 25%-75% (dark grey) of their distribution. In the top-left panel, the CRPSS calculated using the average discharge as reference (i.e., $CRPSS_{ad}$) is shown with a thick dashed line (mean value) together with the corresponding 25%-75% values (dotted lines). Differences between the two methods are the largest for the first lead time, where in many cases the ESP does not bring substantial differences in comparison to a persistence forecast, due to the large weight of the initial model states. On the other hand, the $CRPSS_{ad}$ decreases roughly linearly and suggests the presence of a crossing point for a LT>10 days, when the climatological average discharge seems to become a more skillful benchmark than a persistence forecast. As expected, the CV tends to deteriorate with the lead time, though without a significant increase of the spread of its distribution. Similarly, the mean NS ranges between 0.9 for LT=1 and 0.7 at the end of the forecast range, while in 99% of forecasts NS>0 for LT=10 days. The Pbias shows a rather constant mean under-prediction of 2% to 4%. Its distribution has an increasing spread with the lead time, with 65% to 70% of grid points lying constantly below the zero line.

*Figure 4: CRPSS, CV, NS and Pbias of ESP versus the forecast lead time.*

## 4.2 Performance versus catchment size

**Figure 5** displays the four scores against the upstream area of each grid point, calculated over 1 year ending on 1/10/2013 and for a 5-day lead time. In addition, solid lines indicate the empirical median value (i.e., 50th percentile), in light grey, and the central 90% of the

14

279    distribution (i.e., 5th to 95th percentiles), in dark grey. Largest values on the x-axis

280    correspond to the lower Danube River, with upstream area up to about 800,000 km$^2$. On the

281    left side of each panel, one can note the model grid resolution as limit, with catchments area

282    being always a multiple of 25 km$^2$. Results in **Figure 5** denote a general positive trend of skill

283    scores with increasing upstream area. Indeed, in large rivers, a) the discharge varies more

284    gradually due to the smoothing and averaging effect of the complex river network and b) the

285    influence of the initial discharge, compared to the forecast precipitation input, is larger than in

286    smaller catchments. In detail, as the basin time of concentration increases and approaches the

287    magnitude of the forecast range, a larger proportion of the forecast discharge at the river

288    outlet is made up by a water volume which is already in the model, (i.e., gauged) at the

289    starting time of the forecast run. Therefore the skill of weather forecasts affects that of

290    streamflow forecasts with an average delay increasing with the upstream area, which can be in

291    the order of some days for large European rivers. On the other hand, **Figure 5** shows a clear

292    deterioration of scores for catchments smaller than 300 km$^2$, that is, for a ratio between

293    upstream area and grid size of the weather forecasts of about 0.3. Results are in agreement

294    with those of Pappenberger et al. (2010), though Bartholmes et al. (2006) suggested a

295    minimum threshold of 4000 km$^2$ if extreme values are considered. Indeed, the latter value is

296    used in EFAS as minimum upstream area for flood alerts to be issued to partner institutes.

297    The median value of the Pbias in **Figure 5** indicates that the deterioration of scores can be

298    partly attributed to the underestimation of the discharge for small catchments, which

299    decreases below 2%, in absolute value, for upstream areas larger than 400 km$^2$. As

300    commented in Sect. 4, such trend is to be attributed to the under-prediction of quantitative

301    precipitation in mountain areas and of extreme values in general, not fully captured by the

302    atmospheric circulation model due to its grid size on average coarser than the observation

303    network.

304

305  *Figure 5: $CRPSS_{pf}$, CV, NS and Pbias of ESP versus the upstream area of each river grid*

306  *point.*

307

### 4.3  Evolution of 12-month average performance

309  The evolution of summary scores over the past 5 years is shown in **Figure 6**. Scores are

310  calculated on the 365 days preceding the first day of each month indicated in the x axis. In the

311  top-left panel both $CRPSS_{pf}$ and $CRPSS_{ad}$ are shown, using the same line types as in **Figure 4**.

312  In addition, the average discharge over all grid points of the river network, for each evaluation

313  period, is drawn at the bottom. One can note how the $CRPSS_{ad}$ is largely affected by the

314  magnitude of the observed runoff, so that, in drier years, it gives the impression of increasing

315  forecast performance, and vice-versa. In the $CRPSS_{pf}$, no dependence on the average runoff is

316  visible. The latter shows an improvement of the forecast skills during the year 2013,

317  particularly for the mean of the distribution and for the 75th and 95th quantiles. Such

318  improvement is also pointed out by a reduced mean CV and increased mean NS, where in

319  both cases the central 90% of the distribution becomes narrower since the beginning of 2013,

320  though with a subsequent widening towards the end of the year.

321  Interestingly, the bottom-right panel denotes a slow but constant increase of a negative bias in

322  forecast streamflow over the last years. This appears consistently on all lead times (not

323  shown), though it is more significant towards the end of the forecast range. On the other hand,

324  no corresponding trend was reported in the forecast input precipitation produced by the

325  ECMWF-ENS (personal communication, see some additional details in

326  http://www.ecmwf.int/products/forecasts/d/charts/medium/verification), nor in temperature

327  (possibly inducing a larger snow fraction). Instead, the main reason for such discrepancy is

16

328　most likely due to the progressive increase in the number of stations reporting meteorological

329　observations in recent years. Higher station density leads to a more realistic representation of

330　the input maps to run the EFAS-WB, so that small-scale features such as convective cells are

331　more likely to be better observed quantitatively. In this regard, Kann and Haiden (2005)

332　showed that when high density stations networks are used as reference, the mean absolute

333　error of forecast precipitation tend to increase with the reduction of the aggregation area.

334　Further, some of the stations added recently are located in elevated areas, such as in the Alps

335　and the Pyrenees, where the orography enhances annual rainfall totals and consequently the

336　runoff. Indeed, these areas are where the under-prediction of discharges has become clearer in

337　the recent years, as shown in **Figure 3**.

338

339　*Figure 6: Trend of 12-month average CRPSS, CV, NS and Pbias of ESP from 2009 onwards.*

340

341　**5　Discussion and conclusions**

342　This article presents the current status of the evaluation framework used to monitor and

343　update regularly the forecast performance of the European Flood Awareness System. Results

344　suggest that streamflow forecasts driven by weather predictions provide significant added

345　value to the monitoring of the main European rivers. As expected, performance decreases

346　with lead time, though it remains skilful for the whole 10-day range, in comparison to the use

347　of climatological or persistence forecasts. In large river basins of Europe, the average time lag

348　between weather forcing and runoff is on the order of some days. Hence the real-time

349　hydrological simulation run with meteorological observations gives a significant proportion of

350　the overall predictability, increasing with the basin time of concentration. In smaller river

351　basins, the effect of initial conditions is less important, therefore the predictability is shorter

352 as it mostly depends on that of the weather forecasts. In river basins of size below 300-400

353 $km^2$ forecast skill becomes poorer. Their forecasts show large variability, often even for 1-day

354 lead time, and significant underestimation of the runoff in mountain regions.

355 Being designed on dimensionless scores, the main strength of the proposed verification

356 system is in highlighting relative changes of performance, which can be detected over

357 different regions, forecast lead time, basin size and, most importantly, in time. An evaluation

358 of 12-month average scores over the past 5 years suggests a moderate improvement for all 12-

359 month forecasts ending from the beginning of 2013 onwards. Such improvement occurred

360 notwithstanding an increasing negative forecast bias, especially in mountain regions. This can

361 be attributed to a progressive increase of the meteorological stations used to run the EFAS-

362 WB, which in turn has improved the representation of the runoff dynamics in the presence of

363 pronounced orography. Although the parameterization of the hydrological model was subject

364 to changes and improvements every 1 to 1.5 years on average, the 5 year simulation shown in

365 this study was carried out with a fixed model version, corresponding to the current operational

366 one at the time of writing. Therefore, the positive trend of performance shown in **Figure 6** is

367 likely to underestimate the real improvements which have occurred and rather reflect that of

368 weather forecasts used as input.

369 **5.1 The benchmark of skill scores**

370 The four performance scores presented in the article can be classified into two categories,

371 depending on whether the comparison is carried out against a benchmark or not. On the one

372 hand, the CV and the Pbias give a measure of the RMSE and of the bias of forecasts,

373 respectively. RMSE and bias are commonly used in verification because of their physical

374 meaning, as they quantify the error with the same units of the forecast variable. They are

375 rescaled by the average flow to make them comparable over different regions and along the

18

376 river network. On the other hand, the NS and the CRPSS give a relative performance in

377 comparison to an alternative benchmark forecast. Literature works show a surprising variety

378 of different benchmarks used for comparison (see Pappenberger et al., submitted, for a recent

379 review), sometimes without motivating the choice. Here we argue that, in assessing the

380 predictability of a forecasting system, the benchmark should represent a realistic forecast

381 achievable in case the system was not in place. The use of persistence forecasts is hereby

382 suggested as a suitable benchmark, in that it does not require climatological information of

383 the runoff at the river point, nor additional model runs. In comparison to a benchmark based

384 on the average discharge, persistence acknowledges the role of initial conditions, indicating

385 that the highest value of forecasts corresponds to a balance between the ability to provide

386 accurate forecasts and the ability to detect deviations from an initial state (see $CRPSS_{ad}$ versus

387 $CRPSS_{pf}$ in **Figure 4**). Further, persistence is independent of seasonal variations or trends in

388 the mean value of the forecast variable, as discussed in Sect. 4.3.

389 It is worth noting that the same principle can be applied to the Nash-Sutcliffe efficiency, as

390 suggested by Plate and Lindenmaier (2008), leading to a modified formulation which uses a

391 persistence forecast as reference value:

392 $$NS(LT) = 1 - \frac{\sum_{t=1}^{N}[q_{sim}(t) - q_{fc}(t)]^2}{\sum_{t=1}^{N}[q_{sim}(t) - q_{sim}(t-LT)]^2} \,. \tag{9}$$

393 This formulation was not tested in the present framework, though it may be a valid alternative

394 to the NS for large river basins (see e.g., Pagano, 2013). Its application will be considered for

395 future system developments.

396 **5.2 The EFAS-WB as reference simulations**

397 The main assumption of the presented approach is that the EFAS-WB can be used as a

398 realistic representation of the actual runoff. On the other hand the use of the output of a

399  distributed model as the EFAS-WB allows a performance evaluation over the full

400  computation domain. Moreover, the continuous increase in the number of reporting stations,

401  both for meteorological and hydrological data, is progressively pushing the EFAS-WB closer

402  to the real streamflow conditions in the European rivers. This occurs thanks to a better

403  reproduction of the meteorological input data and to the increase of the number of river

404  stations where the parameters of the hydrological model can be calibrated. Recent advances in

405  the meteorological dataset include the addition of more than 10 high density national

406  networks and an improved approach to interpolating point values into spatial maps (see

407  Ntegeka et al., 2013). This is currently being tested and will be used in the next version of

408  EFAS, together with additional historical observed streamflow at a number of river gauges to

409  improve the model calibration. Similarly, resulting simulated discharges of the EFAS-WB can

410  potentially become a dataset to validate and benchmark a wide range of hydrological models,

411  particularly on large scales. Current main limitations of simulated discharges are at the lower

412  end of the range of the space-time scale of simulated catchments. In fact, the current daily

413  time aggregation of input data induces a smoothing of output discharges, so that simulated

414  extreme values have reported under-estimation issues, relatively to observed values. In

415  addition, the presented scores are not able to capture potential errors in the hydrological

416  model, because both ESP and the EFAS-WB used for validation are generated by the same

417  model. However, this is evaluated separately at those stations where the model parameters are

418  calibrated (see Feyen et al., 2007). Also, an assessment of the total predictive uncertainty is

419  performed at river gauges (currently about 40) where discharge values are received in real-

420  time. The methodology and results are described by Bogner and Pappenberger (2011).

421  **5.3  Concluding remarks**

422  In its current state, the evaluation framework has proved its usefulness in spotting strengths

423  and weaknesses of ensemble forecasts used in EFAS, including trends of performance in time

424  and size limits of river basins under monitoring. In addition, it has pointed out a number of

425  key developments to focus on to improve the evaluation and the diagnostic of the forecasting

426  system:

427  - Implementation of the evaluation framework to streamflow predictions derived from all the

428  different numerical weather predictions used as input in EFAS, including DWD, COSMO-

429  LEPS and products which are foreseen to be tested in the future.

430  - Enlarging the collection of near real time observed discharges for continuous monitoring of

431  the skill scores of both the EFAS-WB and streamflow predictions against observed values.

432  - Comparison of performance scores for updated model versions. A new EFAS version was

433  implemented in January 2014, which includes a more extensive calibration of the

434  hydrological model and an enhanced dataset of meteorological observations.

435  - Complementing the current approach with skill scores targeted to evaluate the performance

436  in forecasting extreme events, including threshold exceedance analyses.

437  - Set up a visualization platform on the web where performance can be monitored by

438  developers, analysts on duty and users, to aid the monitoring of forecasts and the diagnostic of

439  issues.

440

**References**

441  **References**

442  Alfieri, L., Burek, P., Dutra, E., Krzeminski, B., Muraro, D., Thielen, J. and Pappenberger, F.,

443  2013a. GloFAS – global ensemble streamflow forecasting and flood early warning, Hydrol.

444  Earth Syst. Sci., 17, 1161–1175, doi:10.5194/hess-17-1161-2013,.

445  Alfieri, L., Salamon, P., Bianchi, A., Neal, J., Bates, P. and Feyen, L., 2013b. Advances in

446  pan-European flood hazard mapping, Hydrol. Process., doi:10.1002/hyp.9947.

447  Alfieri, L., Salamon, P., Pappenberger, F., Wetterhall, F. and Thielen, J., 2012. Operational

448  early warning systems for water-related hazards in Europe, Environ. Sci. Policy, 21, 35–49.

449  Arheimer, B., Lindström, G. and Olsson, J., 2011. A systematic review of sensitivities in the

450  Swedish flood-forecasting system, Atmos. Res., 100, 275–284.

451  Bartholmes, J. C., Thielen, J., Ramos, M. H. and Gentilini, S., 2009. The European flood alert

452  system EFAS - Part 2: Statistical skill assessment of probabilistic and deterministic

453  operational forecasts, Hydrol. Earth Syst. Sc., 13, 141–153.

454  Bartholmes, J., Thielen, J. and Ramos, M. H., 2006. Quantitative analyses of EFAS forecasts

455  using different verification (skill) scores, The benefit of probabilistic flood forecasting on

456  European scale, edited by: Thielen, J., EUR, 22560, 58–79.

457  Bechtold, P., Semane, N., Lopez, P., Chaboureau, J.-P., Beljaars, A. and Bormann, N., 2014.

458  Representing equilibrium and nonequilibrium convection in large-scale models, J. Atmos.

459  Sci., 71, 734–753, doi:10.1175/JAS-D-13-0163.1.

460  Bogner, K. and Pappenberger, F., 2011. Multiscale error analysis, correction, and predictive

461  uncertainty estimation in a flood forecasting system, Water Resour. Res., 47, W07524,

462  doi:10.1029/2010WR009137.

463  Brown, J. D., Demargne, J., Seo, D.-J. and Liu, Y., 2010. The Ensemble Verification System

464  (EVS): A software tool for verifying ensemble forecasts of hydrometeorological and

465  hydrologic variables at discrete locations, Environ. Modell. Softw., 25, 854–872,

466  doi:10.1016/j.envsoft.2010.01.009.

467  Cloke, H. L. and Pappenberger, F., 2009. Ensemble flood forecasting: A review, J. Hydrol.,

468  375, 613–626.

469  Demargne, J., Brown, J., Liu, Y., Seo, D.-J., Wu, L., Toth, Z. and Zhu, Y., 2010. Diagnostic

470  verification of hydrometeorological and hydrologic ensembles, Atmos. Sci. Lett., 11, 114–

471  122.

472  Feyen, L., Vrugt, J. A., Nualláin, B., van der Knijff, J. and De Roo, A., 2007. Parameter

473  optimisation and uncertainty assessment for large-scale streamflow simulation with the

474  LISFLOOD model, J. Hydrol., 332, 276–289, doi:10.1016/j.jhydrol.2006.07.004.

475  Gourley, J. J., Erlingis, J. M., Hong, Y. and Wells, E. B., 2012. Evaluation of Tools Used for

476  Monitoring and Forecasting Flash Floods in the United States, Weather Forecast., 27, 158–

477  173, doi:10.1175/WAF-D-10-05043.1.

478  Hersbach, H., 2000. Decomposition of the continuous ranked probability score for ensemble

479  prediction systems, Weather Forecast., 15, 559–570.

480  Kann, A. and Haiden, T., 2005. The August 2002 flood in Austria: sensitivity of precipitation

481  forecast skill to areal and temporal averaging, Meteorol. Z., 14, 369–377, doi:10.1127/0941-

482  2948/2005/0042.

483  Van der Knijff, J. M., Younis, J. and de Roo, A. P. J., 2010. LISFLOOD: A GIS-based

484  distributed model for river basin scale water balance and flood simulation, Int. J. Geogr. Inf.

485  Sci., 24, 189–212.

486   Legates, D. R. and McCabe, G. J., 1999. Evaluating the use of "goodness-of-fit" Measures in

487   hydrologic and hydroclimatic model validation, Water Resour. Res., 35, 233–241,

488   doi:10.1029/1998WR900018.

489   Majewski, D., Liermann, D., Prohl, P., Ritter, B., Buchhold, M., Hanisch, T., Paul, G.,

490   Wergen, W. and Baumgardner, J., 2002. The operational global icosahedral-hexagonal

491   gridpoint model GME: Description and high-resolution tests, Mon. Weather Rev., 130, 319–

492   338.

493   Marsigli, C., Boccanera, F., Montani, A. and Paccagnella, T., 2005. The COSMO-LEPS

494   mesoscale ensemble system: Validation of the methodology and verification, Nonlinear Proc.

495   Geoph., 12, 527–536.

496   Miller, M., Buizza, R., Haseler, J., Hortal, M., Janssen, P. and Untch, A., 2010. Increased

497   resolution in the ECMWF deterministic and ensemble prediction systems, ECMWF

498   Newsletter, 124, 10–16.

499   Nash, J. E. and Sutcliffe, J. V., 1970. River flow forecasting through conceptual models part

500   I–A discussion of principles, J. Hydrol., 10, 282–290.

501   Ntegeka, V., Salamon, P., Gomes, G., Sint, H., Lorini, V. and Thielen, J., 2013. A high-

502   resolution European dataset for hydrologic modeling, in EGU General Assembly Conference

503   Abstracts, 15, 7460.

504   Pagano, T. C., 2013. Evaluation of Mekong River Commission operational flood forecasts,

505   2000–2012, Hydrol. Earth Syst. Sci. Discuss., 10, 14433–14461, doi:10.5194/hessd-10-

506   14433-2013.

507  Pappenberger, F., Ramos, M. H., Cloke, H. L., Wetterhall, F., Alfieri, L., Bogner, K.,

508  Mueller, A., and Salamon, P., submitted. How do I know if my forecasts are better? Using

509  benchmarks in Hydrological Ensemble Predictions, J. Hydrol.

510  Pappenberger, F., Thielen, J. and Del Medico, M., 2010. The impact of weather forecast

511  improvements on large scale hydrology: analysing a decade of forecasts of the European

512  Flood Alert System, Hydrol. Process., 25, 1091–1113, doi:10.1002/hyp.7772.

513  Pappenberger, F., Wetterhall, F., Albergel, C., Alfieri, L., Balsamo, G., Bogner, K., Haiden,

514  T., Hewson, T., Magnusson, L., de Rosnay, P., Muñoz Sabater, J. and Tsonevsky, I., 2013.

515  Floods in Central Europe in June 2013, ECMWF Newsletter, 136, 9–11.

516  Plate, E. J. and Lindenmaier, F., 2008. Quality assessment of forecasts, in Mekong River

517  Commission: Sixth Annual Flood Forum, Phnom Penh, May, 10 pp.

518  Randrianasolo, A., Ramos, M. H., Thirel, G., Andreassian, V. and Martin, E., 2010.

519  Comparing the scores of hydrological ensemble forecasts issued by two different hydrological

520  models, Atmos. Sci. Lett., 11, 100–107, doi:10.1002/asl.259.

521  Reed, S., Schaake, J. and Zhang, Z., 2007. A distributed hydrologic model and threshold

522  frequency-based method for flash flood forecasting at ungauged locations, J. Hydrol., 337,

523  402–420.

524  Renner, M., Werner, M. G. F., Rademacher, S. and Sprokkereef, E., 2009. Verification of

525  ensemble flow forecasts for the River Rhine, J. Hydrol., 376, 463–475.

526  Steppeler, J., Doms, G., Schättler, U., Bitzer, H. W., Gassmann, A., Damrath, U. and

527  Gregoric, G., 2003. Meso-gamma scale forecasts using the nonhydrostatic model LM,

528  Meteorol. Atmos. Phys., 82, 75–96.

529   Thielen, J., Bartholmes, J., Ramos, M.-H. and De Roo, A., 2009. The European flood alert

530   system - part 1: Concept and development, Hydrol. Earth Syst. Sc., 13, 125–140.

531   Thirel, G., Rousset-Regimbeau, F., Martin, E. and Habets, F., 2008. On the impact of short-

532   range meteorological forecasts for ensemble streamflow predictions, J. Hydrometeorol., 9,

533   1301–1317.

534   Trinh, B. N., Thielen-del Pozo, J. and Thirel, G., 2013. The reduction continuous rank

535   probability score for evaluating discharge forecasts from hydrological ensemble prediction

536   systems, Atmos. Sci. Lett., 14, 61–65, doi:10.1002/asl2.417.

537   Verkade, J. S., Brown, J. D., Reggiani, P. and Weerts, A. H., 2013. Post-processing ECMWF

538   precipitation and temperature ensemble reforecasts for operational hydrologic forecasting at

539   various spatial scales, J. Hydrol., 501, 73–91, doi:10.1016/j.jhydrol.2013.07.039.

540   Wetterhall, F., Pappenberger, F., Alfieri, L., Cloke, H. L., Thielen-del Pozo, J., Balabanova,

541   S., Daňhelka, J., Vogelbacher, A., Salamon, P., Carrasco, I., Cabrera-Tordera, A. J., Corzo-

542   Toscano, M., Garcia-Padilla, M., Garcia-Sanchez, R. J., Ardilouze, C., Jurela, S., Terek, B.,

543   Csik, A., Casey, J., Stankūnavičius, G., Ceres, V., Sprokkereef, E., Stam, J., Anghel, E.,

544   Vladikovic, D., Alionte Eklund, C., Hjerdt, N., Djerv, H., Holmberg, F., Nilsson, J., Nyström,

545   K., Sušnik, M., Hazlinger, M. and Holubecka, M., 2013. HESS Opinions "Forecaster

546   priorities for improving probabilistic flood forecasts," Hydrol. Earth Syst. Sci., 17, 4389–

547   4399, doi:10.5194/hess-17-4389-2013.

548   Wilks, D. S., 2006. Statistical Methods in the Atmospheric Sciences: An Introduction,

549   electronic version, Elsevier, San Diego, CA.

550   Wittmann, C., Haiden, T. and Kann, A., 2010. Evaluating multi-scale precipitation forecasts

551   using high resolution analysis, Adv. Sci. Res., 4, 89–98.

552   Wood, A. W., Kumar, A. and Lettenmaier, D. P., 2005. A retrospective assessment of

553   National Centers for Environmental Prediction climate model–based ensemble hydrologic

554   forecasting in the western United States, J. Geophys. Res-Atmos., 110,

555   doi:10.1029/2004JD004508.

556

557

558 **Tables**

559

560 Table 1: Summary of performance scores and their information content.

| Score | Short name | Use |
|---|---|---|
| Nash-Sutcliffe efficiency | NS | Normalized measure of the mean squared error of the ensemble mean in comparison to a constant climatological mean |
| Percent bias | Pbias | Dimensionless measure of the forecast bias |
| Coefficient of variation of the Root Mean Squared Error | CV | Dimensionless measure of the Root Mean Squared Error of the ensemble mean |
| Continuous Ranked Probability Skill Score (average discharge as reference) | $CRPSS_{ad}$ | Skill score to compare the distribution of ensemble forecasts around observations, as opposed to using the climatological average discharge |
| Continuous Ranked Probability Skill Score (persistence forecast as reference) | $CRPSS_{pf}$ | Skill score to compare the distribution of ensemble forecasts around observations, as opposed to using the persistence of the initial discharge |

561 **Figure captions**

562

563 Figure 1: Schematic view of the EFAS hydro-meteorological forecasting system.
564

565 Figure 2: stations reporting observed precipitation (left) and average temperature (right) on

566 the 1st October 2013.

567

568 Figure 3: $CRPSS_{pf}$, CV, NS and Pbias over Europe for 1 year of daily forecasts ending on the

569 1st October 2013 (5-day lead time).

570

571 Figure 4: CRPSS, CV, NS and Pbias of ESP versus the forecast lead time.

572

573 Figure 5: $CRPSS_{pf}$, CV, NS and Pbias of ESP versus the upstream area of each river grid

574 point.

575

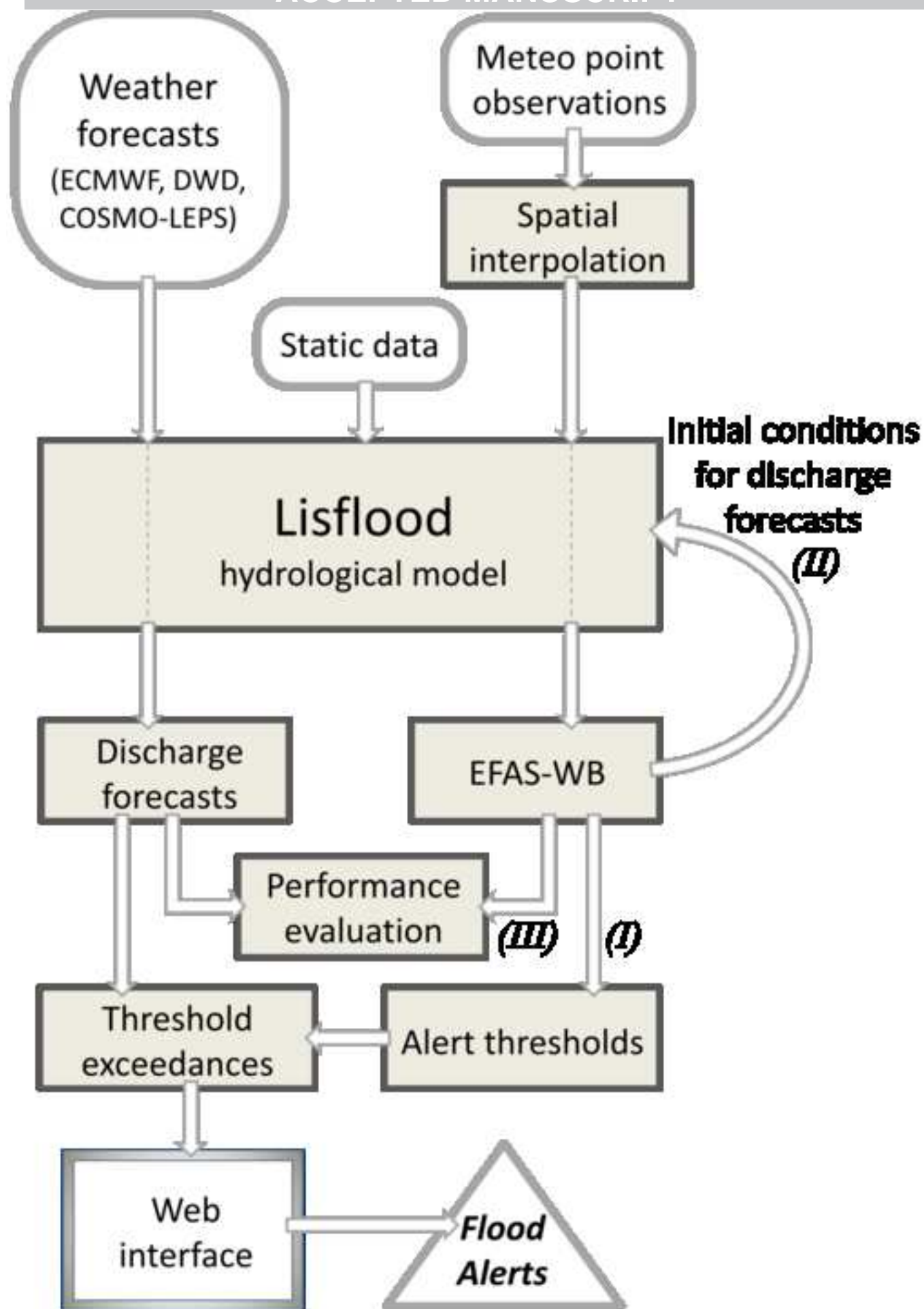576 Figure 6: Trend of 12-month average CRPSS, CV, NS and Pbias of ESP from 2009 onwards.
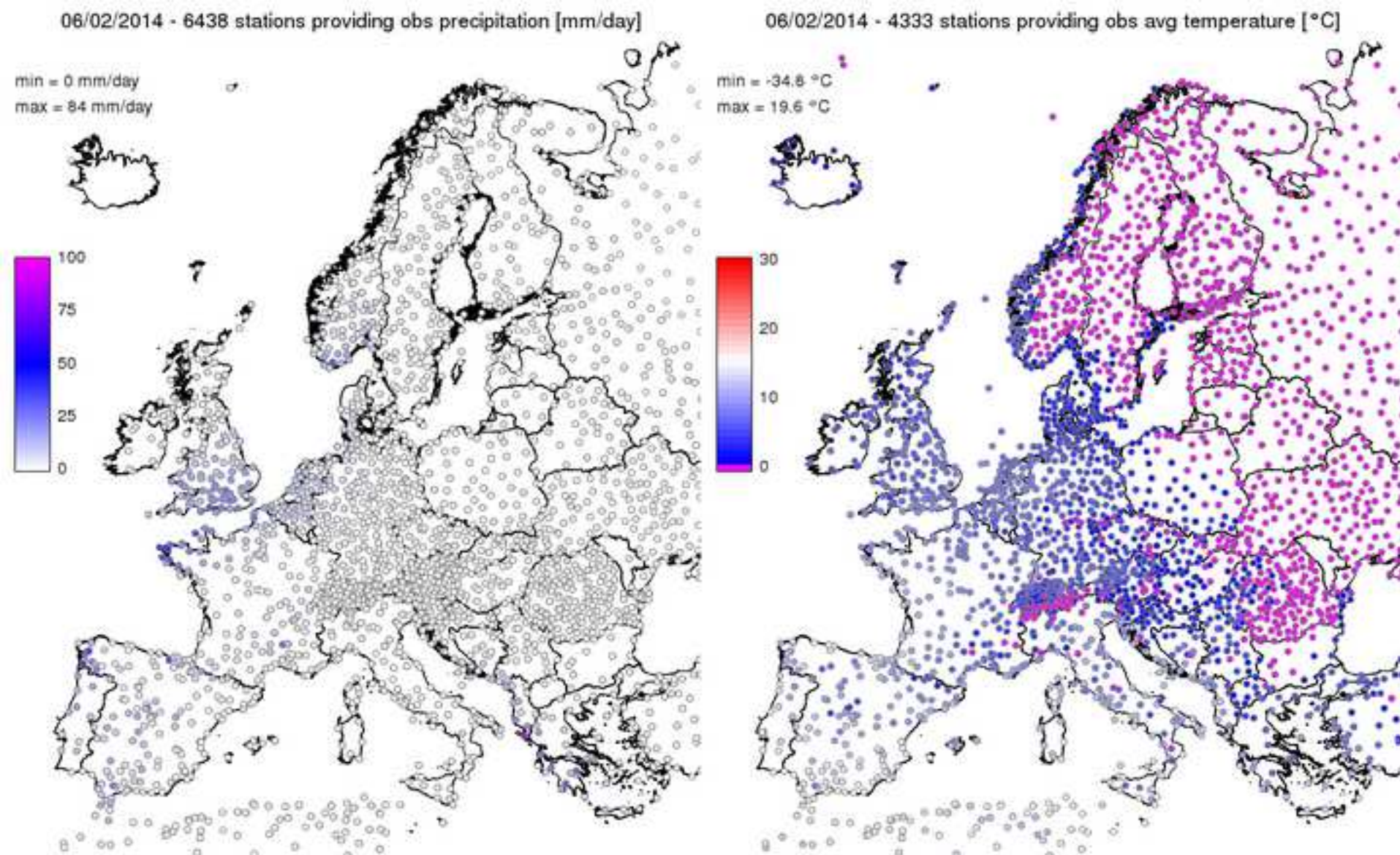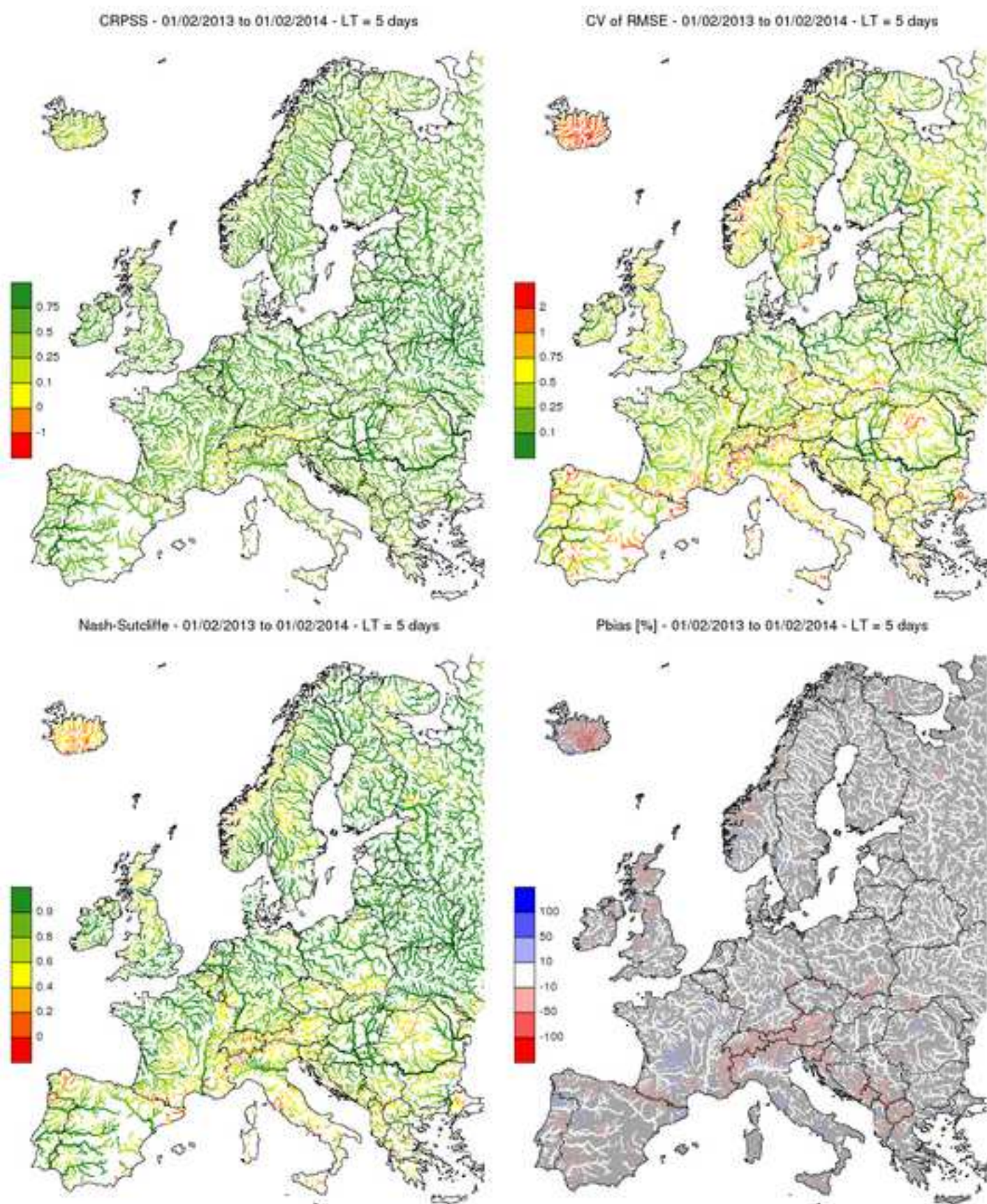
577

**Figure 1**

**Figure 2**

06/02/2014 - 6438 stations providing obs precipitation [mm/day]

min = 0 mm/day
max = 84 mm/day

06/02/2014 - 4333 stations providing obs avg temperature [°C]

min = -34.8 °C
max = 19.6 °C

**Figure 3**

CRPSS - 01/02/2013 to 01/02/2014 - LT = 5 days

CV of RMSE - 01/02/2013 to 01/02/2014 - LT = 5 days

Nash-Sutcliffe - 01/02/2013 to 01/02/2014 - LT = 5 days

Pbias [%] - 01/02/2013 to 01/02/2014 - LT = 5 days

**Figure 4**

**Figure 5**



01/02/2013 to 01/02/2014 · LT = 5 days

01/02/2013 to 01/02/2014 · LT = 5 days

01/02/2013 to 01/02/2014 · LT = 5 days

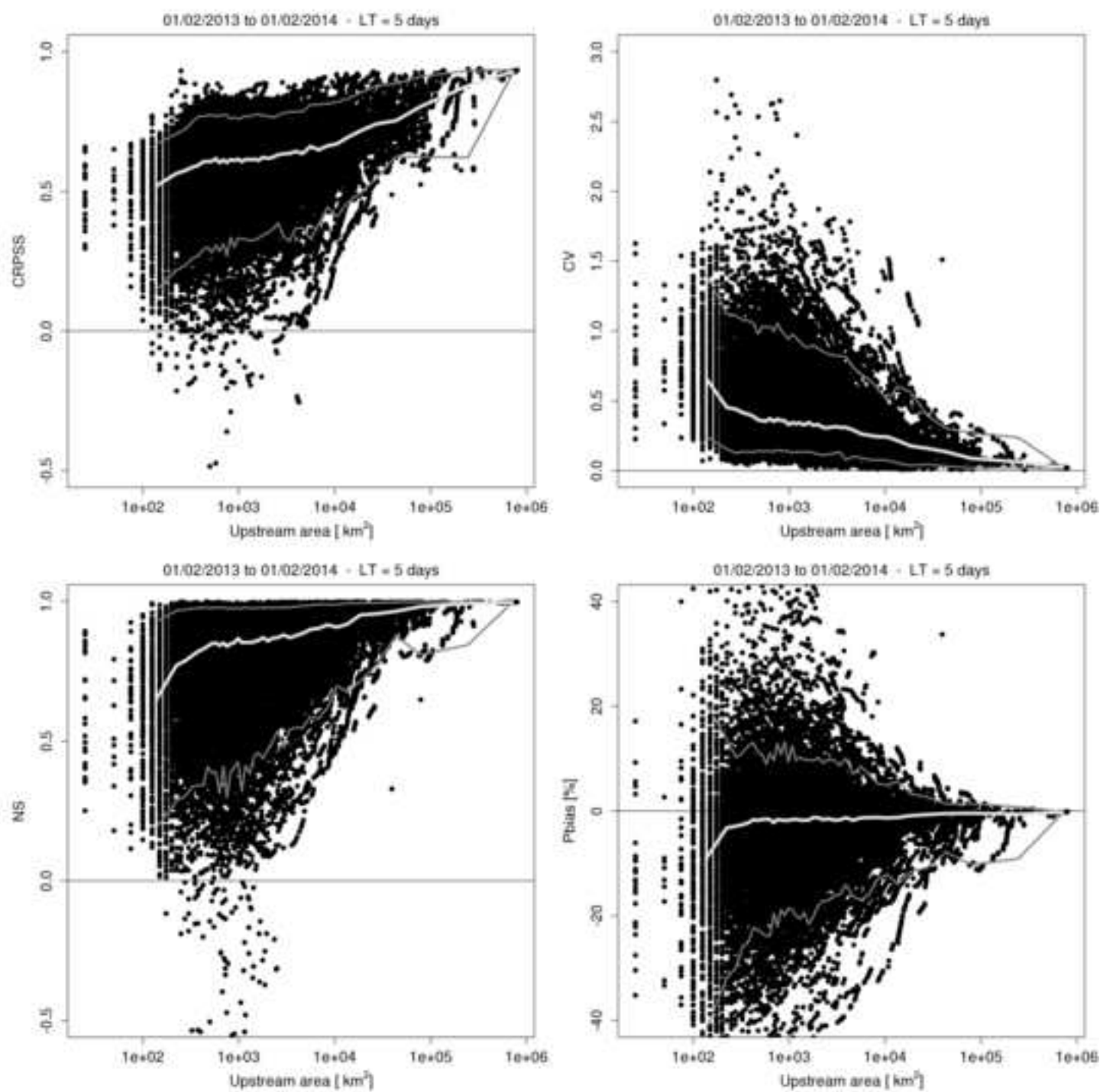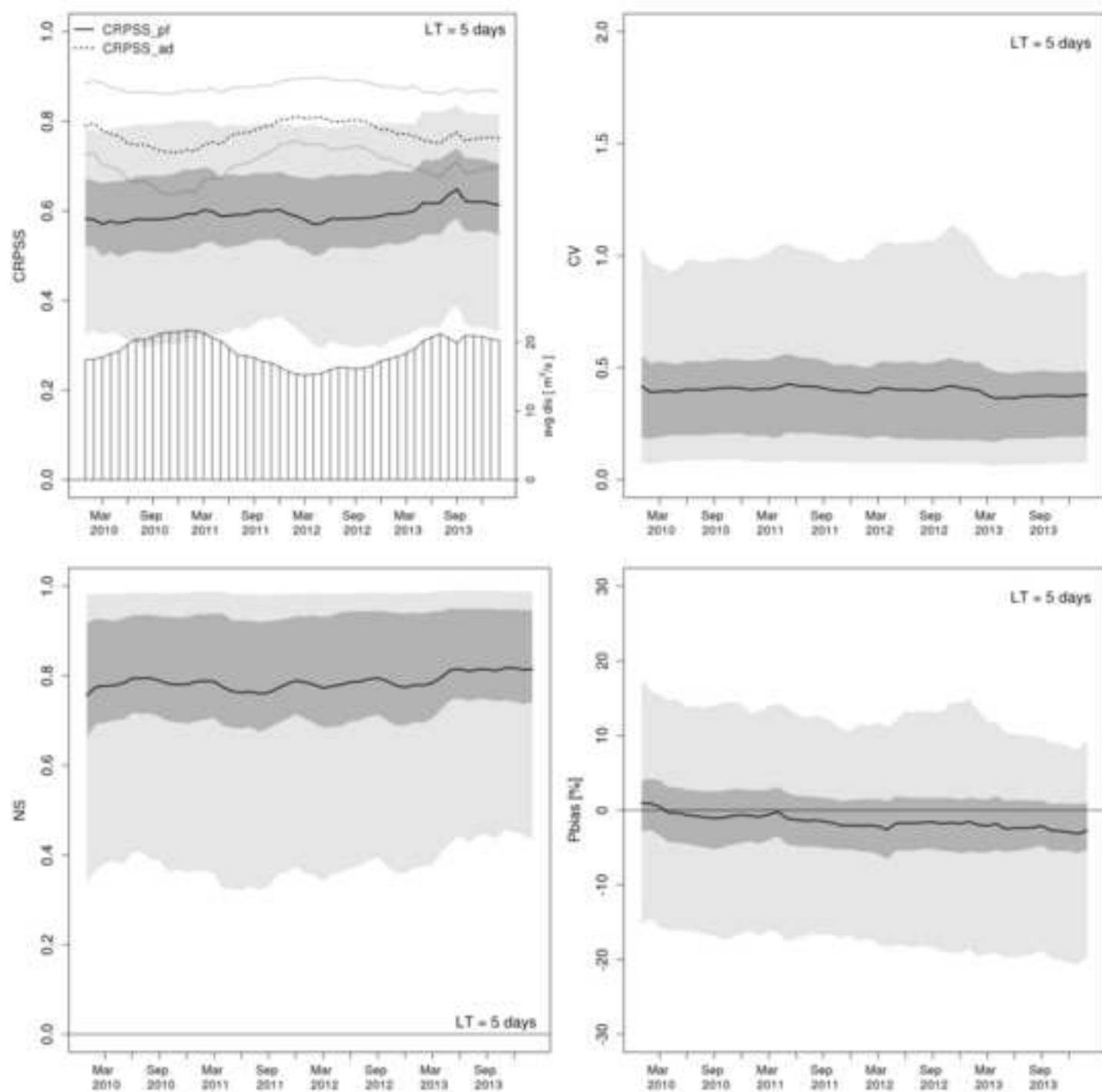01/02/2013 to 01/02/2014 · LT = 5 days

**Figure 6**

578

**Research Highlights**

580

581    − The evaluation framework of the European Flood Awareness
582       System is presented

583

584    − Skill scores of ensemble streamflow predictions over Europe are
585       updated regularly

586

587    − Predictions are skillful in river basins larger than 300 km$^2$ over
588       the 10-day range

589

590    − The use of the CRPSS based on two different references is
591       discussed

592