

PSO (FU 2101)  
Ensemble-forecasts for wind power  
**Comparison of ensemble forecasts with the  
measurements from the meteorological mast  
at Risø National Laboratory**

Henrik Aalborg Nielsen<sup>1</sup>, Henrik Madsen<sup>1</sup>, Torben Skov Nielsen<sup>1</sup>,  
Jake Badger<sup>2</sup>, Gregor Giebel<sup>2</sup>, Lars Landberg<sup>2</sup>  
Kai Sattler<sup>3</sup>, Henrik Feddersen<sup>3</sup>

<sup>1</sup> Technical University of Denmark, Informatics and Mathematical Modelling, DK-2800 Lyngby

<sup>2</sup> Risø National Laboratory, Wind Energy Department, DK-4000 Roskilde

<sup>3</sup> Danish Meteorological Institute, Research and Development, DK-2100 København Ø

23rd March 2004

# Contents

<b>1</b>	<b>Summary</b>	<b>4</b>
<b>2</b>	<b>Introduction</b>	<b>4</b>
<b>3</b>	<b>Ensemble forecasts</b>	<b>5</b>
3.1	NCEP ensemble forecasts . . . . .	5
3.2	ECMWF medium range ensemble forecasts . . . . .	6
3.3	HIRLAM short range ensemble forecasts . . . . .	6
<b>4</b>	<b>Wind data and ensemble forecast data</b>	<b>7</b>
<b>5</b>	<b>Some properties of the ensemble forecasts</b>	<b>8</b>
<b>6</b>	<b>Comparison of ensemble forecasts with measurements</b>	<b>10</b>
6.1	Reliability . . . . .	11
6.1.1	Unadjusted ensemble forecasts . . . . .	11
6.1.2	MOS adjusted ensemble forecasts . . . . .	15
6.2	Resolution . . . . .	23
<b>7</b>	<b>Conclusion and Discussion</b>	<b>29</b>
	<b>References</b>	<b>34</b>
	<b>Appendices</b>	<b>35</b>
<b>A</b>	<b>Plots of MOS adjusted ensemble forecasts</b>	<b>35</b>
A.1	ECMWF corrected using the unperturbed forecast as the dependent variable.	35

A.2 HIRLAM corrected using the unperturbed forecast as the dependent variable. 42

## Preface

This report is mainly based on investigations carried out in the autumn of 2003. The report describes the first attempt of comparing measurements (here of wind speed) with ensemble forecasts. Three types of ensemble forecasts are considered. For one of these (NCEP) it was discovered after completion of the work that the unperturbed forecast with initialization time 12:00 (UTC) is actually a modified version of the other models in this ensemble. This may explain some of the observations reported.

# 1 Summary

Three types of ensemble forecasts of wind speed are compared with measurements 76m a.g.l. Specifically, the comparisons are performed in terms of *reliability* (average correctness in a probabilistic sense) and *resolution* (sharpness of the conditional density function). Ensembles from the experimental HIRLAM ensemble forecast system seem to be reliable w.r.t. the upper quantiles ( $\geq 80\%$ ), but the number of cases is only 86 and on average only 17.2 should exceed the 80% quantile. The resolution of the 80% quantile is good; it varies from 1.7 to 25.7 m/s. The NCEP and ECMWF ensemble forecasts are not reliable for the point measurement.

Since the forecasts should be interpreted as spatial averages the forecasts are adjusted using MOS (Model Output Statistics) and then compared with the point measurements. It is argued that the usual approach to MOS adjustment possess properties which are somewhat in conflict with ensemble forecasting. These problems are demonstrated and alternatives are suggested. However, for the short horizons ( $< 36$  hours), all methods yields adjusted ensemble forecasts for which the spread is too small. The spread of the adjusted NCEP ensembles are in general too small. For ECMWF and HIRLAM and horizons approximately in the range 36–60 hours some of the alternative methods suggested yields reliable ensemble forecasts. For longer horizons the spread of the adjusted ensemble forecasts are too wide.

Unfortunately, the methods which yields reliable ensemble forecasts have problems w.r.t. numerical instabilities, whereby only horizons up to approximately 48 hours have good resolution. The report discuss ways to circumvent the numerical instabilities, which mainly are caused by uncertainty on the unperturbed forecast and inefficient use of the measurements.

## 2 Introduction

This report documents investigations carried out as part of the PSO-funded project *Ensemble-forecasts for wind power* (FU 2101).

In meteorological ensemble forecasting several forecasts are produced and it is hoped that these are informative w.r.t. the predictability of the actual weather situation. For wind power applications it is not sufficient to obtain forecasts of the near-extremes; it is necessary to be able to interpret the ensemble forecasts in a probabilistic sense. In order to investigate the extend to which this is possible three types of ensemble forecasts are compared with actual measurements of wind speed 76m a.g.l. from a mast at Risø National Laboratory.

The wind data and ensembles are briefly described in Section 4 and some background information on the ensemble forecast systems are given in Section 3. Section 5 describes some properties of the ensemble forecasts which can be observed in the forecast-data. In Section 6 the forecasts and MOS adjusted versions of these are compared with the measurements w.r.t. reliability (Section 6.1) and resolution (Section 6.2). Finally, in Section 7, we conclude and discuss on the findings.

### 3 Ensemble forecasts

The investigations in this report make use of ensemble forecast data from numerical weather prediction (NWP) models of three meteorological centres. Two of the NWP models are global, and the third one is a limited-area model. The global ensemble prediction systems have been in operation since about a decade, within which they were continuously further developed. They are especially designed for the medium forecast range with lead times beyond 3 days. The third ensemble is experimental and is designed for the short range below 3 days lead time.

The ensemble prediction systems address the sensitivity of the weather development to the initial condition by making use of perturbed initial conditions in order to account for the uncertainty in the analysis of the atmospheric state, from which the NWP simulations for the forecast are to be started. Common for all three ensemble prediction systems is that they include the “unperturbed” forecast, which takes the original “best-guess” analysis as initial condition. It is referred to as the control forecast.

In addition to the initial condition perturbation the ensemble prediction systems make use of either different components of the NWP model or try to perturb tendencies of model parameters in order to address uncertainties that arise during the model integration and due to the model formulation.

Some major aspects of the ensemble prediction systems are outlined in the following sub-sections.

#### 3.1 NCEP ensemble forecasts

The National Center for Environmental Prediction (NCEP) and the National Weather Service of NOAA (National Oceanic and Atmospheric Administration) in the USA have operated an ensemble prediction system since many years [11]. It includes a global as well as regional NWP models. The set of ensemble members consists of one unperturbed forecast and 5 pairs of perturbed forecasts. For each pair the initial condition is perturbed in the positive and negative direction of bred vectors [12], which are determined within

the analysis/forecasting cycle and which are regarded to reflect the largest state changes of the model atmosphere within the cycle time.

The ensemble system comprises model simulations in different horizontal resolutions and with different lead times of up to 16 days. In this study, the interpolated medium range ensemble data on a  $1^\circ$  horizontal grid are utilized with lead times up to 3.5 days.

### 3.2 ECMWF medium range ensemble forecasts

The European Centre for Medium-Range Weather Forecasts (ECMWF) has operated their Ensemble prediction system<sup>1</sup> (ECMWF-EPS) since the early 1990s [9, 6]. The set of ensemble members consists of one unperturbed forecast and 25 pairs of forecasts, for which the initial conditions are perturbed in the positive and negative direction of the analysis error on basis of singular vectors [2]. The singular vectors used in the ECMWF-EPS describe the largest growth of initial state differences in accordance with the analysis error covariances within a time of 48 hours, in which uncertainties are assumed to grow linearly. The initial perturbations are calculated as linear combinations of the singular vectors such that they cover most of the global area. They are scaled such that their amplitude is comparable to the root-mean-square analysis error of the data assimilation system.

Furthermore, for each model run attempts are made to account for uncertainties in the description and calculation of sub-grid processes by use of stochastic physics [1]. These affect the model dynamics by adding a stochastic component to the tendencies from the sub-grid parameterizations of the model. As a result two runs with the same set of initial conditions will not result in exactly the same output.

The NWP model of the ECMWF-EPS is a global spectral model with truncation at wave number 255 and with 40 vertical levels (TL255L40). The horizontal resolution is comparable to that of a grid with about 80km mesh size. The ensemble simulations are carried out operationally up to a lead time of 10 days.

### 3.3 HIRLAM short range ensemble forecasts

High resolution short range ensemble forecasts using the High Resolution Limited-Area Model HIRLAM have been provided by the Danish Meteorological Institute (DMI). These are experimental ensembles and they are based on the ECMWF-EPS. They adopt a nested model system from a HIRLAM model setup of a previous study with HIRLAM ensembles [10]. The forecast simulations have been performed on a grid with approximately 20km in mesh size on a domain, which includes Denmark and the North Sea.

---

<sup>1</sup>[http://www.ecmwf.int/products/forecasts/guide/The\\_Ensemble\\_Prediction\\_System\\_EPS.html](http://www.ecmwf.int/products/forecasts/guide/The_Ensemble_Prediction_System_EPS.html)

In a limited-area model, the boundary conditions need to be prescribed in addition to the initial condition. As a simple solution, the members from the ECMWF-EPS have therefore been chosen as host model. They can supply initial perturbations with consistent boundary data during the limited-area model simulation. Uncertainties in the initial and in the boundary conditions are addressed this way to some extent. The frequency of availability of the boundary data has been 6 hours. The size of this HIRLAM ensemble is the same as that of the ECMWF-EPS, namely 50 plus the control forecast.

In addition to the ECMWF-EPS members, four further model permutations have been integrated, too. These consist of the control simulation being integrated with different parameterization schemes for convection and condensation in order to address model uncertainties related to these processes. The data from these simulations were, however, not utilized in this study.

## 4 Wind data and ensemble forecast data

The collection of the wind data and the ensemble forecasts (Sec. 3) has been performed over several months. For the high resolution HIRLAM ensembles the period is limited to comprise four months. The large scale ensemble forecasts, however, cover longer periods. All forecast data are gridded and represented in a geographic projection. Moreover, the grid of the HIRLAM ensembles is rotated leading to the grid area being less dependent on the geographic position. Table 1 lists some of the major characteristics of the forecast data used for the investigations at the Risø mast in this study.

The gridded ensemble forecast data has been interpolated horizontally to the location of the Risø mast by bilinear interpolation between the four closest grid points. A vertical interpolation has not been performed.

	NCEP	ECMWF	HIRLAM <sup>2</sup>
wind	10m a.g.l.	10m a.g.l.	10m, ≈70m* a.g.l.
grid spacing	1°	0.75°	0.2°
period start	01/11/2002 00:00 UTC	01/02/2003 12:00 UTC	01/12/2002 12:00 UTC
period end	04/09/2003 00:00 UTC	07/09/2003 12:00 UTC	29/03/2003 12:00 UTC
forecast initiation	00:00 and 12:00 UTC	12:00 UTC	12:00 UTC
lead time	3.5 days	7 days	3 days
data frequency	6 hours	6 hours	1 hour

\* this data is first available from 11/01/2003

Table 1: Overview of the ensemble forecast data

Wind observations from the mast at Risø (UTM zone 32: 694197E, 6176579N) were taken as verifying data. They include measurements of wind speed and direction at 76m

<sup>2</sup>This system is experimental. It is based on the ECMWF-EPS.

a.g.l., and they cover the period from 31/01/2002 23:06 to 01/09/2003 10:59 (UTC). The measurements are averages over the 10 minutes up to the time stamp, and the ensemble forecast data are interpreted the same way. In the case of the data from the NCEP ensembles and for the ECMWF-EPS the time step of the model simulations is larger than the measurement interval. In the case of the HIRLAM ensembles the time step of the model simulations is smaller than the measurement interval. The data was therefore averaged over several time steps to represent an interval of 10 minutes. However, some uncertainty in representation remains for all models, because the ensemble forecasts correspond to spatial averages over the grid size of the respective model. This blurs the picture of temporal representation, even though spatial interpolation as mentioned above is used to obtain ensemble forecasts for the Risø mast.

## 5 Some properties of the ensemble forecasts

Figure 1 shows the mean and quantiles of each of the ensemble forecasts for each horizon. Interestingly the mean and median of the NCEP analysis, i.e. forecast horizon 0 hours, is approximately 2 m/s lower than the remaining mean values.

For ECMWF the mean and median has a cyclic behavior with peaks occurring at horizons 0, 12, 24, 36, ... hours. For the median the distance between the top and bottom is approximately 0.5 m/s and slightly less for the mean. Since ECMWF is only initiated once daily this may just be a consequence of the diurnal variation. Actually, if NCEP forecasts are split in two groups according to the initialization of calculations (00:00 and 12:00 UTC) the same kind of behavior is observed, but to a somewhat less extend (plots not shown).

For HIRLAM within the first 12 hours, which corresponds to the time interval 12:00 – 24:00 UTC, the mean value drops from 9.5 m/s at 0 hours (12:00) to 7.2 m/s at 3 hours (15:00). Hereafter, the mean increase to 8.0 m/s at 8 hours (19:00). Note also that for HIRLAM the 25% and 75% quantiles seems to be in opposite phase, whereas these quantiles seems to be in phase for ECMWF.

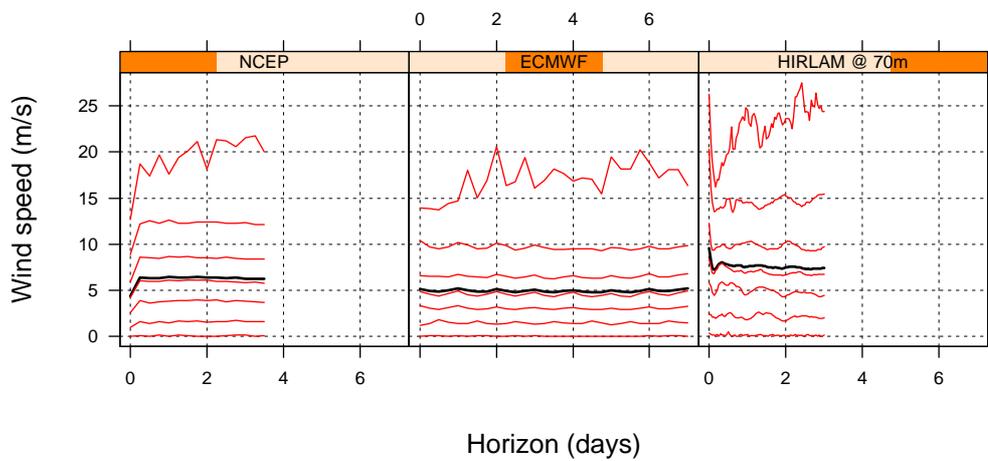
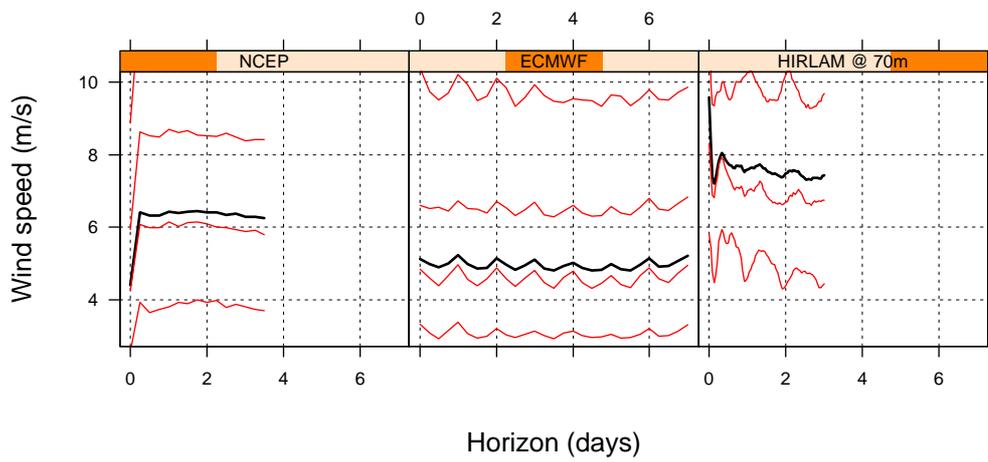


Figure 1: Mean (black) and quantiles (red) of ensemble forecasts for each horizon. The quantiles 0, 0.05, 0.25, 0.50, 0.75, 0.95, and 1 are depicted. The top row displays part of the data in the bottom row (approx. 3-10 m/s).

## 6 Comparison of ensemble forecasts with measurements

For wind power production applications we are interested in interpreting the ensemble forecast in a probabilistic sense. For example if 10 out of 50 ensemble members show wind speeds above 10 m/s at a particular point in the future then we would like to interpret this as a 20% chance of wind speeds above 10 m/s at that particular point in the future. If this property holds on average for all thresholds then the ensemble forecasts are called *reliable*. If an ensemble forecast system is reliable as compared to a particular measurement then the rank of the measurement, when compared with the ensemble members, is uniformly distributed. This can be investigated by plotting a histogram of the ranks which should be fairly flat [13] or by Quantile-Quantile-plots (QQ-plots) [3] where the observed quantiles are plotted against the theoretical (uniform) distribution.

However, if climatological information is available *reliable* (but uninformative for practical use) ensemble forecasts can easily be produced by random sampling from the observed distribution. Figure 2 exemplifies this by using measurements before the first NCEP forecast as climatological information<sup>3</sup>. The obtained rank histogram is fairly flat except that there is a slight over-representation of observations with low ranks; indicating that the observed wind speeds after 1/11/2002 is somewhat lower than before 1/11/2002. This is probably due to a relatively windy period in February and March, 2002. Note also that using the empirical cumulative distribution in the left panel of Figure 2 results in a histogram of essentially the same shape; i.e. the main features of the right panel of the figure is not due to random variation.

The problem with the ensemble forecasts generated using climatology is that the uncertainty indicated by the ensemble is high. For instance the 10% and 90% quantiles in the left panel of Figure 2 is 3.0 and 11.3 m/s, respectively. This feature of “sharpness” of the density indicated by the ensemble of forecasts is called *resolution*. The ultimate goal is to have reliable ensemble forecasts with high resolution.

When comparing ensemble forecasts with a point measurement it should be noted that the forecasts are forecasts of a spatial average corresponding to the horizontal resolution of the model. The measurements at a particular point may well deviate systematically from the spatial average.

### **Remark:**

Here the terms *reliability* and *resolution* are for continuous variables as described above. It seems that in the meteorological literature (e.g. [13]) these words is often used for the more restricted case of binary variables.

---

<sup>3</sup>This definition of climatology is used throughout this report, it covers the period 31/01/2002 23:06 – 31/10/2002 23:57 (UTC).

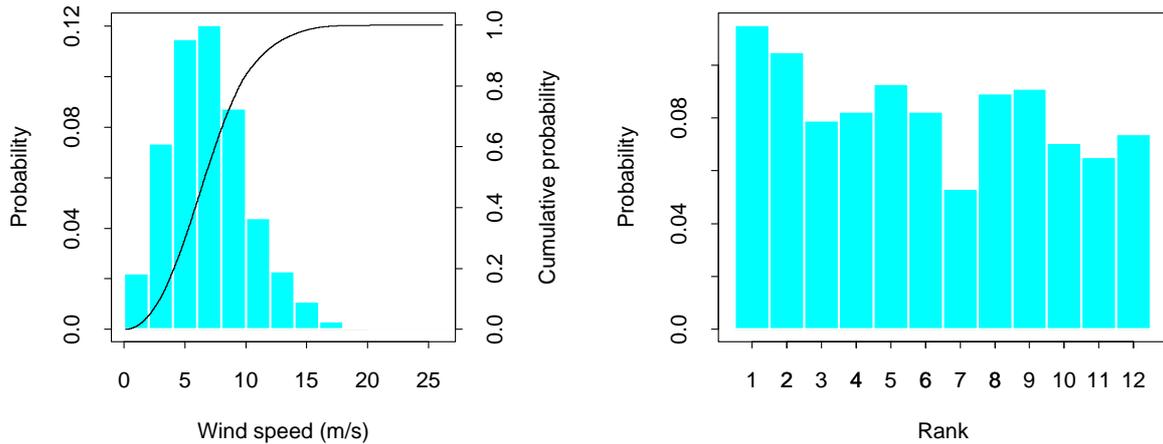


Figure 2: Left: Histogram and empirical cumulative distribution for the measurements from the Risø mast before the first NCEP forecast (1/11/2002). Right: Rank histogram obtained for a climatological ensemble forecast obtained by randomly picking 11 samples in the empirical distribution at the left.

## 6.1 Reliability

In this section the reliability of the ensemble forecasts or ensemble forecasts adjusted using MOS (Model Output Statistics) are investigated. In all cases only forecasts with all members present will be considered, i.e. 11 members for NCEP and 51 members for ECMWF and HIRLAM. The model permutations included in the HIRLAM ensembles are not included in the investigation.

### 6.1.1 Unadjusted ensemble forecasts

Figures 3 and 4 show rank histograms when comparing the three ensemble forecasts to the measurements at the Risø mast. In case of mismatch between time points observations are generated for the forecast time points by linear interpolation between actual time points. To account for unequal number of ensembles between ensemble types the rank is normalized as indicated in the figures. The number of bins are selected as the default for the function `histogram` in S-PLUS.

Generally, for all horizons ECMWF ensembles seem to have an over-representation of high ranks indicating some downward bias of the forecasts as compared to the measurements. For HIRLAM the forecasts for the low horizons are generally too high as compared to the measurements, whereas the rank histograms are fairly flat for the longer horizons. For NCEP, generally, there is a tendency for U-shaped rank histograms, indicating that the spread of the ensemble forecasts are too low as compared to the measurements. However, for the longer horizons the rank histograms corresponding to the NCEP ensembles are fairly flat.

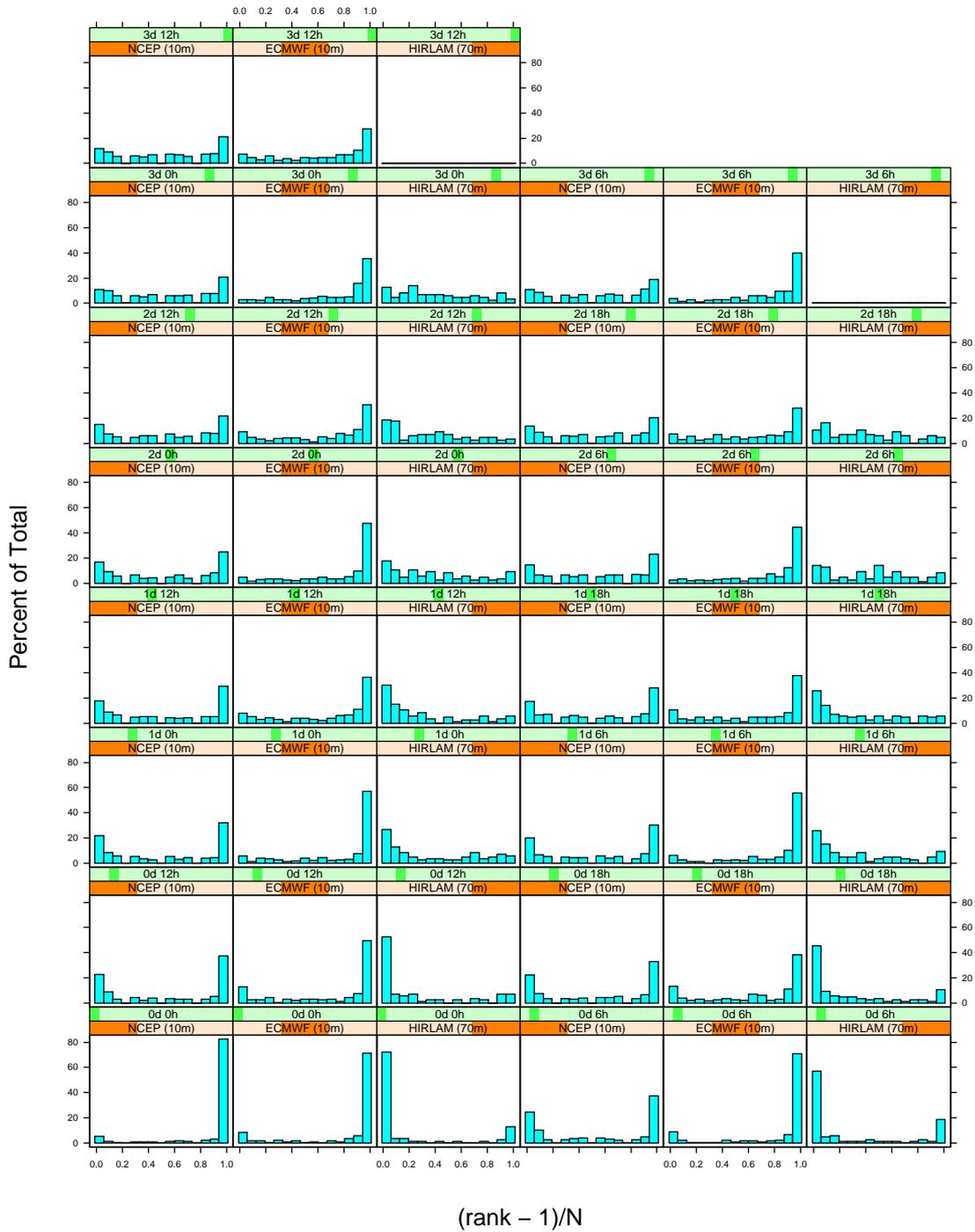


Figure 3: Rank histograms for horizons up to 84 hours. Note that HIRLAM forecasts are not available after 72 hours.

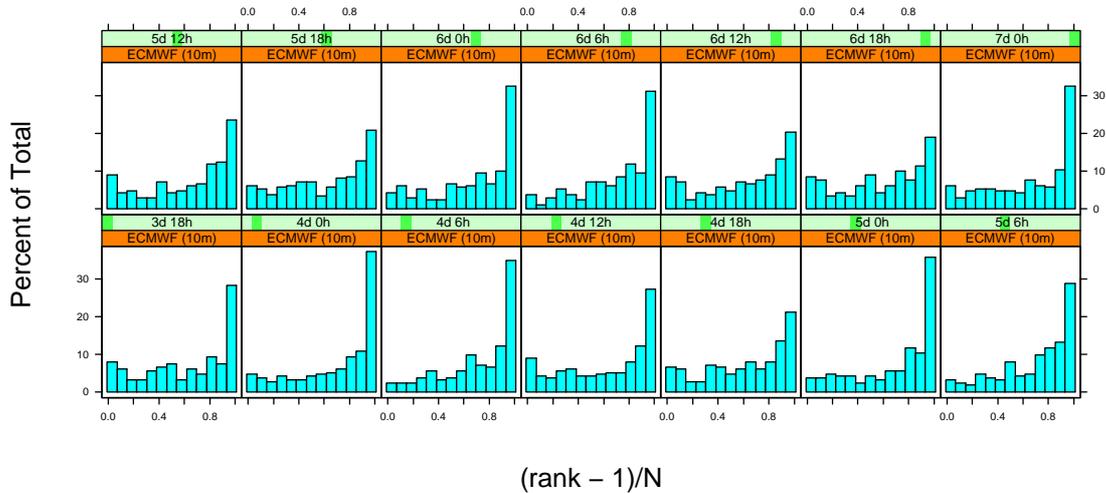


Figure 4: ECMWF rank histograms for horizons longer than 84 hours.

Although, rank histograms give some indication of the overall reliability of the ensemble forecasts it is unclear when the rank histograms are “sufficiently uniform”. For this purpose we propose to use QQ-plots where the normalized rank is plotted against the theoretical quantiles of the uniform distribution on  $[0, 1]$  ( $U(0, 1)$ ). Since the cumulative distribution function of  $U(0, 1)$  is a line connecting  $(0, 0)$  and  $(1, 1)$  the plots are particular simple to interpret in this case.

Figure 5 shows QQ-plots for selected horizons. Consider e.g. the 12h NCEP ensemble forecast; the maximum of this forecast (2nd axis at 1.0) is exceeded in approximately 40% of the cases ( $1.0 - 0.6$  on the 1st axis). Likewise, when the 60h HIRLAM ensemble forecast indicates that there is a 40% chance that a certain threshold will be exceeded ( $1 - 0.6$  on the 2nd axis) then the data suggests that this threshold will only be exceeded in approximately 20% ( $1 - 0.8$  on the 1st axis) of the cases. Of course, such observations are to be regarded as estimates and as such influenced by random variation.

Besides the information gained from the rank histograms the QQ-plots adds the insight that the HIRLAM ensembles seems to represent the upper quantiles (approximately 80% and above) fairly well, i.e. the curves are close to the line of identity. None of the ensemble forecast systems represent the lower quantiles well, except maybe for the 72 hour HIRLAM forecast.

Table 2 summarizes the results when comparing the 80% ensemble quantile to the actual observations and likewise with the 80% quantile based on climatology (as defined on page 10). It is seen that the actual number of cases is rather small. Furthermore, it is seen that the relative frequencies of observations above the 80% ensemble quantile is consequently lower than 20%, i.e. the 80% ensemble quantile seems to be too high.

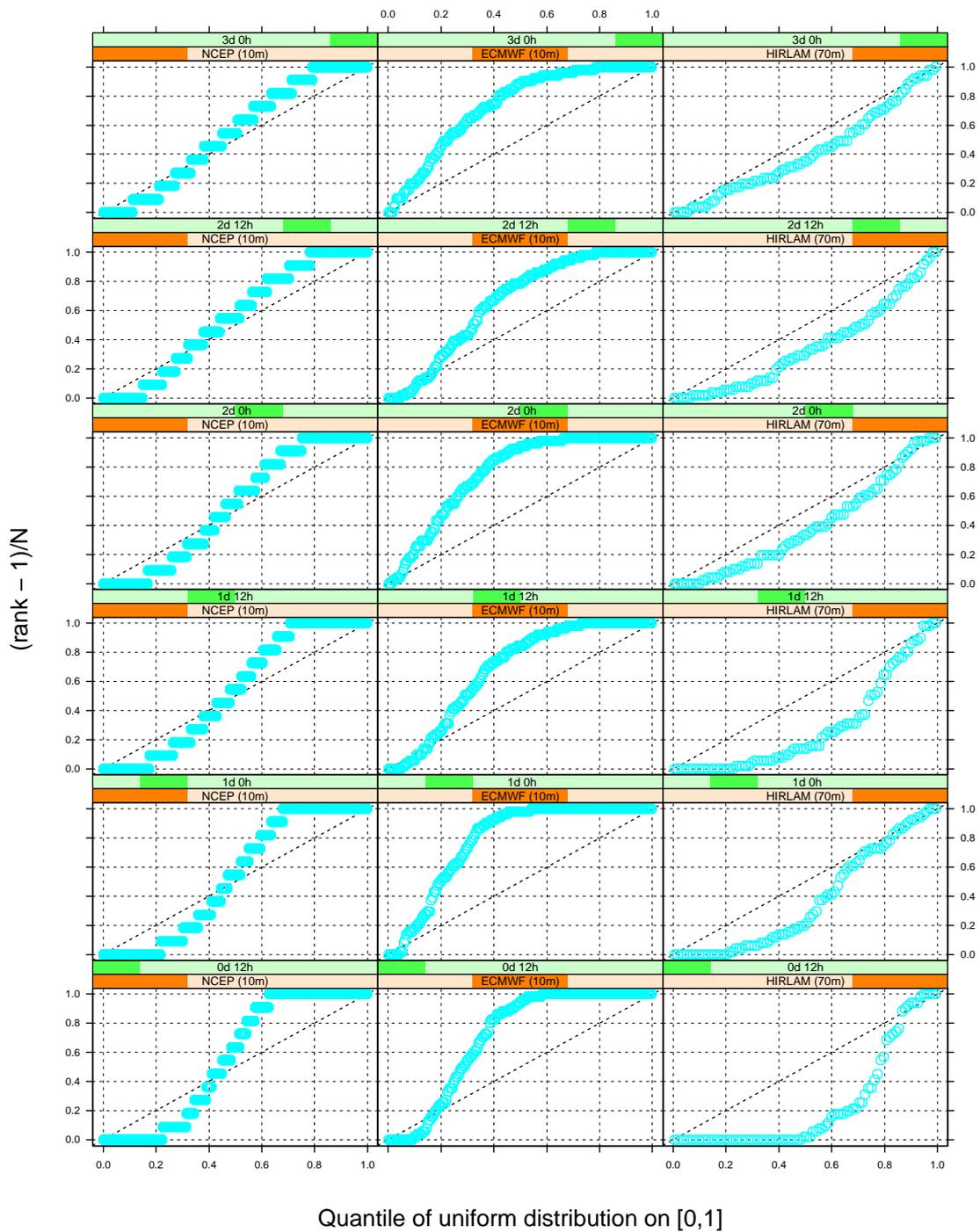


Figure 5: QQ-plots of ranks when comparing measurements with ensemble forecasts

Due to the inherent correlations in the data not attempts has been made to formally test hypotheses.

	$< 80\%E$	$\geq 80\%E$	$\% \geq E$	$< 80\%C$	$\geq 80\%C$	$\% \geq C$
3d 0h 0m 0s OMS	73	13	15.1	68	18	20.9
2d 12h 0m 0s OMS	77	9	10.5	74	12	14.0
2d 0h 0m 0s OMS	73	13	15.1	69	17	19.8
1d 12h 0m 0s OMS	75	11	12.8	74	12	14.0
1d 0h 0m 0s OMS	71	15	17.4	67	19	22.1
0d 12h 0m 0s OMS	74	12	14.0	74	12	14.0

Table 2: HIRLAM (70m): The number of cases where the observation is below ( $<$ ) or above ( $\geq$ ) the ensemble 80% quantile ( $80\%E$ ) or the climatology 80% quantile ( $80\%C$ ). Also the relative frequencies of observations above the 80% quantile are displayed ( $\% \geq E$  /  $\% \geq C$ ).

### 6.1.2 MOS adjusted ensemble forecasts

The results described above might well be influenced by the fact the forecasts are to be interpreted as spatial averages and that these are compared with point measurements. To account for systematic difference we might consider using some kind of MOS (Model Output Statistic) adjustment.

Here it was decided to compare the measurement with the unperturbed forecast and seek a linear transformation of the forecast. This transformation is then applied to all ensemble members and compared to the measurements. In this case the measurements are used both for obtaining the transformation and for calculation of ranks. Hence, the results will be somewhat over-optimistic. Separate transformations are found for each horizon and sector ( $0^\circ-45^\circ$ ,  $45^\circ-90^\circ$ ,  $\dots$ ,  $315^\circ-360^\circ$ ). The reason for estimating a transformation separately for each horizon is that the properties of the ensemble forecast systems varies somewhat over horizons, cf. Section 5. The reason for estimating a transformation separately for each sector is that for the particular location the roughness varies quite significantly, e.g. some sectors contain mostly sea surface.

Let  $y$  denote the measurement and let  $x_{fc}$  denote the forecast. A normal MOS procedure would then use least squares to estimate  $b$  in the model

$$y = bx_{fc} + e_y, \quad (1)$$

where  $e_y$  is the error or noise. The entire ensemble is then transformed using the estimate  $\hat{b}$ , whereby an MOS adjusted ensemble is obtained.

There are some problems with this approach:

- The ranks of the point measurements as compared to the adjusted ensemble will have a tendency of non-uniformity (too few low ranks) due to the fact that  $e_y$  is neglected.
- The forecast can be considered a noisy estimate of the spatial average  $x$ . If the model  $y = b x + e_y$  is the model in which we want to estimate  $b$  then using (1) will result in an estimate of  $b$  biased towards zero. This will also produce a tendency of too few low ranks.

To account for the uncertainty of the forecast we might want to add an intercept  $a$  to the linear model (1), i.e.:

$$y = a + b x_{fc} + e_y . \quad (2)$$

An other solution is obtained by considering the point measurement to be noise free, i.e. only systematic deviations from the spatial average is present in the data. This leads to a model with the forecast as the response:

$$x_{fc} = \alpha + \beta y + e_x , \quad (3)$$

where  $e_x$  is the forecast error and  $\alpha$  and  $\beta$  are coefficients which are estimated using least squares. Estimates of intercept  $a$  and slope  $b$  by which to transform the ensemble forecasts are then obtained as:

$$\hat{a} = -\hat{\alpha}/\hat{\beta} \quad \text{and} \quad \hat{b} = 1/\hat{\beta} . \quad (4)$$

However, in reality there may both be uncertainty on the forecast and random deviation between the spatial average and the point measurement. Technically, model (2) should then be treated as an *errors-in-variables* problem, but this requires knowledge about the ratio between the standard deviation of the forecast error (when comparing the forecast to the spatial average) and the standard deviation of the point measurement (the random deviation from  $E(y|x) = a + b x$ ). However, this ratio is not known and therefore we just use orthogonal regression, which corresponds to assuming the ratio to be one.

QQ-plots of the ranks obtained when using these different possibilities are shown in Figures 6, 7, and 8. The labels used in the plots are:

**“Intercept Forecast”**: Linear regression including intercept and with the forecast as the response, i.e.  $x_{fc} = a + b y + e_x$  fitted by least squares.

**“Intercept Orthogonal”**: Linear relation found by minimizing the sums of squares of the (orthogonal) distances between the line (including intercept) and the data points.

**“Intercept Measurement”**: Linear regression including intercept and with the measurement as the response, i.e.  $y = a + b x_{fc} + e_y$  fitted by least squares.

**“NO intercept Forecast”**: Linear regression excluding intercept and with the forecast as the response, i.e.  $x_{fc} = b y + e_x$  fitted by least square.

“NO intercept Measurement”: Linear regression excluding intercept and with the measurement as the response, i.e.  $y = bx_{fc} + e_y$  fitted by least squares.

From the QQ-plots it is concluded that, with respect to *reliability*:

- None of the QQ-plots obtained for NCEP ensembles are close to the line of identity.
- For ECMWF and HIRLAM the methods labeled “Intercept Forecast” and “Intercept Orthogonal” perform well whereas the remaining methods results in QQ-plots deviating from the line of identity. There is a tendency of “Intercept Forecast” performing better than “Intercept Orthogonal”.
- For the low horizons the ensemble spread is too small (S-shaped curves).

Figure 9 shows the QQ-plots for ECMWF forecasts for horizons longer than the ones covered in the plots mentioned above. The corresponding rank histograms are displayed in Figure 10. It is seen that the conclusions listed above are valid for the 4 day horizon also. However, for horizons of 5 days or longer especially “Intercept Forecast” the adjusted ensemble forecasts are too wide, i.e. there are too many ranks near the center. Since the standard deviation of the point measurement can not depend on the forecast horizon the results indicate that the underlying spread of the ECMWF ensemble forecast are too wide for horizons 5 days and above.

Overall for all horizons up to 7 days the method “Intercept Orthogonal” perform well, but for horizons up to approximately 3 days the method “Intercept Forecast” is better in terms of reliability. The possibility exists that averaging of point measurements can reduce the random variation and thereby improving the estimates. .

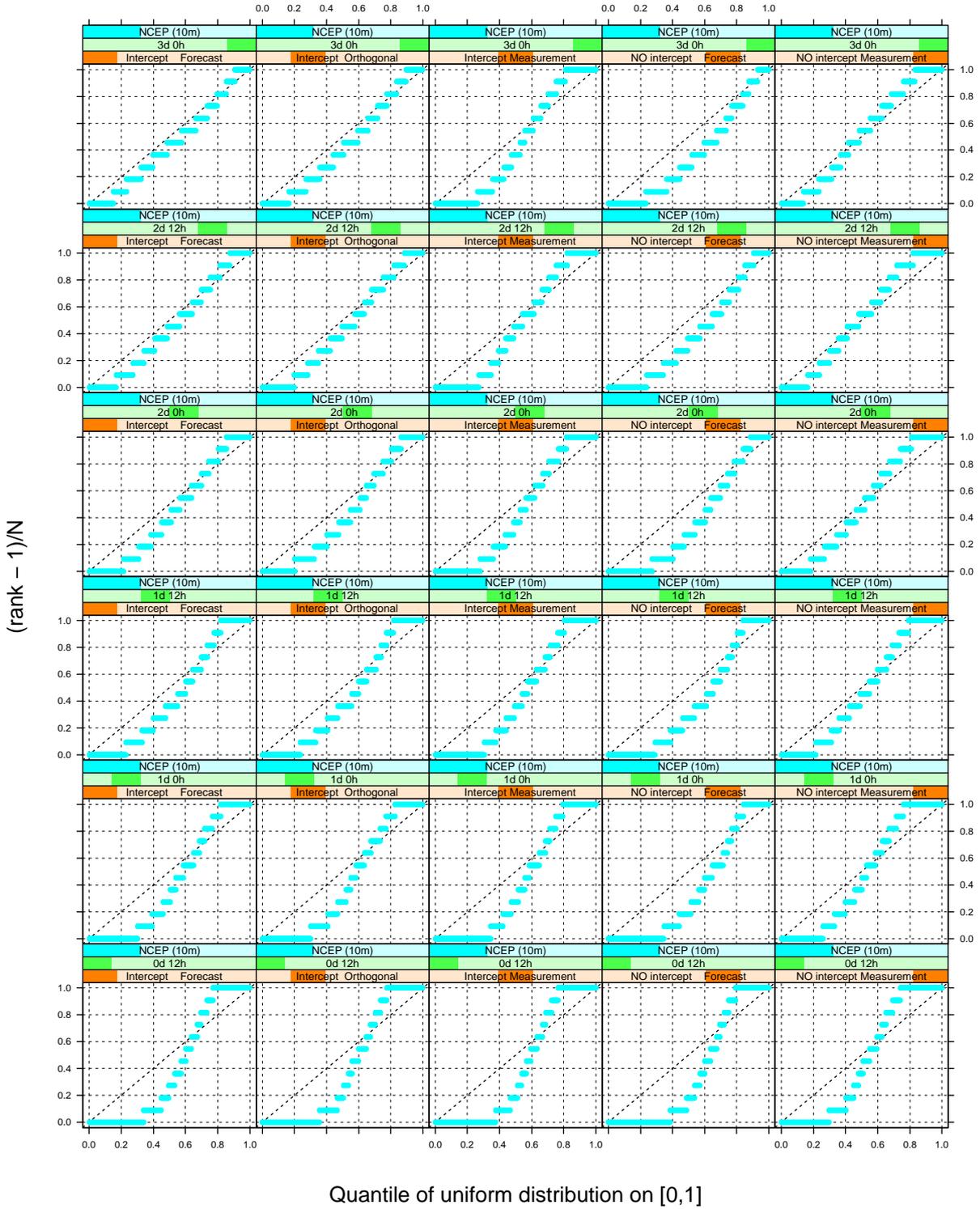


Figure 6: QQ-plots of ranks when comparing measurements with MOS-adjusted NCEP ensemble forecasts.

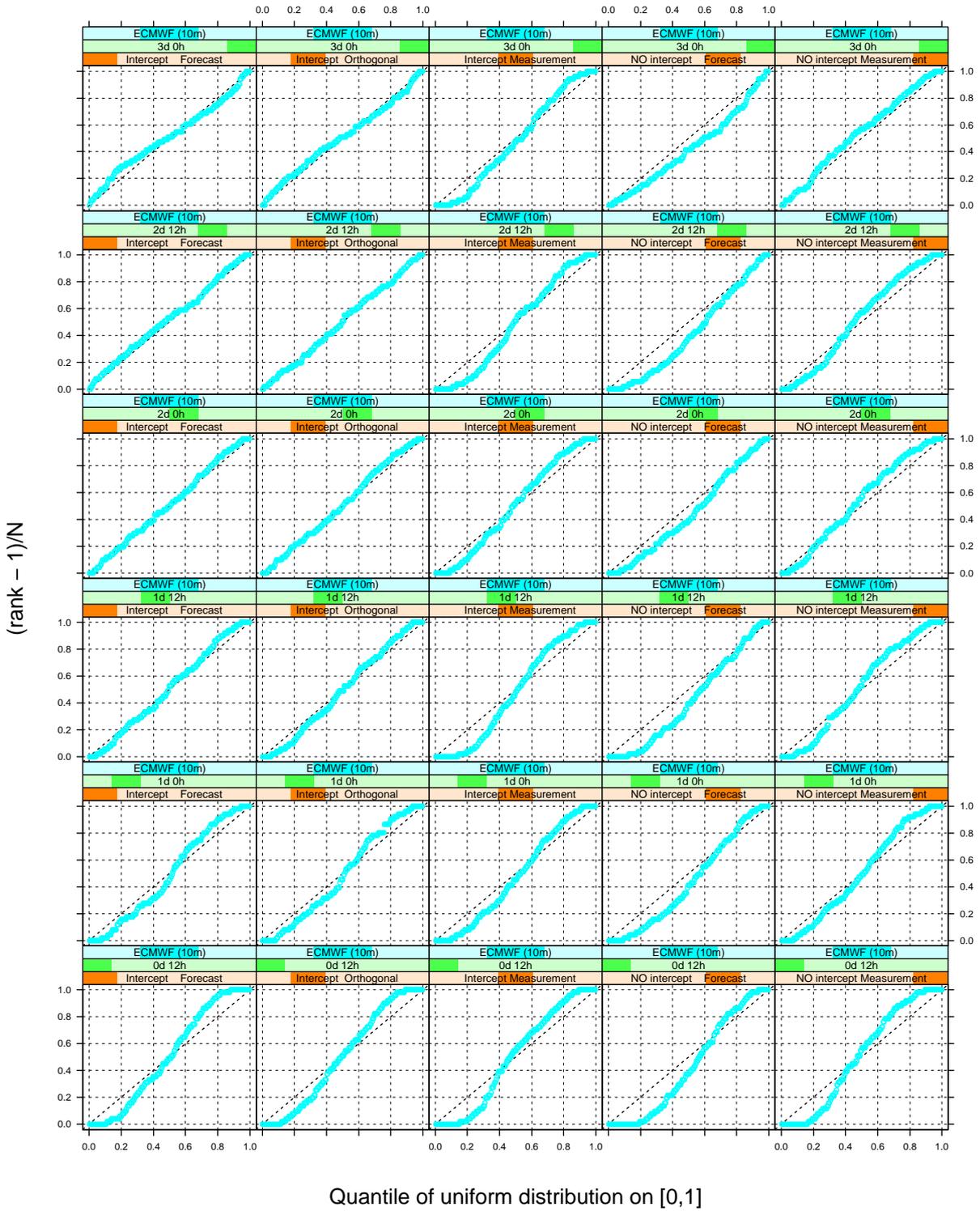


Figure 7: QQ-plots of ranks when comparing measurements with MOS-adjusted ECMWF ensemble forecasts.

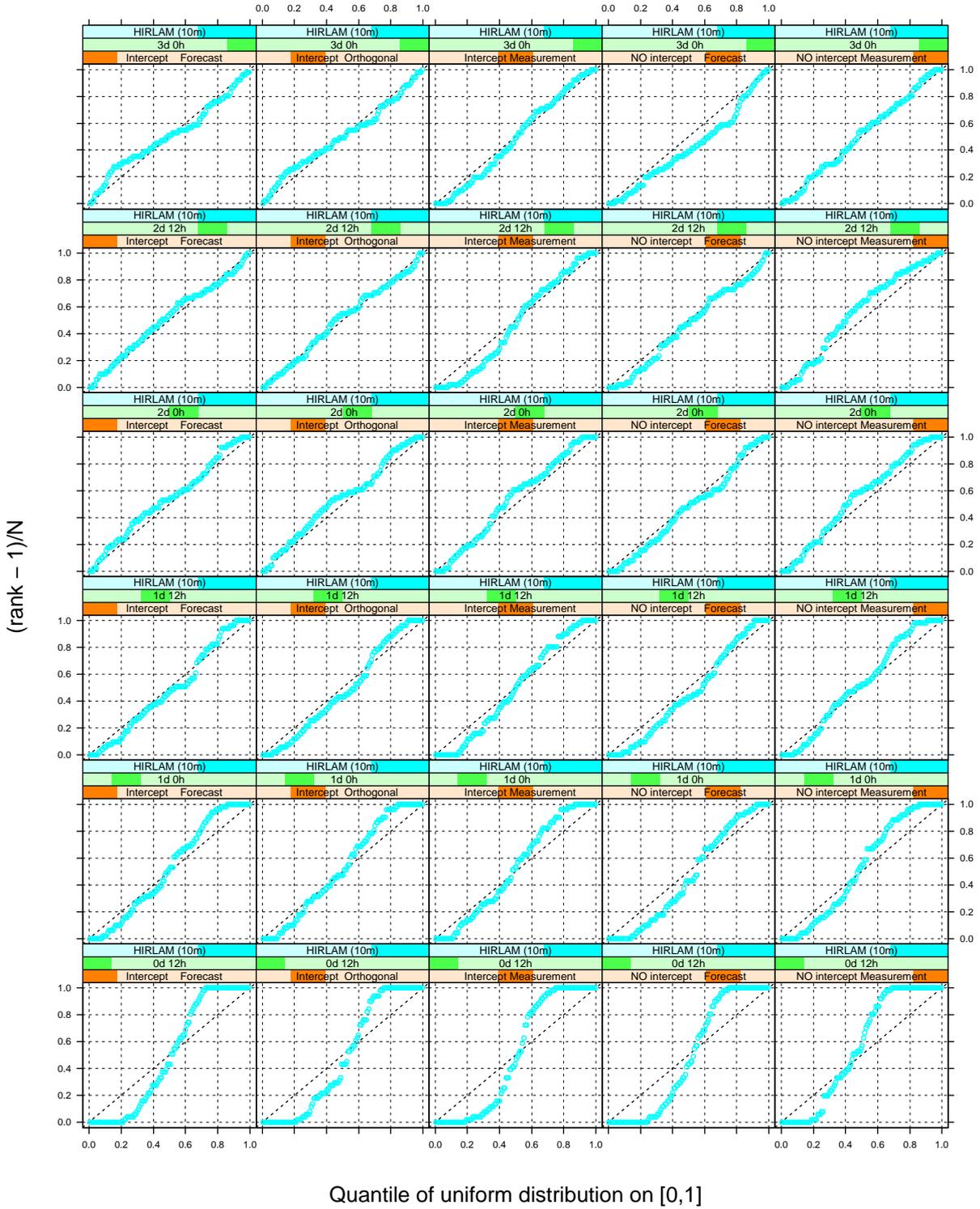


Figure 8: QQ-plots of ranks when comparing measurements with MOS-adjusted HIRLAM ensemble forecasts.

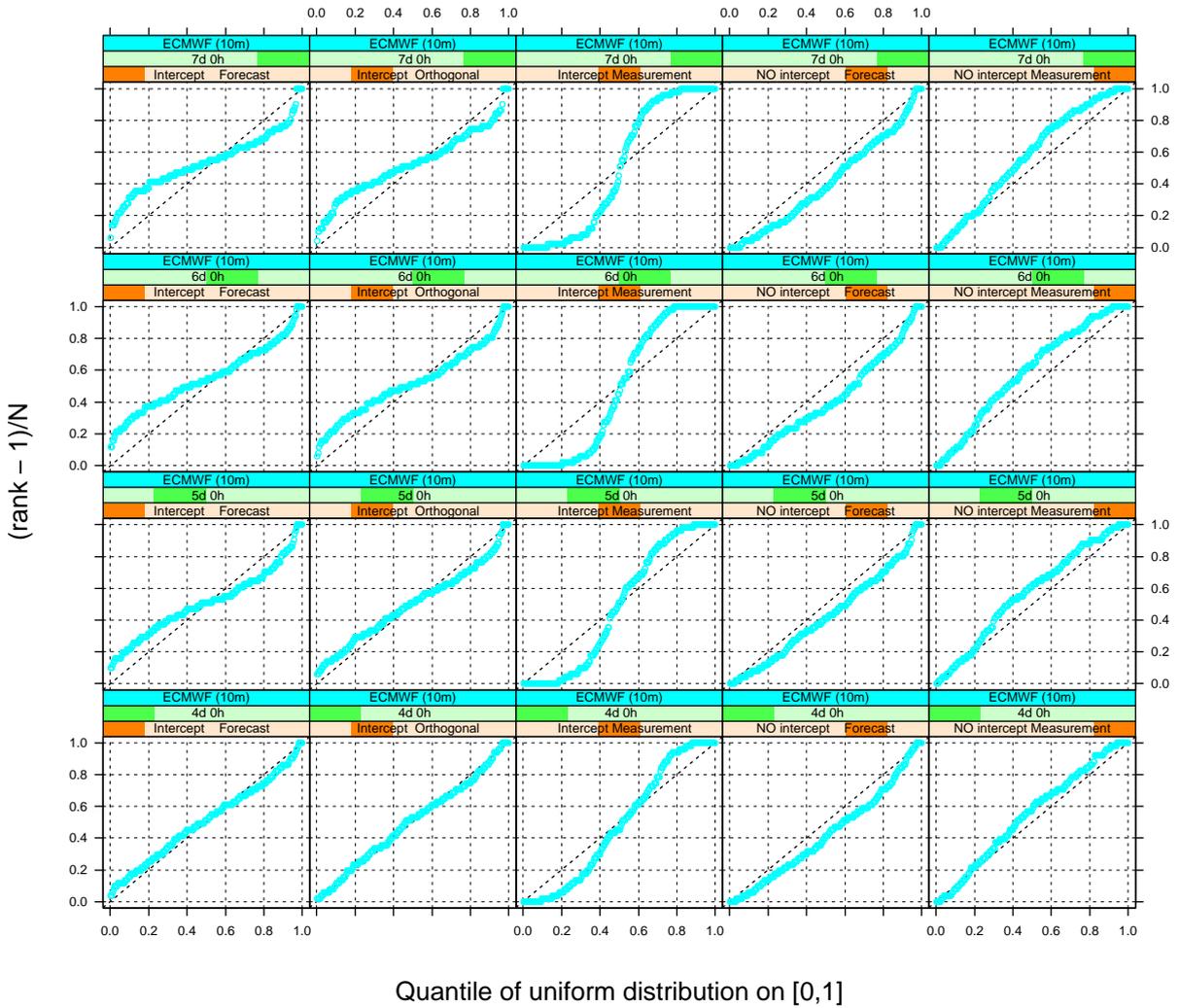


Figure 9: QQ-plots of ranks when comparing measurements with MOS-adjusted ECMWF ensemble forecasts for horizons 4, 5, 6, and 7 days.

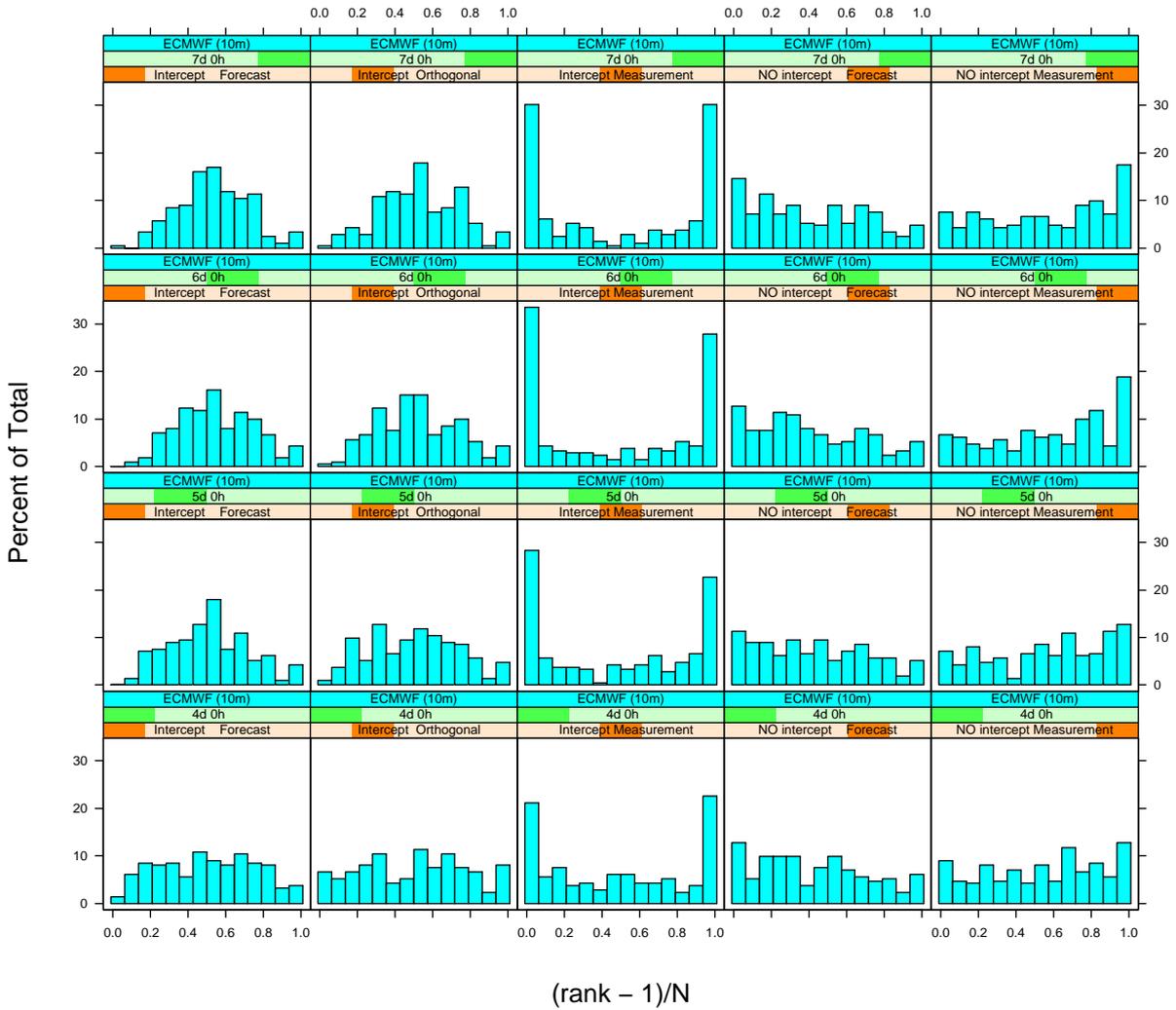


Figure 10: Rank histograms when comparing measurements with MOS-adjusted ECMWF ensemble forecasts for horizons 4, 5, 6, and 7 days.

## 6.2 Resolution

As noted on page 13, the unadjusted HIRLAM ensemble forecast at 70m seems to represent the upper quantiles (80% and above) fairly well. The quantiles and the observations are depicted in Figure 11. The climatological (page 3) 80% quantile is 9.5 m/s, whereas the ensemble quantile ranges from 1.7 m/s to 25.7 m/s. This indicates good resolution of the 80% HIRLAM ensemble quantile, see also Figure 12 where histograms of the quantiles are displayed.

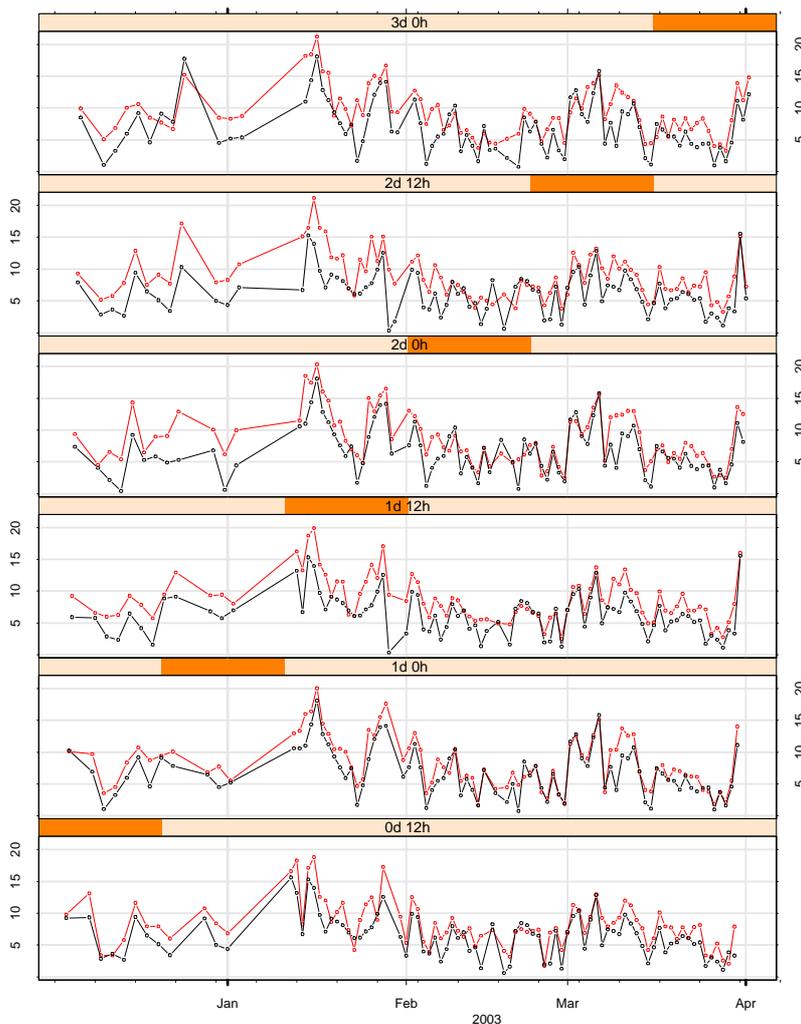


Figure 11: Observations the wind speed in 76m (black) and 80% quantiles of HIRLAM (70m) ensembles (red) for the horizons indicated on the plots. Only time points where all members are present are included in the plots. In case of mismatch between time points for observations and time points for forecasts linear interpolation is used to obtain appropriate observations.

The resolution of the MOS adjusted ECMWF ensembles is considered in the following. Also, the HIRLAM ensembles are briefly discussed. In Section 6.1 it was concluded that

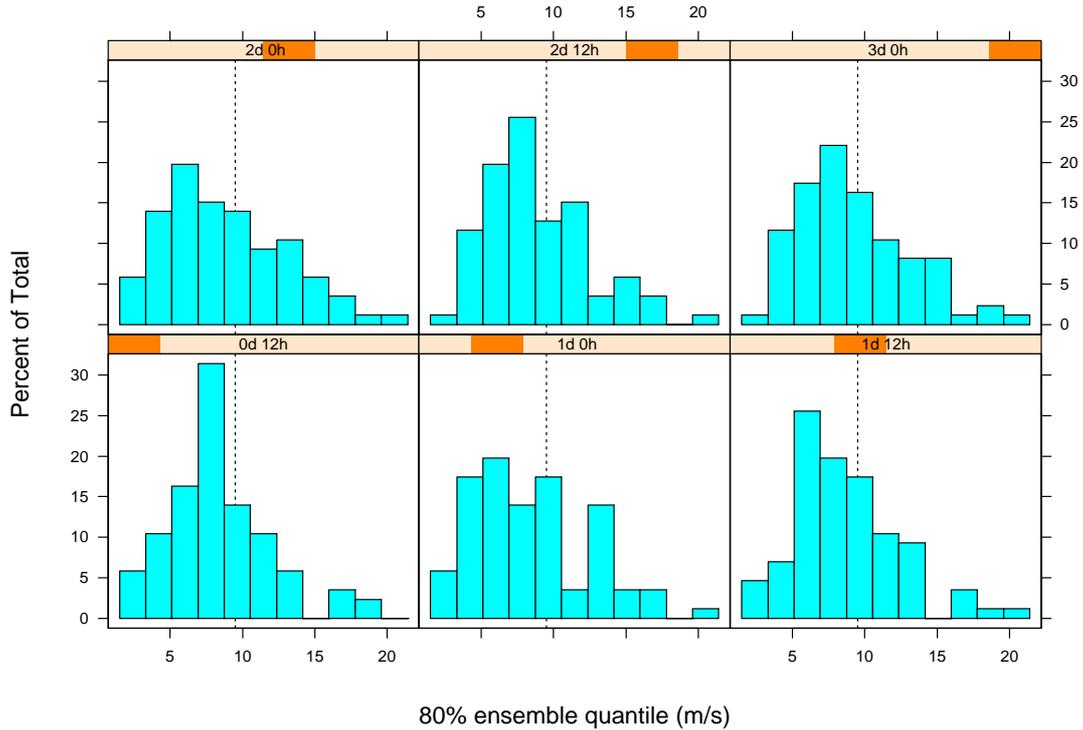


Figure 12: Histograms of the 80% HIRLAM ensemble quantile for the 70m wind speed. The 80% climatology quantile is indicated as a dotted line and the forecast horizon is displayed on top of each plot.

using the unperturbed forecast as the dependent variable and the point measurement as the independent variable and including an intercept in the model (i.e. **Intercept Forecast** on page 16) yielded good results w.r.t. reliability. More, specifically:

1. For horizons of 24 hours or lower the MOS adjusted ECMWF ensembles are too narrow.
2. For horizons between 36 and 60 hours the MOS adjusted ECMWF ensembles are reliable.
3. For horizons of 3 days or longer the MOS adjusted ECMWF ensembles are too wide.

Note that all coefficients (intercept and slope) used in the correction are found using the same data as the data for which ranks etc. are calculated.

Since the ensembles are too narrow it does not make sense to consider the resolution of item 1. For item 3 the resolution may still be better than for climatology and hence it may be beneficial to consider these as reliable. The plots in Appendix A.1 shows the MOS adjusted ensembles considered above. And in Appendix A.2 corresponding plots for HIRLAM are included.

It is seen that negative forecasts occur, which indicates that the linear approximation is inappropriate at low wind speeds. However, as shown in Section 6.1 fixing the intercept at zero does not yield reliable ensemble forecasts. Since only non-negative measurements of wind speed are possible the reliability plots, i.e. rank histograms and QQ-plots, will not be affected by setting negative forecasts to zero. This approach will be used below.

Also, the plots show that for some horizons numerical instability seems to occur. Closer investigations reveal that these instabilities occur for certain sectors where the sector  $180^\circ - 225^\circ$  is the one most frequently occurring. The reason for the instabilities is that  $\beta$  in (3) (page 16) is estimated as being close to zero whereby the absolute values of (4), which are used for transforming the ensembles, increase dramatically. These results indicate that in practice more advanced methods of MOS transformation is required. One possibility would be to estimate the coefficients in (3) as smooth functions of the wind direction and the forecast horizon. Furthermore, all observations of the wind speed should be used. One way to do this is to define the forecasts as continuous curves by use of e.g. linear interpolation.

Table 3 shows the number of cases where the ensemble IQR (Inter Quantile Range)<sup>4</sup> is below the IQR of climatology. For the horizons mentioned in item 2 above the MOS adjusted ensemble forecasts are reliable. It is seen that although the 60 hour (2d 12h) forecast has numerical instable values the ensemble IQR is below the climatology IQR in nearly 60% of the cases. For the horizons mentioned under item 3 some of the lower horizons have near 50% cases where the ensemble IQR is lower than the climatology IQR.

For horizons 36 and 48 hours Figure 13 shows histograms of the IQR of the MOS adjusted ECMWF ensemble compared to the IQR of the climatological distribution. The figure also contains plots based on the MOS adjusted HIRLAM ensembles, but these are available for a shorter period. Considering this the adjusted HIRLAM ensembles do not seem to differ much in resolution as compared to the ECMWF ensembles.

For ECMWF it is seen that the majority of the ensemble IQR values are well below the climatology value. This indicates a clear benefit from using the ensemble forecasts. For the 36 hour forecast it turns out that in approximately 50% of the cases the IQR is below 2 m/s. A IQR of 2 m/s indicates that in 50% of these cases the forecast is within 1 m/s.

As mentioned in Section 6.1 the MOS adjustment of the ECMWF forecasts obtained using orthogonal regression perform also fairly well in terms of reliability. For MOS performed using orthogonal regression Table 4 show the number of cases where the ensemble IQR is below the climatology IQR. Comparing with Table 3 it is clearly seen that the resolution of the MOS-adjusted ensemble forecasts obtained using orthogonal regression is higher than the resolution obtained using the unperturbed forecast as the dependent variable. However, as discussed above the method primarily discussed in this section may possibly be improved. This may be preferable due to it's simplicity.

---

<sup>4</sup>The inter quantile range is the difference between the 75% and the 25% quantiles and hence is a measure of how wide the distribution is.

	Horizon	Total	No. below	Relative (%)
Item 1	0d 0h	212	210	99.1
	0d 6h	212	211	99.5
	0d 12h	212	208	98.1
	0d 18h	212	207	97.6
	1d 0h	212	209	98.6
	1d 6h	212	195	92.0
Item 2	1d 12h	212	188	88.7
	1d 18h	212	191	90.1
	2d 0h	212	181	85.4
	2d 6h	212	176	83.0
	2d 12h	212	122	57.5
Item 3	2d 18h	212	145	68.4
	3d 0h	212	104	49.1
	3d 6h	212	91	42.9
	3d 12h	212	83	39.2
	3d 18h	212	56	26.4
	4d 0h	212	51	24.1
	4d 6h	212	35	16.5
	4d 12h	212	7	3.3
	4d 18h	212	4	1.9
	5d 0h	212	1	0.5
	5d 6h	212	2	0.9
	5d 12h	212	0	0.0
	5d 18h	212	0	0.0
	6d 0h	212	0	0.0
	6d 6h	212	0	0.0
	6d 12h	212	0	0.0
	6d 18h	212	0	0.0
7d 0h	212	0	0.0	

Table 3: MOS adjustment using `Intercept Forecast` as defined on page 16: Number of cases where the ensemble IQR is below the climatology IQR. The item refers to the three groups of horizons mentioned on page 24. Note however that the groups were set up by considering horizons only horizons 12, 24, 36, ... hours.

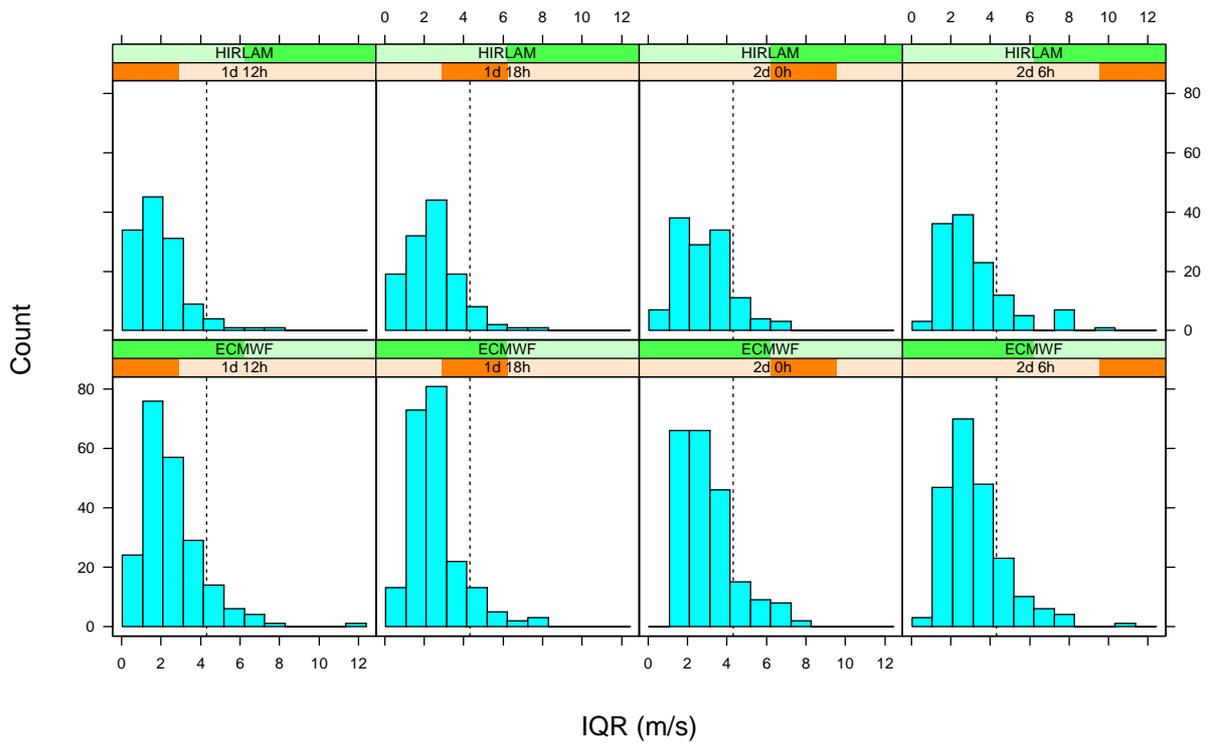


Figure 13: Histograms (counts) of ensemble IQR (MOS adjusted ECMWF/HIRLAM with the unperturbed forecast as the dependent variable). The IQR corresponding to climatology, as defined on page 3, is shown as a vertical line on the plot.

	Horizon	Total	No. below	Relative (%)
Item 1	0d 0h	212	210	99.1
	0d 6h	212	211	99.5
	0d 12h	212	207	97.6
	0d 18h	212	209	98.6
	1d 0h	212	211	99.5
	1d 6h	212	201	94.8
Item 2	1d 12h	212	199	93.9
	1d 18h	212	198	93.4
	2d 0h	212	190	89.6
	2d 6h	212	193	91.0
	2d 12h	212	175	82.5
Item 3	2d 18h	212	176	83.0
	3d 0h	212	148	69.8
	3d 6h	212	154	72.6
	3d 12h	212	147	69.3
	3d 18h	212	121	57.1
	4d 0h	212	127	59.9
	4d 6h	212	125	59.0
	4d 12h	212	58	27.4
	4d 18h	212	52	24.5
	5d 0h	212	39	18.4
	5d 6h	212	74	34.9
	5d 12h	212	31	14.6
	5d 18h	212	1	0.5
	6d 0h	212	2	0.9
	6d 6h	212	0	0.0
	6d 12h	212	7	3.3
6d 18h	212	0	0.0	
7d 0h	212	0	0.0	

Table 4: MOS adjustment using *orthogonal regression*: Number of cases where the ensemble IQR is below the climatology IQR. The item refers to the three groups of horizons mentioned on page 24. Note however that the groups were set up by considering horizons only horizons 12, 24, 36, ... hours.

## 7 Conclusion and Discussion

This report compares three types of ensemble forecasts of wind speed to a wind speed measurement 76m a.g.l. from a mast at Risø National Laboratory, Denmark. The ensemble forecasts considered are

- NCEP (National Centers for Environmental Prediction) ensemble forecasts from the National Weather Service of NOAA (National Oceanic and Atmospheric Administration) in U.S. This set of ensembles consists of one unperturbed forecast and 5 pairs of forecasts for which the initial conditions are perturbed in the positive and negative direction of bred vectors. Horizontal resolution: 1°.
- ECMWF (European Centre for Medium-Range Weather Forecasts) ensemble forecasts from the Ensemble Prediction System<sup>5</sup>. This set of ensembles consists of one unperturbed forecast and 25 pairs of forecasts for which the initial conditions are perturbed in the positive and negative direction of singular vectors. Furthermore, for each model run attempts are made to account for sub-grid processes by use of stochastic physics [1]. As a result two runs with the same set of initial conditions will not result in exactly the same output. Horizontal resolution: 75km. The unperturbed forecast is not influenced by stochastic physics.
- HIRLAM ensembles from DMI (Danish Meteorological Institute). These are experimental ensembles based on the ECMWF ensembles. Also a few model permutations are included, but these are not used in this report. Horizontal resolution: 20km.

For NCEP and ECMWF the 10m a.g.l. wind speed and direction are used and for HIRLAM the 70m a.g.l. wind speed and direction are used. For the analysis of the MOS adjusted forecasts (see below) the 10m a.g.l. HIRLAM wind speed is used since more data are available for this level. Bilinear interpolation is used to obtain forecasts valid for the location of the mast. The period for which measurements and forecasts are available are listed in Section 4.

In Section 5 the ensemble forecasts are investigated without comparing them to the measurements. It is shown that the NCEP analysis (0 hour forecast) averaged over all ensemble members and time points is approximately 2 m/s lower than the remaining mean values. For ECMWF the average and quantiles over all ensemble members and time points varies with a period of 24 hours, but since the calculations are only initiated once a day, this may just be a consequence of the diurnal variation forecasted by the model. In fact if the NCEP forecasts are split by initialization time the same kind of behavior is observed, but to a lesser extent. For the HIRLAM ensembles the average drops by more than 2 m/s within the first three hours, i.e. comparing the analysis and the 3 hour forecast. This time interval corresponds to 12:00 to 15:00 UTC and hence the drop can not be attributed

---

<sup>5</sup>[http://www.ecmwf.int/products/forecasts/guide/The\\_Ensemble\\_Prediction\\_System\\_EPS.html](http://www.ecmwf.int/products/forecasts/guide/The_Ensemble_Prediction_System_EPS.html)

to diurnal variation (Denmark is one hour ahead of UTC). Consequently, there is some indication that the bias of the forecasts depends on the horizon.

The ensemble forecasts are compared w.r.t. *reliability* (average correctness of the forecasted distributions) and *resolution* (sharpness of the forecasted distributions), cf. [13] or the beginning of Section 6. When interpreting ensemble forecasts in a probabilistic sense reliability is a requirement and resolution is an indicator of performance. In this report reliability is addressed by use of rank histograms and QQ-plots. Rank histograms, also called Talagrand diagrams, are useful in order to obtain an overview, but QQ-plots (quantile-quantile-plots) are generally preferable in that they contain additional information. As an example QQ-plots may indicate that overall a particular ensemble forecast system is not reliable, but the upper quantiles are reliable. Only the overall conclusion can be reached by use of rank histograms. Resolution is addressed by comparing the IQR (Inter Quantile Range), the difference between the 75% and the 25% quantiles, with the IQR obtained using climatological information. Here the climatological information is obtained by simply using the available measurements with time stamps before any of the ensemble forecasts, i.e. 9 months of measurements. See also Section 4.

When comparing the forecasts to the point measurement differences in temporal and spatial resolution should be taken into account. Both measurements and forecasts are 10 minute averages (over the interval up to the time stamp), however the spatial resolutions of the measurements and forecasts are very different. This may have the consequence that (i) the point measurement may differ systematically from the (unobservable) spatial averages although the forecasts may agree well with the spatial averages, and (ii) the point measurement may differ non-systematically or stochastically from the spatial averages. Using MOS (Model Output Statistics) it should be possible to correct for (i). However, (ii) will influence the reported reliability of the ensemble forecasts, even when these are corrected by MOS.

As expected without MOS correction none of the ensemble forecast systems are reliable w.r.t. the measurement. Nevertheless, the upper quantiles (from 80%) of the HIRLAM ensembles seems to be reliable, but the actual number of times in which the 80% quantile is exceeded is very small, cf. Table 2 page 15. The 80% ensemble quantile range from 1.7 to 25.7 m/s, indicating good resolution.

If (ii) above can be neglected then some kind of MOS correction should be able to correct for systematic differences. If further, the ensemble forecast system produce forecasts which are reliable w.r.t. the spatial averages then the MOS adjusted ensembles should also be reliable and this can be investigated using rank histograms or QQ-plots. In the report different ways of performing the MOS adjustment are considered. Due to the limited data the coefficients used for correction are determined on the same set of data as the one for which the correction is performed. To account for a bias which may depend on the horizon and on the sector ( $0^{\circ}$ – $45^{\circ}$ ,  $45^{\circ}$ – $90^{\circ}$ ,  $\dots$ ,  $315^{\circ}$ – $360^{\circ}$ ) the MOS correction is determined separately for each combination of sector and horizon.

Usually, a linear regression through the origin with the unperturbed forecast as the independent variable and the measurement as the dependent variable is used for MOS adjustment. However, due to the uncertainty of the forecast the estimate will be biased, but the point predictions of the wind speed at the mast will be good in e.g. mean square sense [5, 7, 8]. The point predictions can probably be further improved by not forcing the linear regression through the origin. At least two reasons exist why this approach to MOS adjustment must be expected to result in *unreliable* ensemble forecasts; (a) even if the estimates are unbiased the error-term of the linear regression is not included in the MOS adjustment, and (b) the uncertainty of the unperturbed forecast results in a downward bias on the slope and an upward bias on the intercept. The consequence of both (a) and (b) are that the spread of the MOS adjusted ensemble forecasts will be too small. This is demonstrated in Section 6.1, see e.g. Figure 7 on page 19. Actually, w.r.t. reliability it seems to be best to force the linear regression through the origin.

If there is no non-systematic difference between the point measurement (10 minute average) and the spatial average, then only the unperturbed forecast is associated with uncertainty. Consequently, using the unperturbed forecast as the dependent variable and the measurement as the independent variable will yield unbiased estimates.

Horizons 12, 24, 36, ... of the MOS adjusted ensembles are checked w.r.t. reliability and resolution. The spread of the NCEP ensembles adjusted in this way is still too small. For ECMWF and HIRLAM the adjusted ensembles seem to be reliable for horizons 36–60 hours, whereas the spread is too small for shorter horizons and too large for longer horizons.

For horizons 36–48 hours the resolution is good in that the ensemble IQR, in most cases, is smaller than the IQR corresponding to climatology. For 60 hours numerical instabilities result in unrealistic MOS adjusted ensemble forecasts. For horizons of 60 hours or longer QQ-plots indicate that the spread of the MOS adjusted ensemble forecasts is too large, also for some of these horizons numerical instabilities are observed.

The cause of the numerical instabilities seems to be a rather large uncertainty of the estimates associated with the MOS correction, cf. Section 6.2, caused by the low number of observations in each sector and the fact that the ensemble calculations are only initiated at 12:00 UTC. When the unperturbed forecast is rather uncertain the estimate of slope may occasionally be close to zero, possibly just due to random variation<sup>6</sup> and this results in a very large slope for in the linear transformation which is used to adjust the ensemble forecasts, cf. (4) on page 16.

To avoid the numerical instabilities the following issues should be considered:

- More effective use of the observations. One way to accomplish this is to obtain

---

<sup>6</sup>By swapping  $x$  and  $y$  in the regression we have replaced bias by variance. This is required since we can not accept bias in this case.

forecasts for all observations, e.g. by linear interpolation, and then let the coefficients be smooth functions of the horizon.

- More effective use of the information contained in the wind direction by modeling the coefficients as smooth functions of the wind directions.
- Derive the MOS transformation based on a short forecast horizon and use this for all horizons. However this is difficult, at least for HIRLAM (cf. Section 5).

If both uncertainty on the forecast and stochastic variations between the point measurement and the unobservable spatial average is to be taken into account the problem must be treated as an errors-in-variables problem [4]. This requires knowledge about the ratio of the two variances involved. The ratio is not known in this case, but orthogonal regression has been applied which corresponds to a ratio of one. The method results in adjusted ensemble forecasts with similar reliability as when the unperturbed forecast is used as the dependent variable. Furthermore, the numerical instabilities of the method is somewhat smaller.

In [13, p. 141] it is stated that for statistically stationary forecast and observation systems and given a large enough sample, perfect reliability can be achieved by a simple statistical calibration. The corresponding transformation can be obtained by fitting a smooth function to the data on a QQ-plot, with the observed quantile/probability  $((rank - 1)/N)$  as the explanatory variable. This transformation is then used to transform the ensemble quantiles to agree with the observed quantiles.

Note that this transformation is solely focusing on reliability; as an extreme example, if the ensemble forecasts are not related to the observations then the quantiles obtained are simply quantiles in the climatological distribution. Also, in practice the transformation can not repair any QQ-plot. As an example consider the 12h NCEP ensemble forecast for which the QQ-plot is displayed in the lower left corner of Figure 5 on page 14. It is seen that in approximately 20% / 40% of the cases the observation is below / above the lowest / highest ensemble member. There is no way to distinguish between these quantiles; the best we can hope for is values near the center of the two extreme intervals. Consequently, in the example, we should not expect to gain information about quantiles below 10% or above 80%. The same kind of problems occur when many of the observations are in the high end of the ensemble members as e.g. for ECMWF on the figure just mentioned.

Because of the above and because the transformation is not directly based on any physical aspects (e.g. spatial averages versus point measurements) we suggest that a MOS transformation based on the unperturbed forecast is applied before the transformation based on the QQ-plot is applied. In this setting it may not be a requirement that the MOS adjustment is based on unbiased estimates. Hence, one solution might be to use a standard MOS adjustment derived on basis of a low forecast horizon and then apply the transformation based on the QQ-plots, for each horizon separately. By substituting the MOS adjustment by the estimation of a power curve for wind farms, this points towards a way of producing reliable wind power ensembles.

Finally, we note that the observations mentioned in Section 5 do have consequence for how to estimate and use a power curve model for producing ensembles of wind power. Since we, for some setups, want an estimate without bias we might want to select the most precise forecast for building the power curve model. If e.g. for NCEP the analysis is selected for this purpose then using other horizons result in a large upward shift in power production.

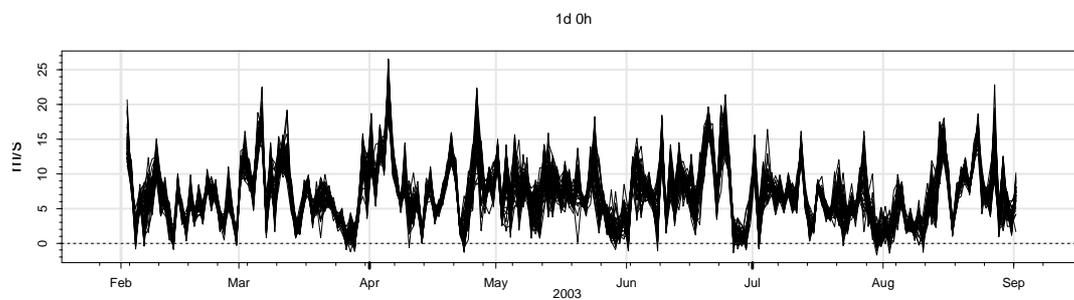
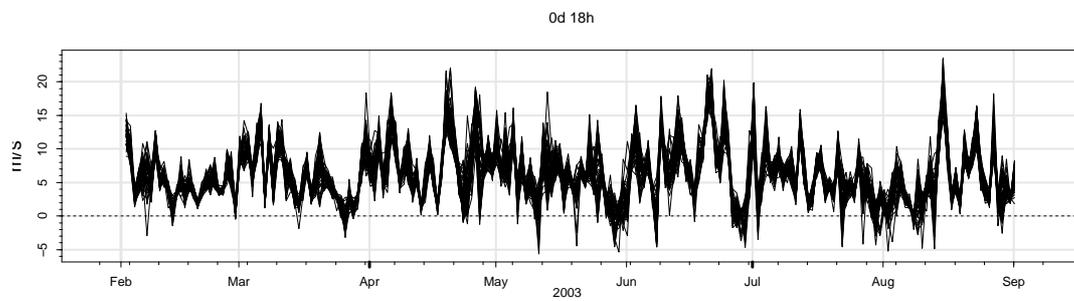
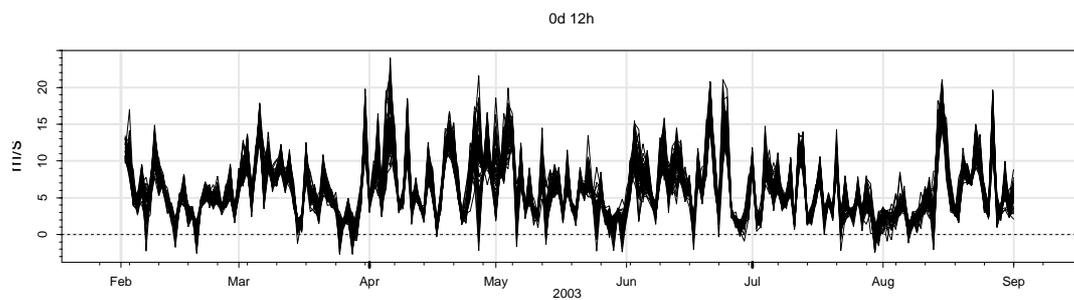
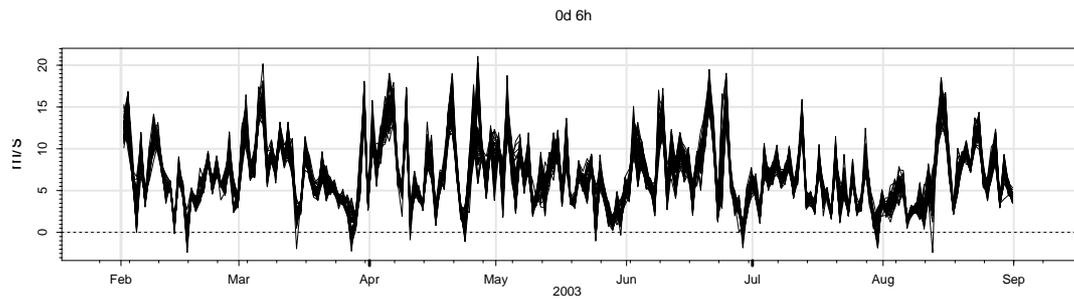
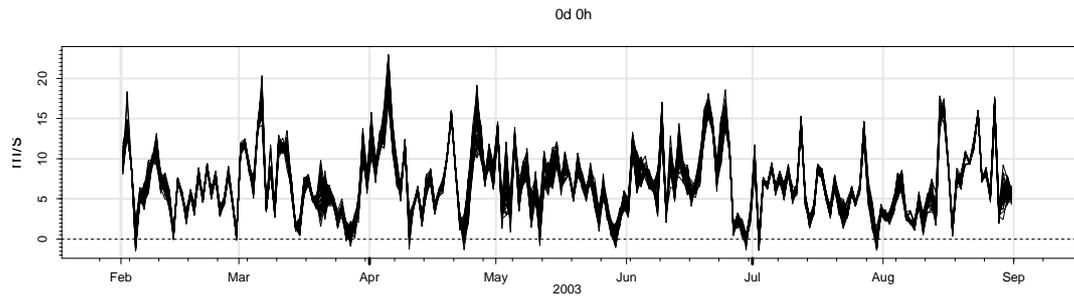
## References

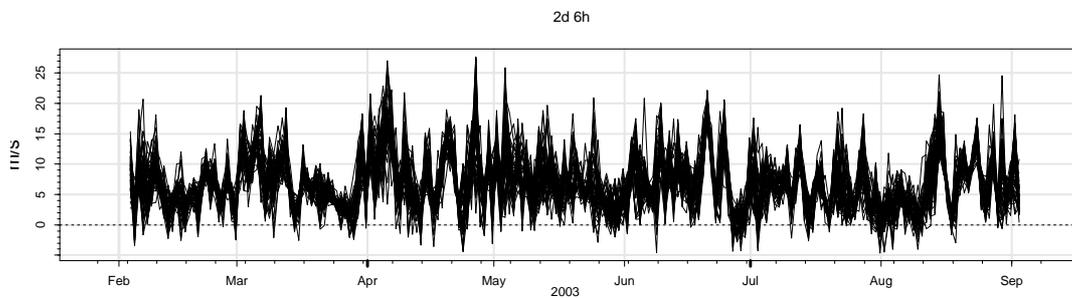
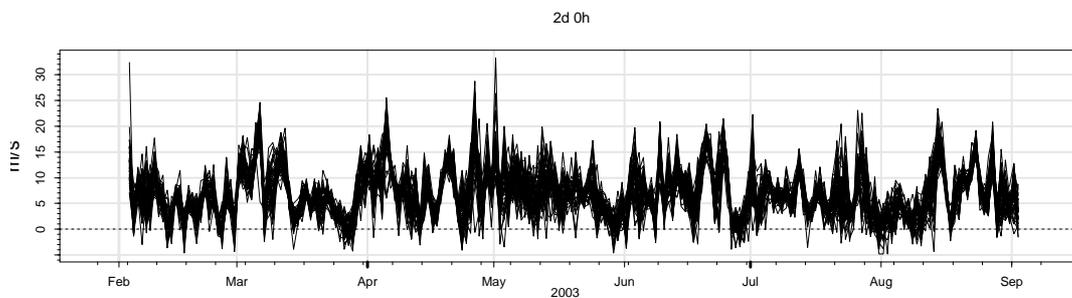
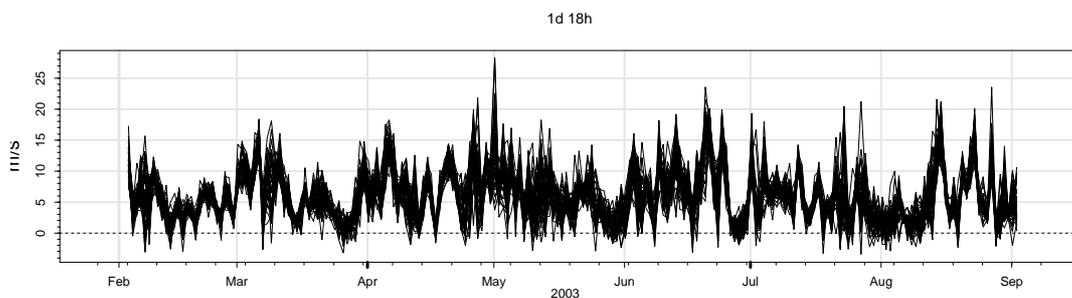
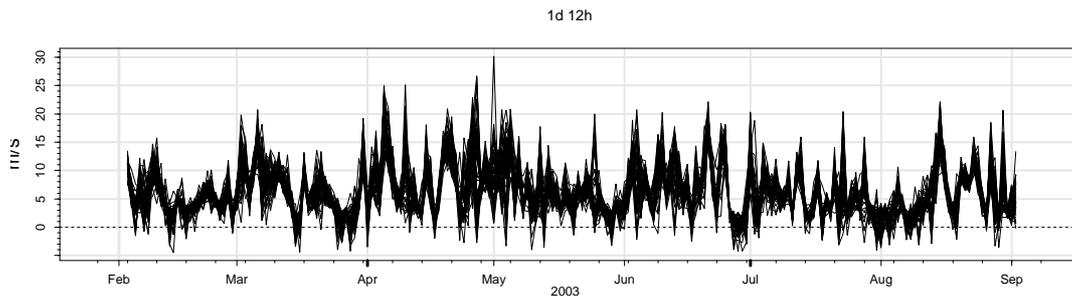
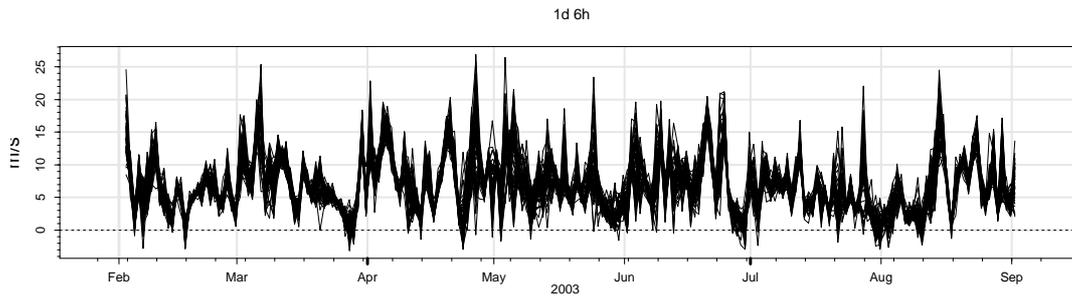
- [1] R. Buizza, M. Miller, and T.N. Palmer. Stochastic representation of model uncertainties in the ECMWF Ensemble Prediction System. *Quarterly Journal of the Royal Meteorological Society*, 125(560):2887–2908, 1999.
- [2] R. Buizza and T.N. Palmer. The singular vector structure of the atmosphere global circulation. *Atmos. Sci.*, 52:1434–1456, 1995.
- [3] J.M. Chambers, W.S. Cleveland, B. Kleiner, and P.A. Tukey. *Graphical methods for data analysis*. Wadsworth Publishing Co Inc, 1983.
- [4] Wayne A. Fuller. *Measurement error models*. John Wiley & Sons, 1987.
- [5] Bo Jonsson. Prediction with a linear regression model and errors in a regressor. *International Journal of Forecasting*, 10(4):549–555, 1994.
- [6] F. Molteni, R. Buizza, T.N. Palmer, and T. Petroliagis. The ecmwf ensemble prediction system: methodology and validation. *Quart. J. Roy. Meteor. Soc.*, 122:73–119, 1996.
- [7] Henrik Aalborg Nielsen, Torben Skov Nielsen, and Henrik Madsen. On on-line systems for short-term forecasting for energy systems. In *Proceedings of the OR 2002 conference*, pages 265–271, Klagenfurt, Austria, 2002. Springer.
- [8] Henrik Aalborg Nielsen, Torben Skov Nielsen, and Henrik Madsen. Using meteorological forecasts for short term wind power forecasting. In *Proceedings of the IEA R&D Wind Annex XI Joint Action Symposium on Wind Forecasting Techniques*, pages 49–58, Norrköping, Sweden, December 3–4 2002.
- [9] T.N. Palmer, F. Molteni, R. Mureau, R. Buizza, P. Chapelet, and J. Tribbia. Ensemble prediction. In *Proceedings of the ECMWF Seminar on Validation of models over Europe*, volume 1, Shinfield Park, Reading RG2 9AX, UK, 1993. ECMWF.
- [10] K. Sattler and H. Feddersen. Treatment of uncertainties in the prediction of heavy rainfall using different ensemble approaches with dmi-hirlam. *DMI Sci. Rep.*, **03-07**, 62pp, 2003. Available at [www.dmi.dk](http://www.dmi.dk).
- [11] Z. Toth and E. Kalnay. Ensemble forecasting at nmc: The generation of perturbations. *Bull. Amer. Meteor. Soc.*, 74:1485–1490, 1993.
- [12] Z. Toth and E. Kalnay. Ensemble forecasting and the breeding method. *Mon. Wea. Rev.*, 12:3297–3319, 1997.
- [13] Zoltan Toth, Oliver Talagrand, Guillem Candille, and Yuejian Zhu. *Forecast verification – a practitioner’s guide in atmospheric science*, chapter Probability and ensemble forecasts. Wiley, 2003.

# Appendices

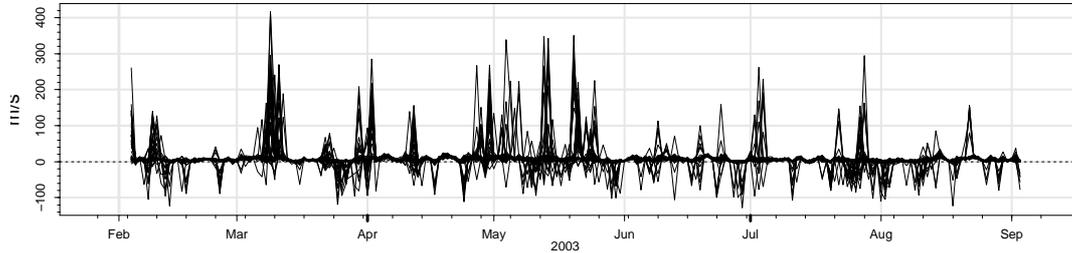
## A Plots of MOS adjusted ensemble forecasts

### A.1 ECMWF corrected using the unperturbed forecast as the dependent variable.

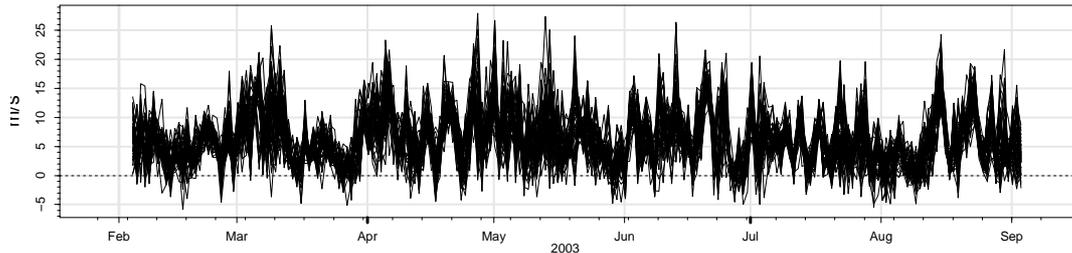




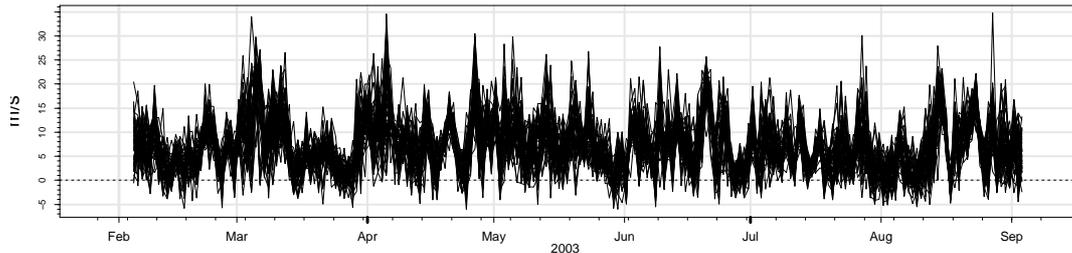
2d 12h



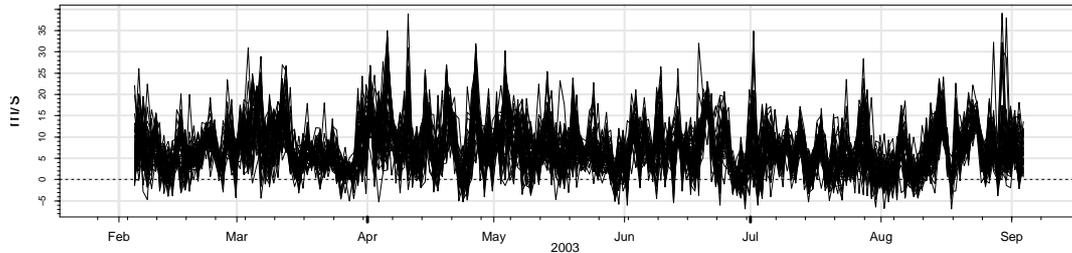
2d 18h



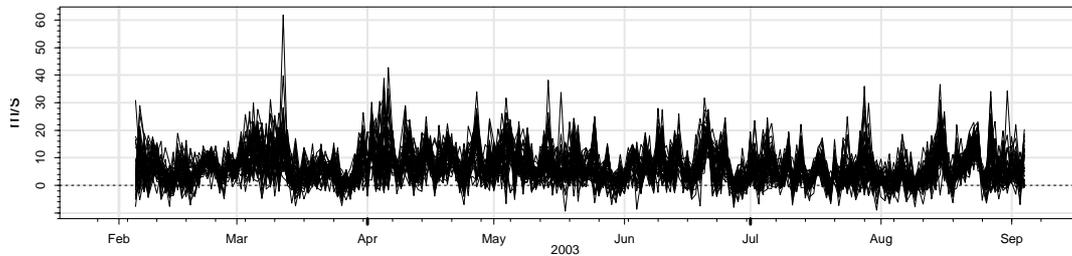
3d 0h

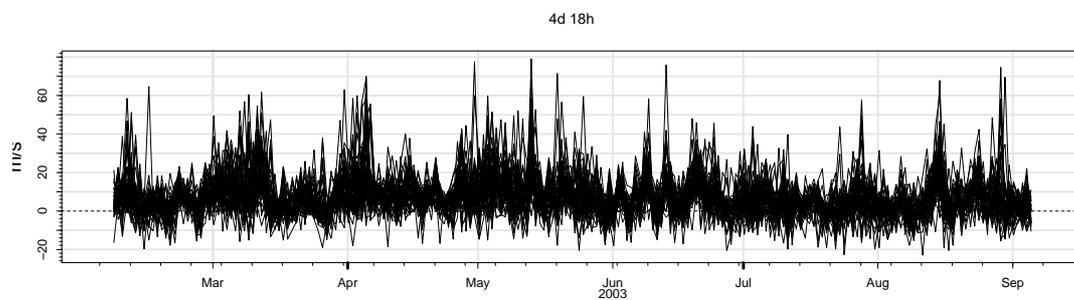
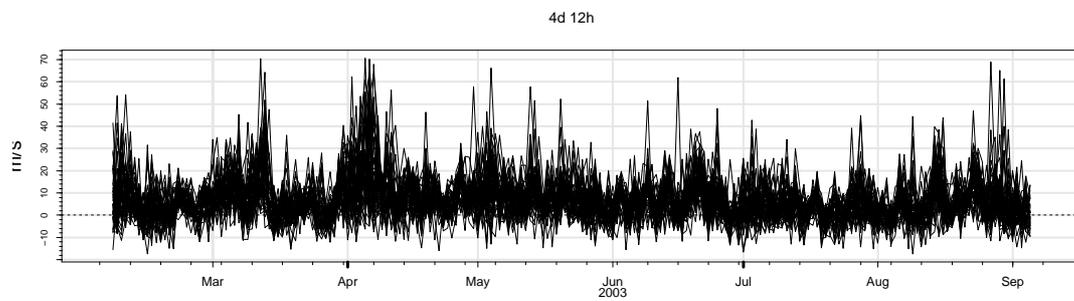
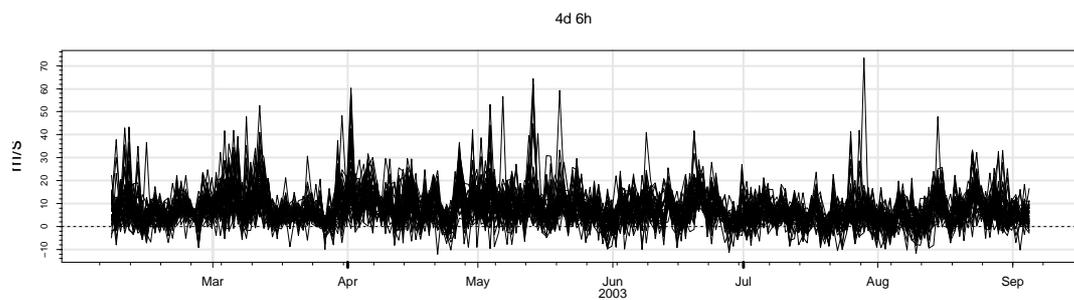
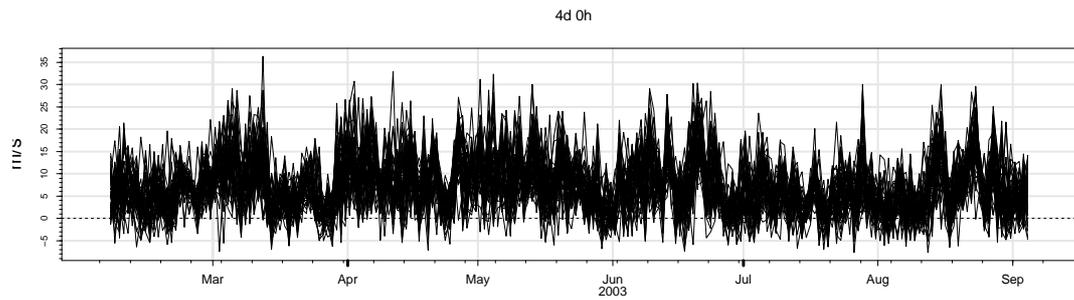
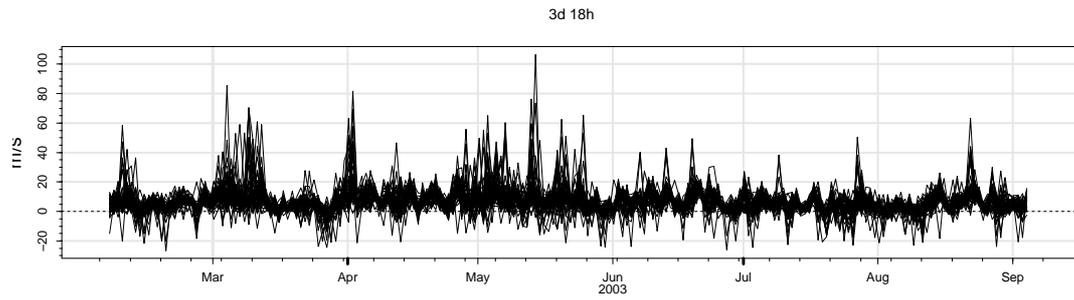


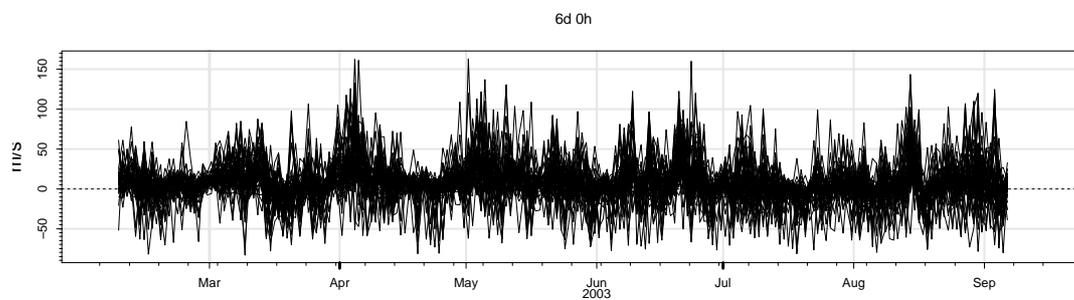
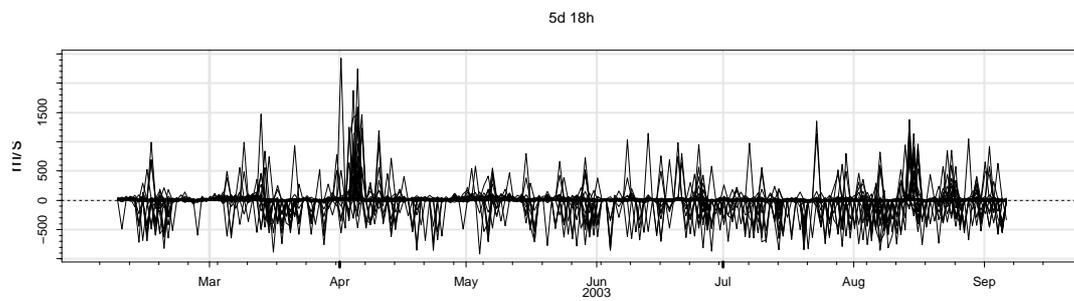
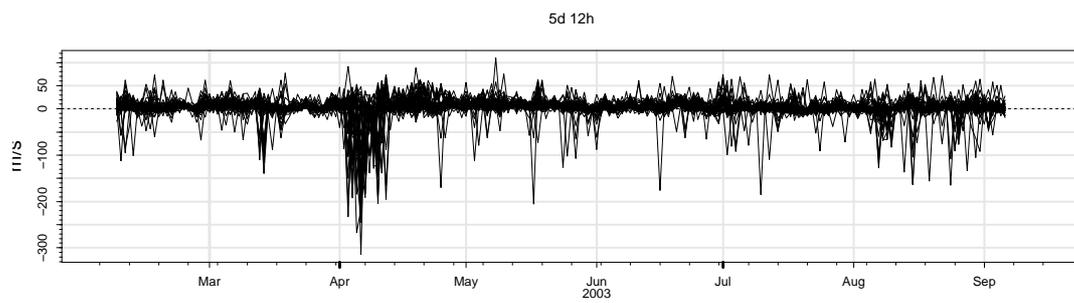
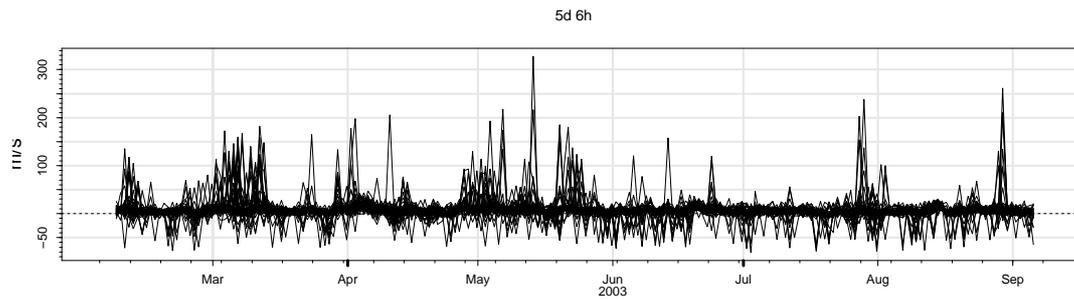
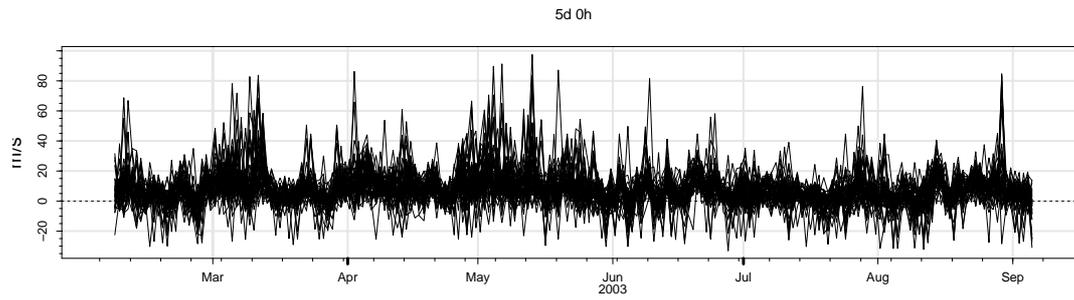
3d 6h



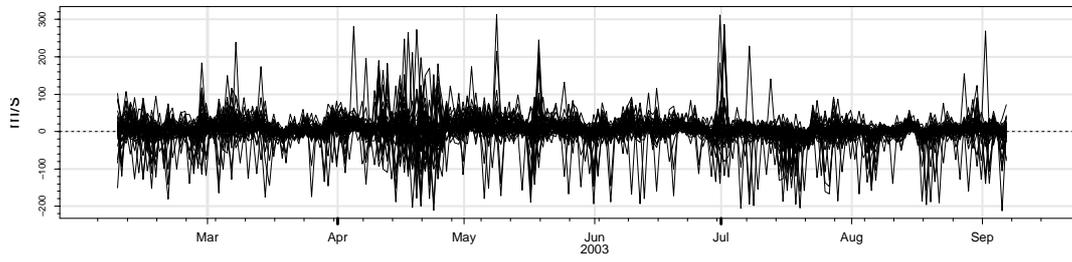
3d 12h



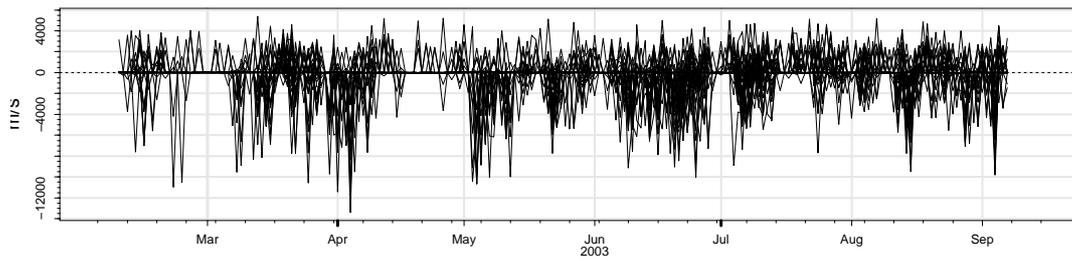




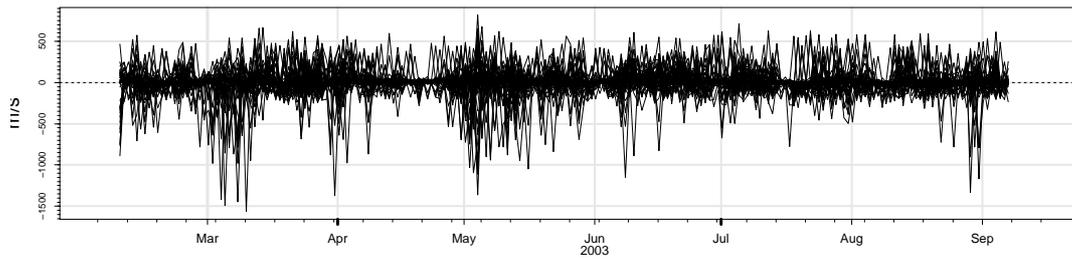
6d 6h



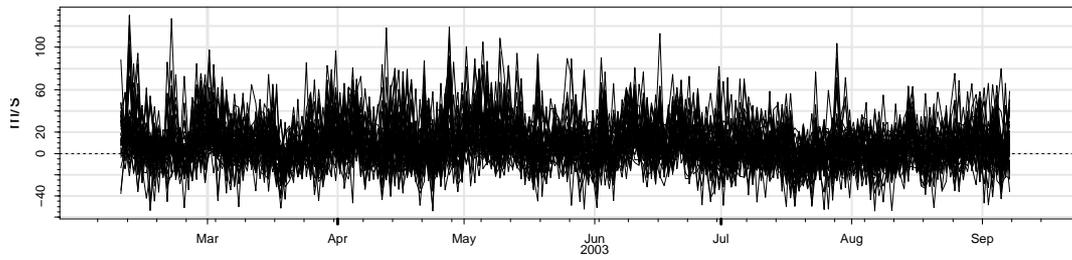
6d 12h



6d 18h



7d 0h



## A.2 HIRLAM corrected using the unperturbed forecast as the dependent variable.

