

# Calibration of medium-range metocean forecasts for the North Sea

Conor Murphy<sup>a</sup>, Ross Towe<sup>b</sup>, Philip Jonathan<sup>a,\*</sup>

<sup>a</sup>*School of Mathematical Sciences, Lancaster University LA1 4YF, United Kingdom.*

<sup>b</sup>*Shell Information Technology International Ltd., London SE1 7NA, United Kingdom.*

---

## Abstract

We assess the value of calibrating forecast models for significant wave height  $H_S$ , wind speed  $W$  and mean spectral wave period  $T_m$  for forecast horizons between zero and 168 hours from a commercial forecast provider, to improve forecast performance for a location in the central North Sea. We consider two straightforward calibration models, linear regression (LR) and non-homogeneous Gaussian regression (NHGR), incorporating deterministic, control and ensemble mean forecast covariates. We show that relatively simple calibration models (with at most three covariates) provide good calibration and that addition of further covariates cannot be justified. Optimal calibration models (for the forecast mean of a physical quantity) always make use of the deterministic forecast and ensemble mean forecast for the same quantity, together with a covariate associated with a different physical quantity. The selection of optimal covariates is performed independently per forecast horizon, and the set of optimal covariates shows a large degree of consistency across forecast horizons. As a result, it is possible to specify a consistent model to calibrate a given physical quantity, incorporating a common set of three covariates for all horizons. For NHGR models of a given physical quantity, the ensemble forecast standard deviation for that quantity is skilful in predicting forecast error standard deviation, strikingly so for  $H_S$ . We show that the consistent LR and NHGR calibration models facilitate reduction in forecast bias to near zero for all of  $H_S$ ,  $W$  and  $T_m$ , and that there is little difference between LR and NHGR calibration for the mean. Both LR and NHGR models facilitate reduction in forecast error standard deviation relative to naive adoption of the (uncalibrated) deterministic forecast, with NHGR providing somewhat better performance. Distributions of standardised residuals from NHGR are generally more similar to a standard Gaussian than those from LR.

*Keywords:* metocean, forecast, calibration, weather, linear regression, non-homogeneous Gaussian regression.

---

## 1. Introduction

Safe execution of offshore activities requires the forecasting of environmental time series to improve decision making, for e.g. platform evacuation in severe weather, wind farm maintenance scheduling, and on- and off-loading from floating LNG facilities. Good forecast performance for a range of environmental conditions is particularly important, as is reliable quantification of forecast uncertainty; in principle, we are interested in forecasting the full joint spatio-temporal distribution of metocean sea state variables (incorporating sea state significant wave height  $H_S$ , mean wave period  $T_M$  and wind speed  $W$ ) well, including marginal and joint extremes.

Weather-forecasting organisations now routinely provide metocean forecasts to inform offshore activities. Usually, for use at specific locations, it is possible to calibrate the original forecasts further, so that the calibrated forecast exhibits smaller bias and uncertainty than the original forecast. Moreover, modern forecasts tend to come in the form of a combination of different components, including e.g. a deterministic forecast, a control forecast and an ensemble of forecasts representing a range of possible future metocean temporal trajectories; all of these are available to calibrate the original forecast.

There is a large literature on forecast calibration. Gneiting et al. (2005) presents a calibration method for probabilistic forecasts, based on multiple linear regression, to address both forecast bias and under-dispersion. Pinson et al. (2009) discusses the generation of statistical scenarios of short-term wind generation that accounts for interdependence of prediction errors and predictive distributions of wind power production. Gneiting (2011) provides a framework for estimation and evaluation of point forecasts. Gneiting (2014) discusses calibration of medium-range weather forecasts, presenting alternative strategies including Bayesian model averaging (BMA), non-homogeneous Gaussian regression (NHGR) and empirical copula coupling (ECC). The latter incorporates marginal calibration of raw ensembles together with a reordering to retain rank correlation structure across variates. Schefzik et al. (2013) discusses ECC

---

\*Corresponding author p.jonathan@lancaster.ac.uk

for uncertainty quantification in computer models. Gneiting and Katzfuss (2014) provides a review of probabilistic forecasting. Bessa et al. (2017) and Sweeney et al. (2020) review the challenges facing forecasters in electric power and renewable energy, emphasising the advantages of probabilistic methods. Gilbert et al. (2021) proposes boosted semi-parametric models for probabilistic forecasts which outperform those estimated via maximum likelihood. Heinrich et al. (2021) proposes probabilistic post-processing of multivariate forecasts, incorporating moving averages and covariance matrix regularization, allowing for non-stationary, non-isotropic and negative correlations in the forecasting error. van der Meer (2021) proposes a multivariate probabilistic ensemble model to forecast solar irradiance. Bjerregard et al. (2021) provides an introduction to multivariate probabilistic forecast evaluation. Gao et al. (2022) consider probabilistic forecasting of Arctic Sea ice thickness. Astfalck et al. (2023) emphasises that proper scoring rules (see e.g. Winkler 1969, Gneiting et al. 2007) are necessary for the evaluation of probabilistic forecasts, and illustrates their performance in an operational maritime engineering context. Proper scoring rules evaluate forecasting performance in terms of forecast sharpness and calibration, such that a model score is optimized when the reported forecast distribution is equal to the true predictive distribution given information (from e.g. data). Allen et al. (2023) considers the adoption of transformed kernel scores to emphasise the importance of forecasting events with high impact. Adnan et al. (2023) uses BMA in conjunction with various machine learning methods to provide short-term probabilistic forecasting of  $H_S$ . Cerqueira and Torgo (2024) discusses direct forecasting of exceedance probability for  $H_S$ . Hoehlein et al. (2024) examines the use of permutation-invariant neural networks for ensemble-based forecasting, that treat forecast ensembles as a set of unordered member forecasts, learning link functions that are by design invariant to permutations of the member ordering. Tyralis and Papacharalampous (2024) provides a recent review of predictive uncertainty estimation using machine learning.

### *Objectives and outline*

The typical practising metocean engineer uses the simplest tools (e.g. linear regression) for calibration, rarely exploiting the richness of data provided by modern forecast models. The aspiration of the current paper is to demonstrate that there is material benefit from exploiting the output of modern forecasts more fully within pragmatic forecast calibration procedures for realistic offshore application. Specifically we do not claim that the calibration procedures considered are the best currently available, but we do contend that they provide useful tools with which the metocean engineer can reasonably expect to improve the performance of their forecast calibrations.

The objective of the current work is to evaluate the performance of simple methods to calibrate forecasts for  $H_S$ ,  $T_M$  and  $W$ , based on forecasts from a commercial forecast provider and offshore measurements for a central North Sea (CNS) location. Specifically, we evaluate the relative performance of different marginal linear regression (LR) and non-homogeneous Gaussian regression (NHGR) calibration models for a given measured metocean variable (one of  $H_S$ ,  $T_M$  and  $W$ ) in terms of deterministic, control and ensemble forecasts for the variable, and also potentially forecasts for other variables.

The layout of the article is as follows. Section 2 describes the motivating CNS application, and Section 3 outlines the calibration methodologies employed. Section 4 provides results from the calibration studies, including assessment of both within-sample and out-of-sample performance. Section 5 provides discussion and conclusions. Online Supplementary Material (SM) provides supporting figures.

## **2. Motivating application**

We consider the calibration of forecast data, from weather forecast provider StormGeo, for a location in the central North Sea. At this location, we have access to in-situ measured data for sea state  $H_S$  and mean spectral wave period  $T_m$  (measured using a downward-looking Saab Rex wave radar) and wind speed  $W$  (measured using a Gill ultrasonic wind sensor) for the interval 17 May 2022 to 6 September 2023.  $W$  is sampled at 10 minute intervals, whereas  $H_S$  and  $T_m$  estimates are provided every 30 minutes. These data were averaged to provide hourly values.

Forecast data for the same period was issued every 6 hours (from mid-night), providing forecasts for horizons at three-hourly resolution (for forecast horizons  $\leq 72$  hours) and at six-hourly resolution for longer forecast horizons  $\leq 168$  hours. At each issue time, the forecast data consist of a deterministic forecast, a control forecast and an ensembles forecast with 50 exchangeable members (see Section 3.1) for all of  $H_S$ ,  $W$  and  $T_m$ . These will be referred to henceforth as “forecast components” for definiteness.

Figure 1 provides illustrations of the data for three forecast issue times as described in the figure caption. For  $H_S$  and  $W$ , the correspondence between reality and forecast components is generally good. For  $T_m$ , there is evidence that ensemble forecasts are somewhat larger than both reality and the deterministic forecast.

Figure 2 gives scatter plots of the deterministic forecast on reality, and individual ensemble member forecasts on reality, for three forecast horizons, for the full period of observation. As would be expected, the now-cast prediction (forecast horizon = 0 hours) provides the least scatter about the  $y = x$  line of agreement; scatter about this line

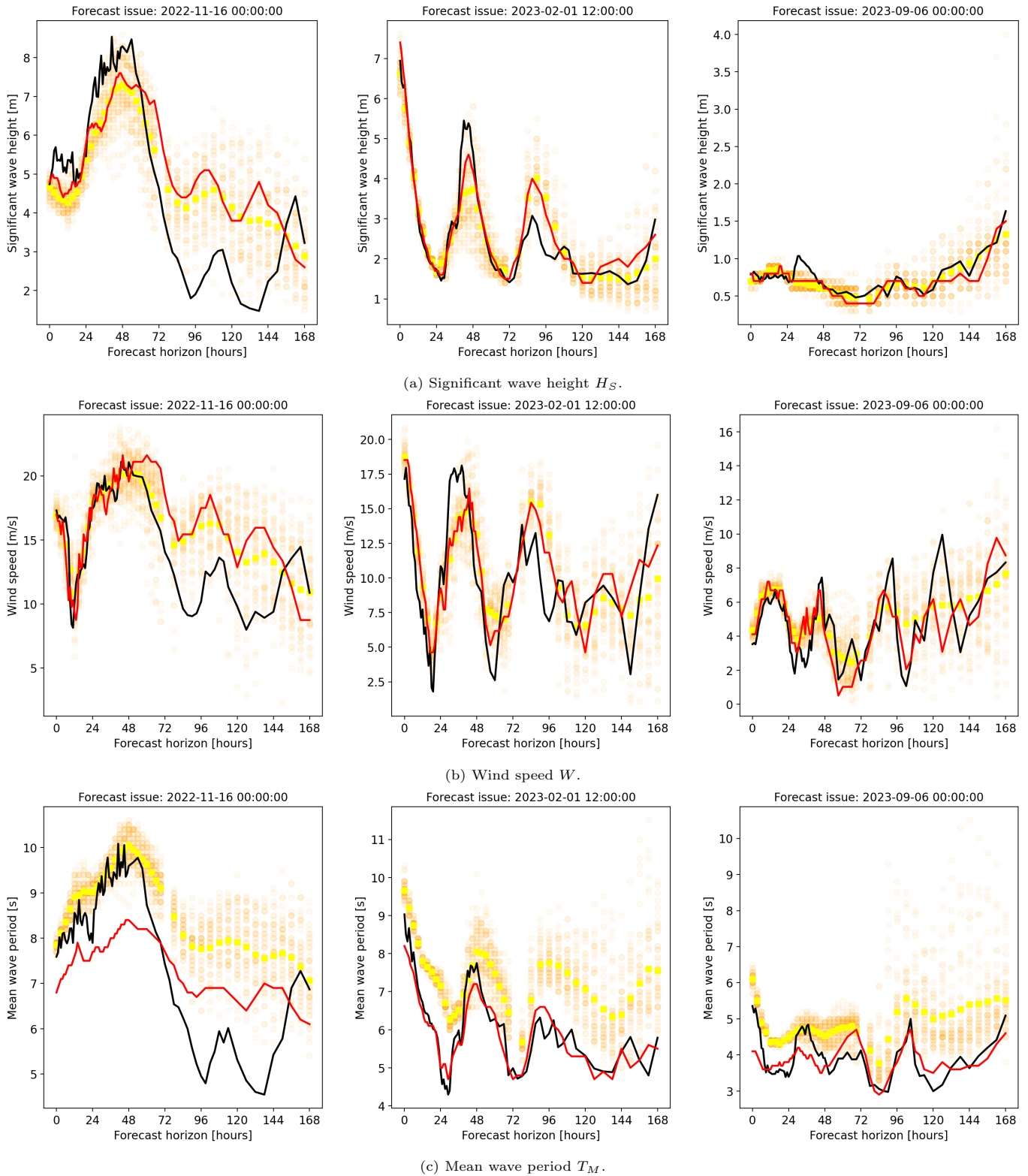


Figure 1: Forecasts and (future) reality for forecast components at three illustrative forecast times. Columns show forecasts issued on 16 November 2022 (left), 1 February 2023 (centre) and 6 September 2023 (right) with forecast horizons  $\in [1, 168]$  hours for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ . Title of each column gives the forecast issue time. (Future) reality illustrated using black line. Deterministic forecast illustrated using a red line. Forecast ensemble members shown as orange circles, and the ensemble mean as yellow circles. Note that the future reality and deterministic forecast lines are shown as continuous for ease of inspection, whereas in fact they are discrete as described in the text.

increases with forecast horizon for all variables. For  $H_S$  and  $W$ , the clouds of red and yellow points are centred on the line of agreement, indicating that the bias in these forecasts is low. The bias of the deterministic forecast for  $T_m$  is also relatively low, but considerable bias for the ensemble forecast is present, suggesting that a simple offset (or linear regression calibration) model would improve forecast performance. There is some evidence also, for all variables, that forecasts of the highest values may be biased low (see e.g. the forecast of  $H_S$  at forecast horizons of 78 and 168 hours.) The figure therefore suggests that the characteristics of the calibration model should vary with both the value of the actual response, and the forecast horizon.

Figure 3 quantifies these findings, by summarising the bias and standard deviation of forecast components for the full period of observation. There is a small positive bias in the deterministic and ensemble mean forecast for  $H_S$  of around 0.05 m for all forecast horizons. For  $W$ , the absolute mean ensemble bias is  $\leq 0.1$  m/s over all horizons, whereas bias for the deterministic forecast shows a more systematic trend with absolute mean ensemble bias is  $\leq 0.3$  m/s everywhere. For  $T_m$ , the deterministic forecast shows a typical bias of around -0.1 s everywhere, whereas as the ensemble bias is relatively large at around 1.1 s everywhere. In terms of forecast standard deviation, the expected trend of increasing forecast error with horizon is observed for all variables. At longest horizons, the ensemble mean is seen to provide lower forecast error than the deterministic forecast for all variables; this effect is particularly strong for  $H_S$  and  $W$ . For  $T_m$ , the deterministic forecast yields lowest standard deviations at short horizons. Again as would be expected, the forecast error from individual (exchangeable) forecast ensemble members is poorer than from the ensemble mean.

### 3. Methodology

In this section we introduce two straightforward methods used to establish calibration models, which provide best estimates of (future) reality (i.e. of each of measured  $H_S$ ,  $W$  and  $T_m$  independently) as a function of forecast components (i.e. the deterministic, control and ensemble forecasts) for the same set of metocean variables. The first approach is simple linear regression (LR), outlined in Section 3.2; the second is an extension of linear regression to accommodate non-stationary error variance, namely non-homogeneous Gaussian regression (NHGR), outlined in Section 3.3. The performance of calibration models with different levels of complexity is evaluated as a function of forecast horizon by comparison of values for the Akaike Information Criterion (AIC), as explained in Section 3.4. Further, for ease of interpretation of regression output, we choose to standardise the covariates in a particular manner, as discussed in Section 3.5. We stress that these choices of calibration models are meant to represent the simplest approaches the metocean engineer might consider for practical application. Results of the modelling exercise are then presented in Section 4. First, we outline the manner in which we accommodate the exchangeable nature of the ensemble members in the calibration models.

#### 3.1. Accommodating exchangeable ensemble members in a calibration model

A forecast ensemble of 50 members is available for each of  $H_S$ ,  $W$  and  $T_m$  for each combination of forecast issue time  $t$  hours and forecast horizon  $\tau$  hours, with  $t, \tau \in \mathbb{Z}_{\geq 0}$ . These ensemble members are exchangeable, in the sense that for any two forecast issue times  $t$  and  $t'$ , the association between the values  $E_j(t)$  and  $E_{j'}(t')$  for ensemble members  $E_j$  and  $E_{j'}$  is no different when  $j = j'$  to when  $j \neq j'$ ,  $j, j' = 1, 2, \dots, 50$ . Hence it makes no sense to include individual ensemble members directly as covariates in the calibration model. However, we are free to use summary statistics of the 50-member ensemble, such as the ensemble mean  $M_E(\tau; t)$  and the ensemble standard deviation  $S_E(\tau; t)$ , the values of which are unaffected by permutations of the ensemble members (see e.g. Gneiting 2014). For this reason, in the LR and NHGR models below, the only ensemble covariates we consider are the ensemble mean  $M_E(\tau; t)$  (in the LR and NHGR mean), and the ensemble standard deviation  $S_E(\tau; t)$  (for the NHGR error variance).

#### 3.2. Linear regression (LR)

We estimate a linear regression model for a directly-measured metocean variable  $Y \geq 0$  (i.e. one of  $H_S$ ,  $W$  and  $T_m$ ) at time  $t + \tau$  based on forecast components for the *same* metocean quantity; i.e. the deterministic, control and ensemble mean forecasts for that quantity, written  $\mathbf{X}(\tau; t) = (X_1(\tau; t), X_2(\tau; t), \dots, X_j(\tau; t), \dots, X_{n_x}(\tau; t))$ , and forecast components  $\mathbf{Z}(\tau; t) = (Z_1(\tau; t), Z_2(\tau; t), \dots, Z_j(\tau; t), \dots, Z_{n_z}(\tau; t))$  for *other* metocean quantities at forecast issue time  $t$  and forecast horizon  $\tau$ . With  $X_j(\tau; t), Z_k(\tau; t) \geq 0 \forall j, k = 1, 2, 3, \dots$ , we write

$$Y(t + \tau | \mathbf{X}(\tau; t) = \mathbf{x}, \mathbf{Z}(\tau; t) = \mathbf{z}) \sim N(a(\tau) + \sum_{j=1}^{n_x} b_j(\tau)x_j + \sum_{k=1}^{n_z} c_k(\tau)z_k, s^2) \quad (1)$$

for real-valued parameters  $a(\tau)$ ,  $b_j(\tau)$  and  $c_j(\tau)$ , for values  $\mathbf{x} = (x_1, x_2, \dots, x_j, \dots, x_{n_x})$  of covariate vector  $\mathbf{X}(\tau; t)$ , and values  $\mathbf{z} = (z_1, z_2, \dots, z_j, \dots, z_{n_z})$  of covariate vector  $\mathbf{Z}(\tau; t)$ . The regression error variance  $s^2$  (and standard deviation



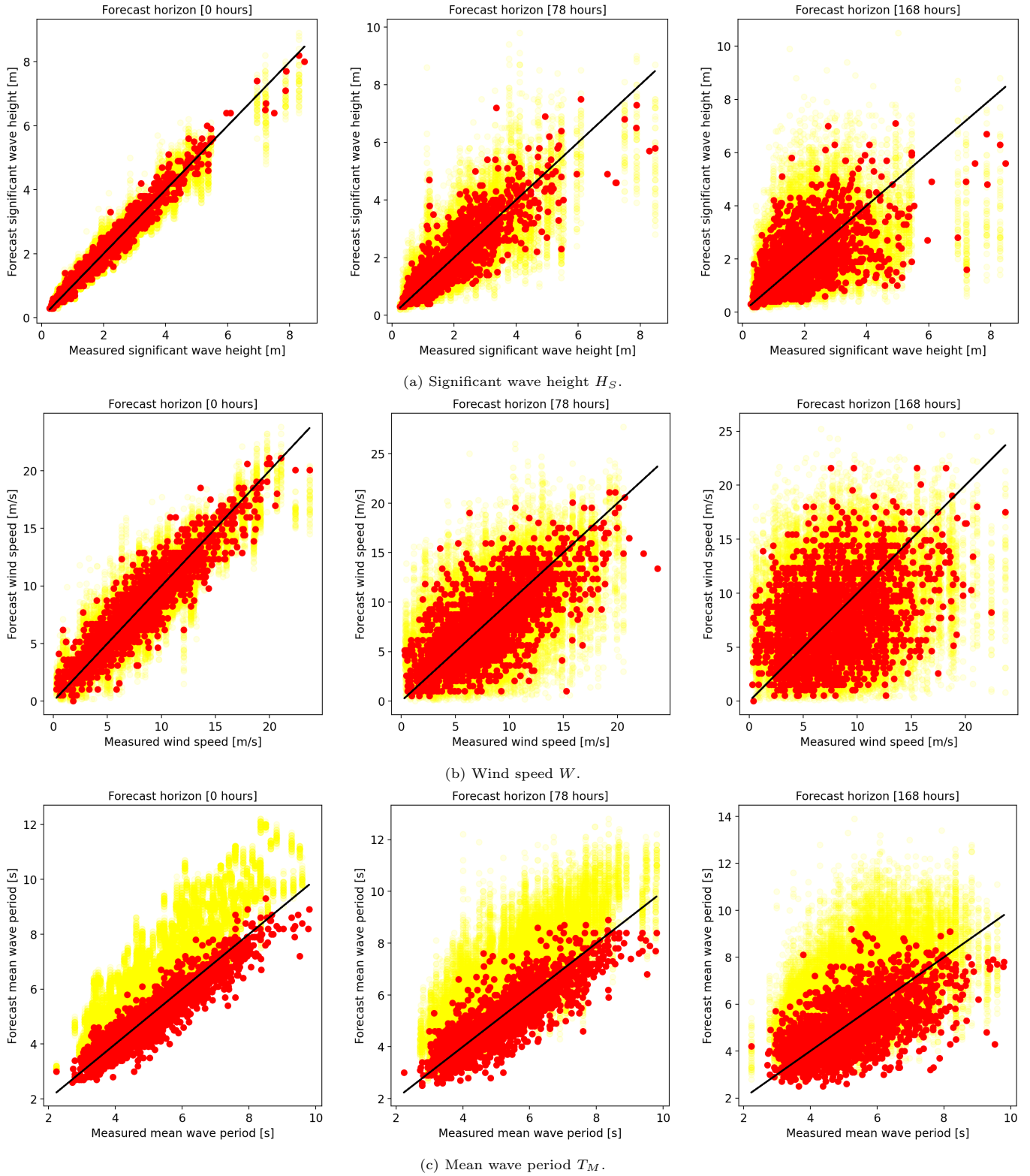
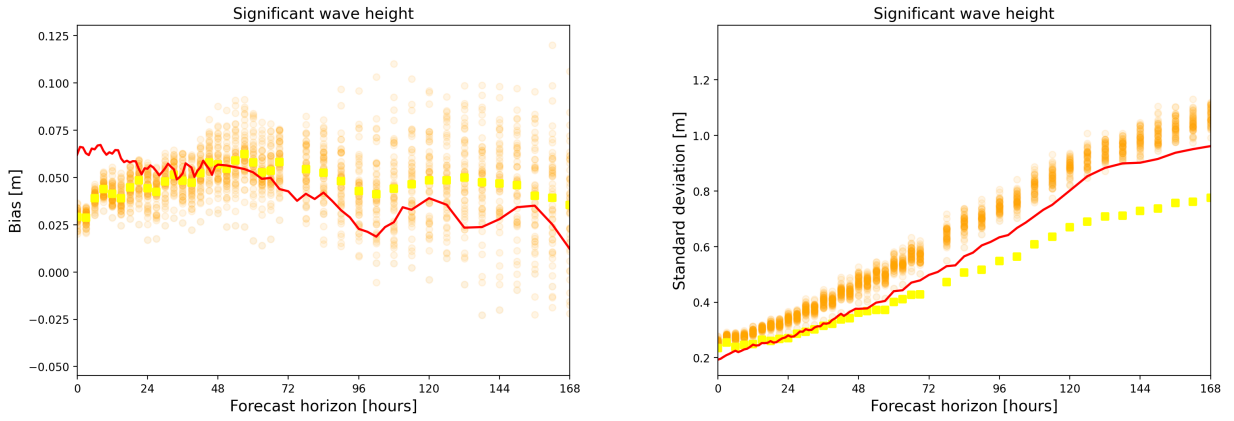
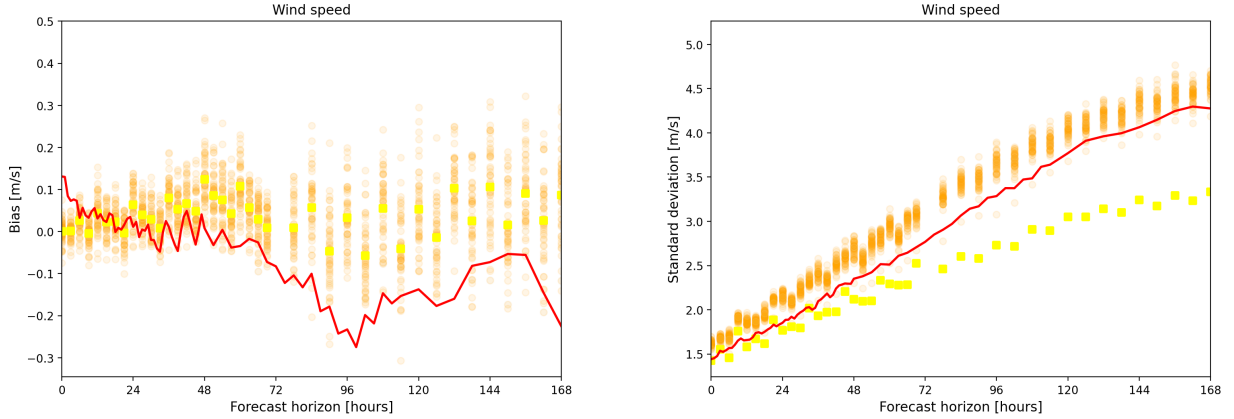


Figure 2: Scatter plots of forecast values ( $y$ ) on measured ( $x$ ) corresponding to three forecast horizons (columns) for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ . Title of each panel gives the forecast horizon. Deterministic forecasts shown in red. Individual ensemble member forecasts shown in yellow.

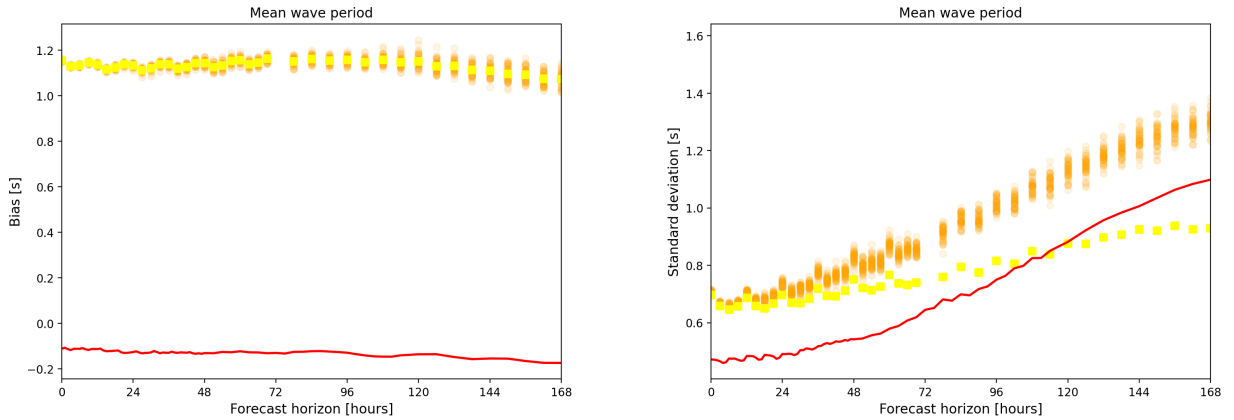
$s > 0$ ) is assumed fixed for all forecast horizons  $\tau$ . We estimate the parameters  $a$ ,  $b$  and  $c$  using maximum likelihood (or equivalently least squares). Once the parameter estimates are available, we estimate  $s^2$  as the ratio of the residual sum of squares for the fitted model, and the number of degrees of freedom in the model.



(a) Significant wave height  $H_S$ .



(b) Wind speed  $W$ .



(c) Mean wave period  $T_M$ .

Figure 3: Bias (left) and standard deviation (right) of forecast error as a function of horizon for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ . Bias for deterministic forecast (red line), individual ensemble members (orange circles) and ensemble mean (yellow discs). Note that deterministic forecast bias line is shown as continuous for ease of inspection, whereas in fact it is only available at a discrete set of forecast horizons.

The covariate vector  $\mathbf{X}(\tau; t)$  always consists of one or more of the deterministic, control and ensemble mean forecasts for forecast horizon  $\tau$  issues at time  $t$ , for the metocean quantity being predicted. Thus the components of  $\mathbf{X}(\tau; t)$  in a model for measured  $H_S$  will consist of one or more of the  $H_S$  forecast components only; analogously, the covariate vector  $\mathbf{Z}(\tau; t)$  will consist of one or more of the forecast components for  $W$  and  $T_m$  only. We will refer to “type- $X$ ” and “type- $Z$ ” covariates in subsequent sections for clarity and brevity of description.

### 3.3. Non-homogeneous Gaussian regression (NHGR)

Using the notation of Section 3.2, non-homogeneous Gaussian regression (NHGR) can be seen as an extension of the linear regression in Equation 1 to non-stationary error variance. Thus, in the current work, the NHGR model form is

$$Y(t + \tau | \mathbf{X}(\tau; t) = \mathbf{x}, \mathbf{Z}(\tau; t) = \mathbf{z}, S_E(\tau; t) = s_E) \sim N(a(\tau) + \sum_{j=1}^{n_X} b_j(\tau)x_j + \sum_{k=1}^{n_Z} c_k(\tau)z_k, (d(\tau) + e(\tau)s_E)^2) \quad (2)$$

That is, the NHGR forecast error standard deviation  $d(\tau) + e(\tau)s_E > 0$  is assumed to be a linear function of the ensemble forecast standard deviation  $s_E > 0$ ; thus, when the ensemble forecast variability is large (e.g. for longer forecast horizons), a larger NHGR forecast error variance is implied (for the same choice of model parameters  $d(\tau)$  and  $e(\tau)$ ). As for LR, NHGR parameters  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$  are estimated using maximum likelihood estimation.

### 3.4. Identifying best-performing models using AIC

We assess the performance of models using AIC, defined as  $2p + 2\hat{\ell}$ , where  $p$  is the number of parameters in the model ( $p = n_X + n_Z + 1$  for LR, and  $p = n_X + n_Z + 3$  for NHGR), and  $\hat{\ell}$  is the value of the negative log likelihood of the sample for the model evaluated at its minimum with respect to model parameters; that is,  $\hat{\ell}$  is simply the sum of negative log Gaussian densities of the sample points at the maximum likelihood solution. Better performing models provide lower values for AIC (see Akaike 1974, Emiliano et al. 2014). For all responses, and for both LR and NHGR methodologies, we find that including a total of at most three covariates (i.e.  $n_X + n_Z \leq 3$  in Equations 1 and 2) reduces AIC sufficiently that addition of further covariates cannot be justified.

### 3.5. Standardisation of covariates

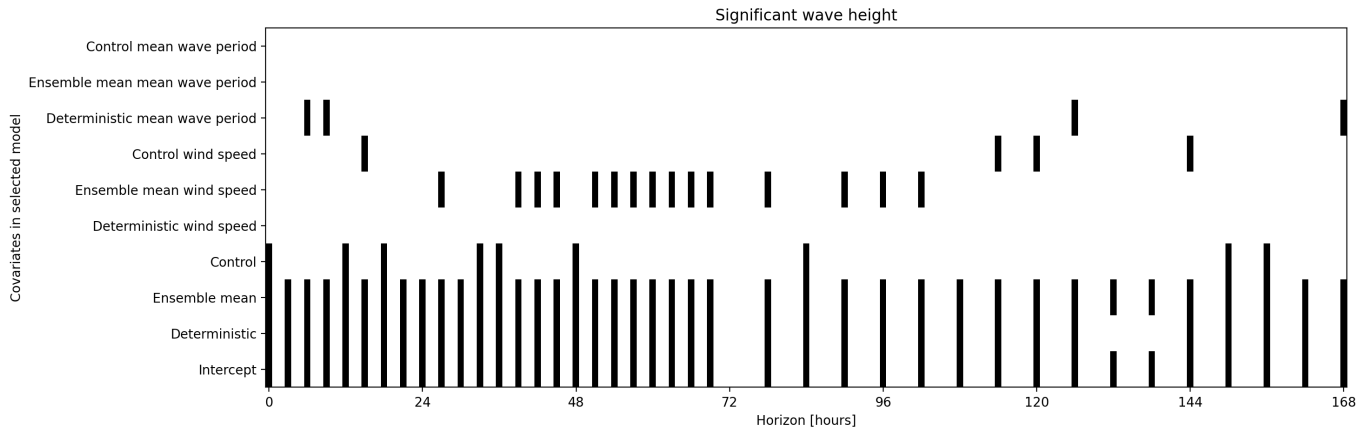
Standardisation of covariates often facilitates more intuitive interpretation of regression analysis. For this reason, in the current work, covariates for forecast components  $Z_k(\tau; t)$  ( $k = 1, 2, \dots, n_Z$ ) representing different (“type- $Z$ ”) physical quantities to the response  $Y(t + \tau)$  are standardised prior to model fitting so that their sample variance is equal to that of the response. Using this standardisation, the estimated regression coefficient  $c_k(\tau)$  indicates exactly the fraction of the response variance explained by covariate  $Z_k(\tau; t)$ ; e.g.  $c_k(\tau) = 1$  indicates that  $Z_k(\tau; t)$  explains exactly 100% of the response variance. Specifically, to achieve standardisation, we simply substitute the expression

$$z_k^* = \frac{\sigma_Y}{\sigma_{Z_k}}(z_k - \mu_{Z_k}) \quad (3)$$

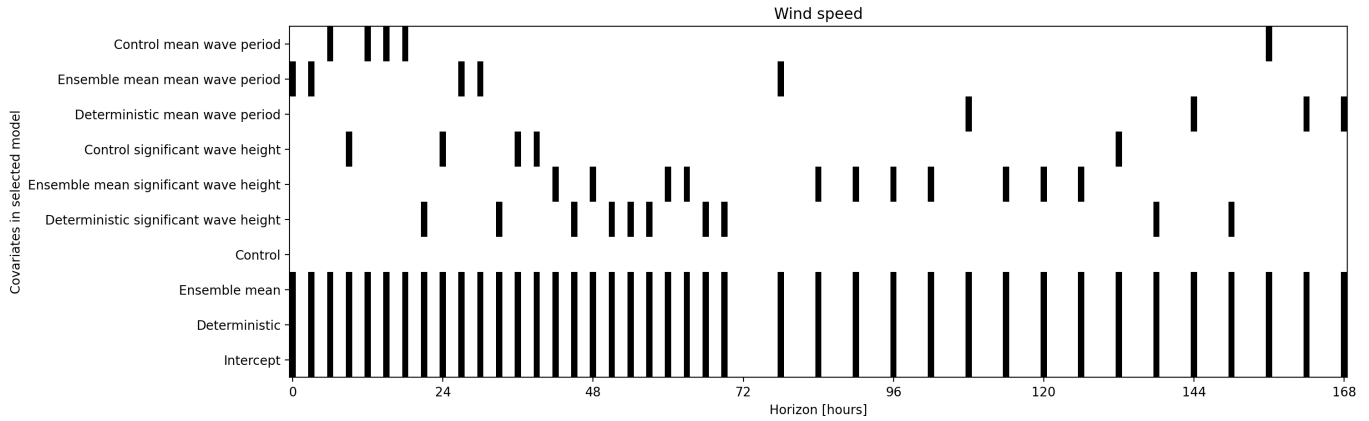
for  $z_k$  in Equations 1 and 2, where for  $\sigma_Y > 0$  and  $\sigma_{Z_k} > 0$  we use sample estimates for the standard deviation of  $Y(t + \tau)$  and  $Z_k(\tau; t)$ , and for  $\mu_{Z_k}$  we use the sample mean of  $Z_k(\tau; t)$ . However, covariates  $X_j(\tau; t)$  ( $j = 1, 2, \dots, n_X$ ) representing the same (“type- $X$ ”) physical quantity as the response are not standardised. In this case, an estimate for parameter  $b_j > 1$  would indicate that the variance of the corresponding forecast variable is smaller than that of the response. We stress that standardisation of covariates does not affect predictions made using the LR and NHGR models, nor the choice of optimal models; it is merely a convenient transformation to aid interpretation of regression coefficients. (For comparison, the Supplementary Material provides plots of relative contributions of covariates to the regression in which both response and all covariates are standardised to zero mean and unit standard deviation; see Figures SM1 and SM3).

## 4. Results

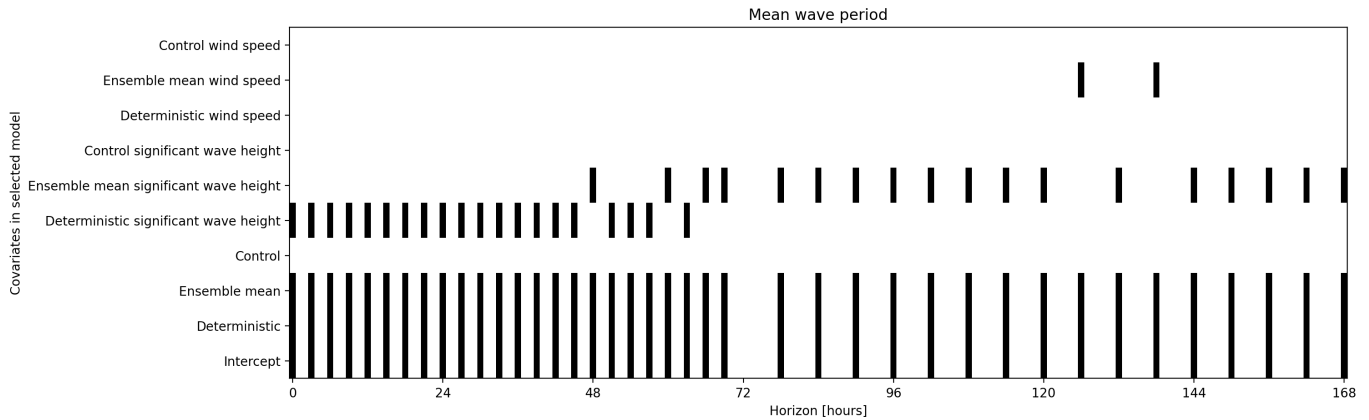
For each metocean response variable  $Y(t + \tau)$  (i.e. measured  $H_S$ ,  $W$  and  $T_m$ ), we estimate LR and NHGR calibration models independently, using covariates  $X(\tau; t)$  (i.e. one or more of the deterministic, ensemble mean and control forecasts for the same metocean variable) and covariates  $Z(\tau; t)$  (i.e. one or more of the deterministic, ensemble mean and control forecasts for different metocean variables). The estimated LR and NHGR models are then illustrated in more detail in Sections 4.1 and 4.2. Supporting plots are provided in on-line Supplementary Material (SM). Where possible, uncertainties in parameters and predictions are quantified using 95% confidence and Gaussian forecast uncertainty intervals available in closed form. Otherwise bootstrap resampling is used to estimate corresponding 95% uncertainty bands (estimated from repeated model fitting to bootstrap resamples of the original sample for model fitting). The characteristics of the original uncalibrated deterministic forecast, and those of LR- and NHGR-calibrated forecasts, are assessed within-sample in Section 4.3. In Section 4.4, we assess the performance of the three forecasts using data for two subsequent time periods, not used for calibration model estimation.



(a) Significant wave height  $H_S$ .



(b) Wind speed  $W$ .



(c) Mean wave period  $T_M$ .

Figure 4: Bar code plot of covariates included in optimal linear regression models for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_m$  for each forecast horizon. Terms included in consistent LR calibration model for each response are given in Table 1.

#### 4.1. Linear regression

For each combination of metocean response variable and forecast horizon, the optimal LR model was identified as that which minimises AIC, with the constraint that  $n_X + n_Z \leq 3$ . Bar code plots of the covariates present in optimal LR models for each forecast horizon are shown in Figure 4. Inspection of the figure suggests that the ensemble mean and deterministic forecasts (for the same physical quantity as the response) occur almost always in calibration models; the regression intercept is always imposed. A “consistent” model (with the same model form over all forecast horizons for a given variable) was then selected as that which minimises AIC most often over all forecast horizons. The chosen consistent model forms are listed in Table 1.

$H_S$	$W$	$T_m$
Intercept	Intercept	Intercept
Deterministic $H_S$	Deterministic $W$	Deterministic $T_m$
Ensemble mean $H_S$	Ensemble mean $W$	Ensemble mean $T_m$
Ensemble mean $W$	Ensemble mean $H_S$	Deterministic $H_S$

Table 1: Terms included in the consistent LR calibration models for significant wave height  $H_S$  (left), wind speed  $W$  (centre) and mean wave period  $T_m$  (right).

Figure 5 then illustrates the growth of AIC with forecast horizon of LR calibration models for each of  $H_S$ ,  $W$  and  $T_m$ . To provide a scale for comparisons, each panel of the figure gives AIC values for a calibration model using the deterministic forecast only (red), using the consistent calibration model for each response (cyan, see Table 1), and using the optimal calibration model for the combination of response and forecast horizon (orange, see Figure 4). We see that LR calibration using additional ensemble and deterministic covariates improves performance for all responses and forecast horizons relative to calibration using the deterministic forecast (for the same physical quantity as the response) only. Only for  $H_S$  and short forecast horizons is the deterministic forecast competitive. We note also that the consistent calibration model is competitive with the optimal model for all variables and forecast horizons, indicating that there is little benefit in selecting different LR calibration model forms for different forecast horizons. More generally, as might be expected, calibration model performance degrades with increasing forecast horizon.

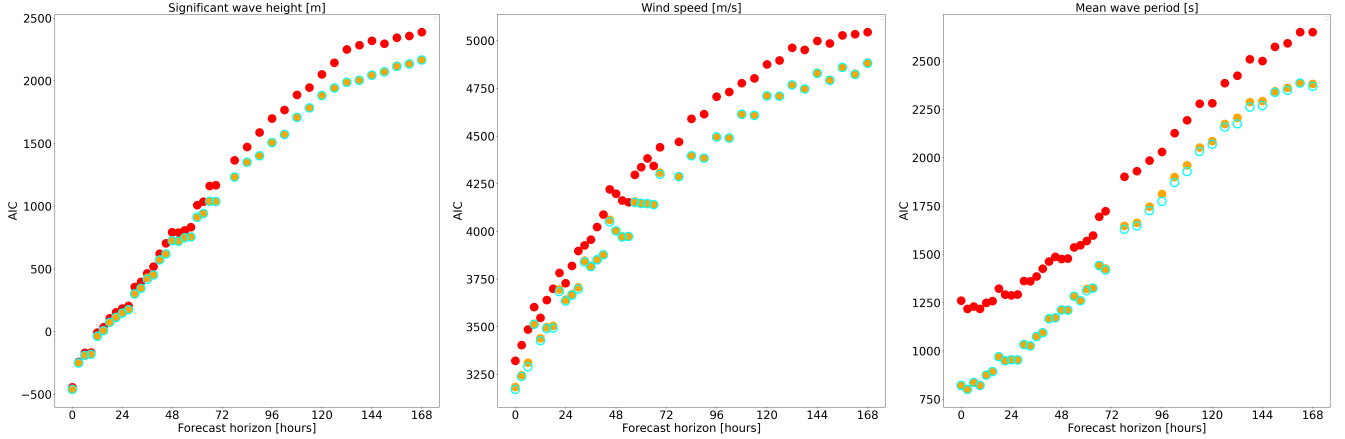
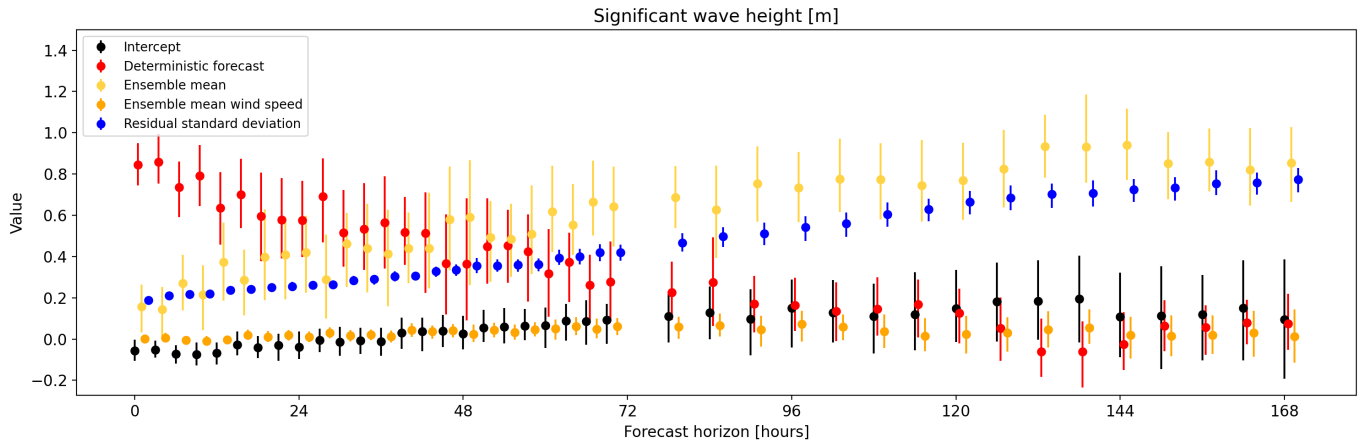


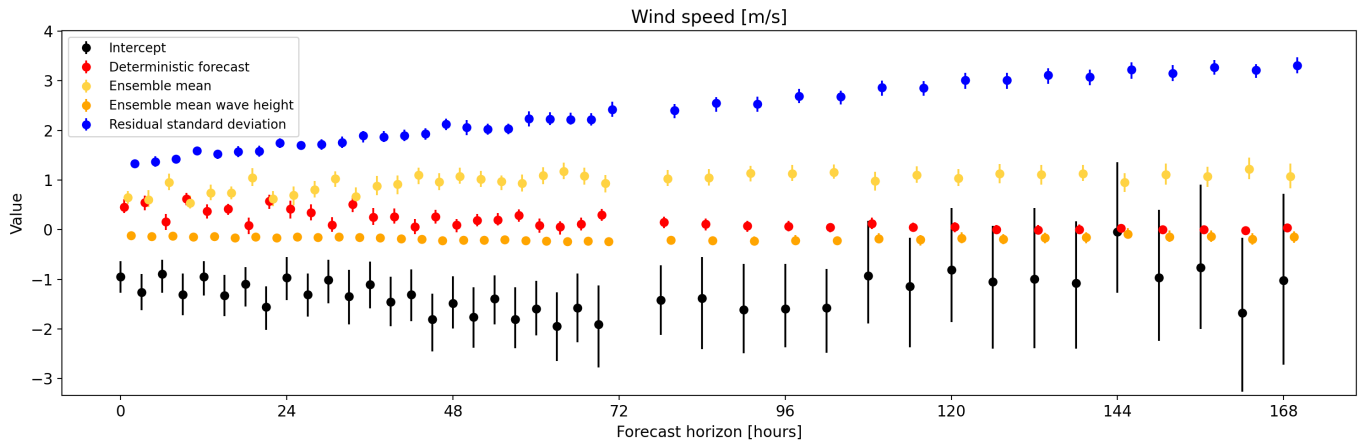
Figure 5: Variation of AIC with forecast horizon for significant wave height ( $H_S$ , left), wind speed ( $W$ , centre) and mean wave period ( $T_m$ , right). Each panel shows the growth of AIC with forecast horizon, for three calibration models: based on the deterministic forecast only (red disc), based on the consistent LR calibration model (cyan circle) and based on the optimal LR calibration model choice for that forecast horizon (orange disc). The consistent model form for each of  $H_S$ ,  $W$  and  $T_m$  is given in Table 1.

Parameter estimates for the consistent LR calibration models are illustrated in Figure 6. Because of the standardisation of covariates employed, we expect parameter estimates for covariates to lay in the interval  $(0, 1)$  generally, which is the case. We also report 95% confidence bands for the parameter estimates. For  $H_S$  in Figure 6(a), the deterministic forecast (for  $H_S$ , in red) dominates for short forecast horizons, but the ensemble mean forecast (for  $H_S$ , in light orange) dominates for long forecast horizons. The roles of the intercept (black) and ensemble mean  $W$  (dark orange) are minor. The residual standard deviation (blue) grows from 0.2 m for shortest forecast horizons to approximately 0.7 m at 168 hours forecast horizon. For  $W$  in Figure 6(b), the ensemble mean forecast  $W$  plays a more important role at all forecast horizons. The intercept provides a correction of around -1 m/s on average over all forecast horizons. Again the effect of the covariate corresponding to a Z-type physical quantity different to the response is small. For  $T_m$  in Figure 6(c), the contributions of all covariates are similar for short forecast horizons, but the ensemble mean forecast for  $T_m$  dominates at larger horizons. The intercept provides a correction of around 1.2 m/s on average over all forecast horizons. These trends are further summarised in Figures SM1 and SM3 in terms of plots of variable contributions.

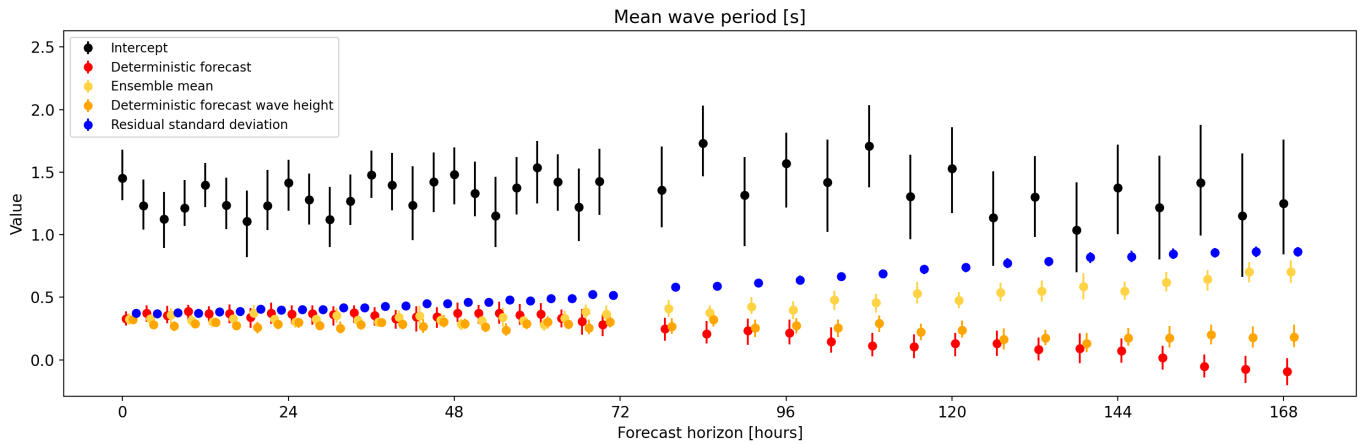
Figure 7 provides LR calibrated forecasts and (future) reality for each metocean response variable at three given forecast issue times. The actual measured response is shown in black, and the corresponding deterministic forecast in red. The ensemble mean forecast is in yellow. Note that both the deterministic and ensemble mean forecasts are uncalibrated. The consistent LR calibrated forecasts are shown as the mean (cyan disc) and estimated 95% Gaussian forecast uncertainty band (calculated using the estimated model error standard deviation, see Equation 1; cyan vertical line). For both  $H_S$  and  $W$ , there is little evidence that either the ensemble mean forecast or the mean calibrated LR



(a) Significant wave height  $H_S$



(b) Wind speed  $W$



(c) Mean wave period  $T_M$

Figure 6: Variation of estimated parameters with forecast horizon from consistent LR calibration models for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ . Consistent LR model forms are given in Table 1. For each parameter, the mean parameter estimate is indicated by a disc, and the 95% confidence interval by a vertical line. The intercept is given in black, the deterministic and ensemble mean forecasts (for the same physical quantity) in red and light orange. The remaining covariate (in dark orange) is specified in the plot legend for each panel. Also shown is the regression residual standard deviation (blue). Note that covariates are standardised as described in Section 3.5.

forecast provide much improvement over the deterministic forecast. The previously-noted bias in the ensemble mean forecast for  $T_M$  is clear. Visual inspection suggests that the prediction band, increasing in width with forecast horizon, generally does a relatively good job of including actual (future) reality.

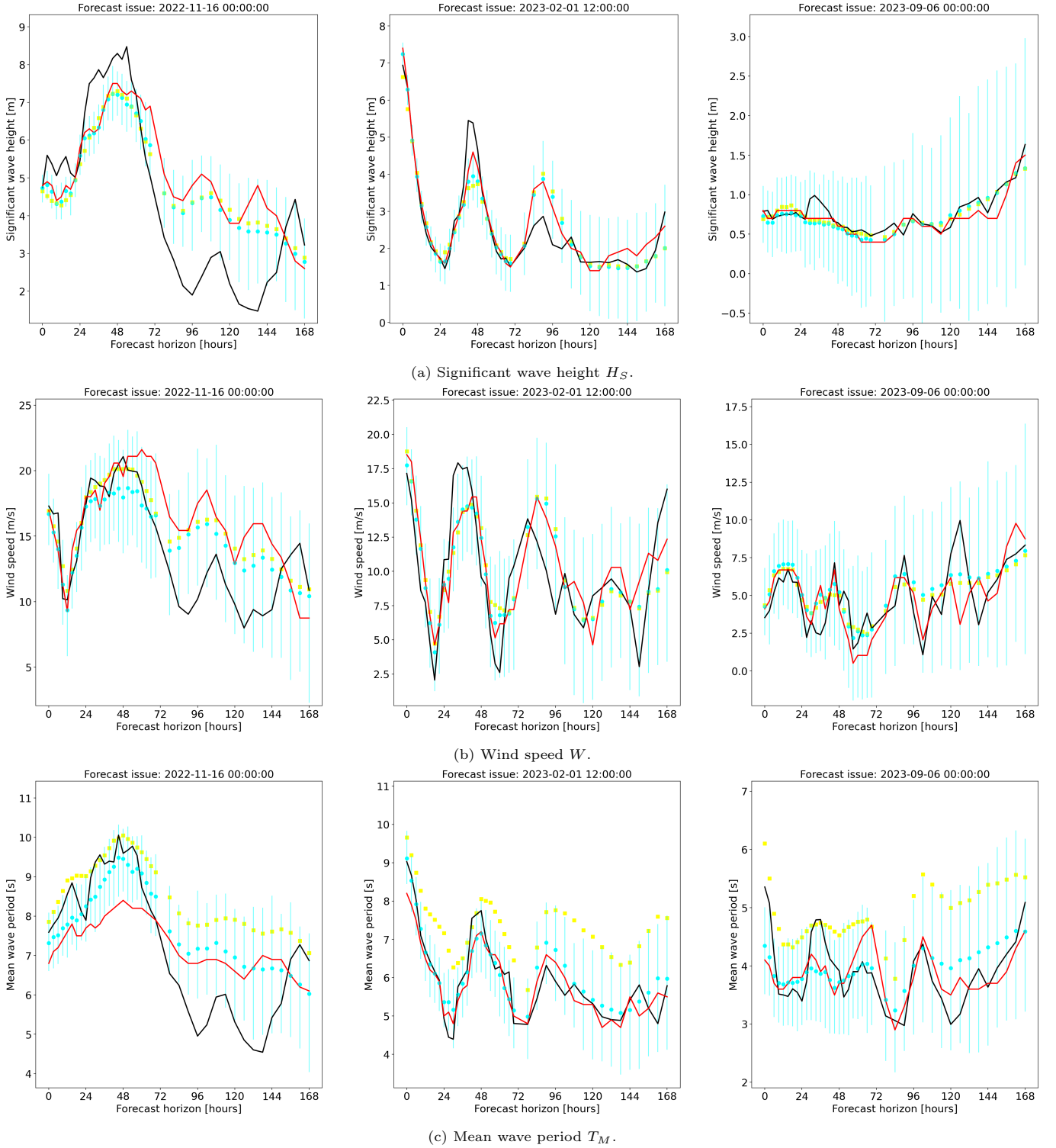


Figure 7: Consistent LR-calibrated forecasts and (future) reality for given variable at given forecast time. Three examples (columns) of forecasts on horizons  $\in [1, 168]$  for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ . The title of each column gives the time of forecast issue. (Future) reality illustrated using black line. The (uncalibrated) deterministic forecast illustrated using a red line. The (uncalibrated) ensemble mean forecast is shown as yellow discs. Optimal calibrated forecast given in cyan, in terms of the mean (disc) and 95% Gaussian forecast uncertainty band (vertical line). Note that (future) reality and deterministic forecast are shown as lines (rather than as discrete time points) for ease of interpretation. The forms for the consistent LR-calibrated models are given in Table 1.

## 4.2. NHGR

A model selection procedure similar to that described in Section 4.1 was used to identify optimal NHGR calibration models for each combination of metocean response variable and forecast horizon. The resulting optimal models are illustrated in the bar plot in Figure SM2, and the consistent model forms given in Table 2. As for optimal LR calibration, inspection of Figure SM2 shows that both the deterministic and ensemble mean forecast (type- $X$  covariates, for the same metocean quantity, see Section 3.2) are again included in the model, but that the type- $Z$  covariate differs for NHGR calibration; e.g. for  $H_S$ , LR calibration favours ensemble mean  $W$ , whereas NHGR calibration favours ensemble mean  $T_m$ . Note also that for all forecast horizons, a term in the ensemble standard deviation  $s_E$  (for the same metocean quantity) is included.

$H_S$	$W$	$T_m$
Intercept	Intercept	Intercept
Deterministic $H_S$	Deterministic $W$	Deterministic $T_m$
Ensemble mean $H_S$	Ensemble mean $W$	Ensemble mean $T_m$
Ensemble mean $T_m$	Deterministic $H_S$	Deterministic $H_S$

Table 2: Terms included in the consistent NHGR calibration models for significant wave height  $H_S$  (left), wind speed  $W$  (centre) and mean wave period  $T_m$  (right).

Model performance for NHGR calibration is again assessed using AIC, and illustrated in Figure 8. As a reference, AIC from consistent LR calibration is provided (as cyan discs), together with AIC for consistent NHGR calibration (orange discs) and for optimal NHGR calibration per forecast horizon (red circles). We see from the three panels that (i) consistent NHGR calibration is an improvement on consistent LR calibration, and (ii) there is minimal difference between the performance of consistent and optimal NHGR calibration. There is some evidence however that a different choice of type- $Z$  covariate (specifically, ensemble mean  $H_S$  forecast instead of deterministic  $H_S$  forecast; see Figure SM2) would improve performance for  $T_m$  at long forecast horizons.

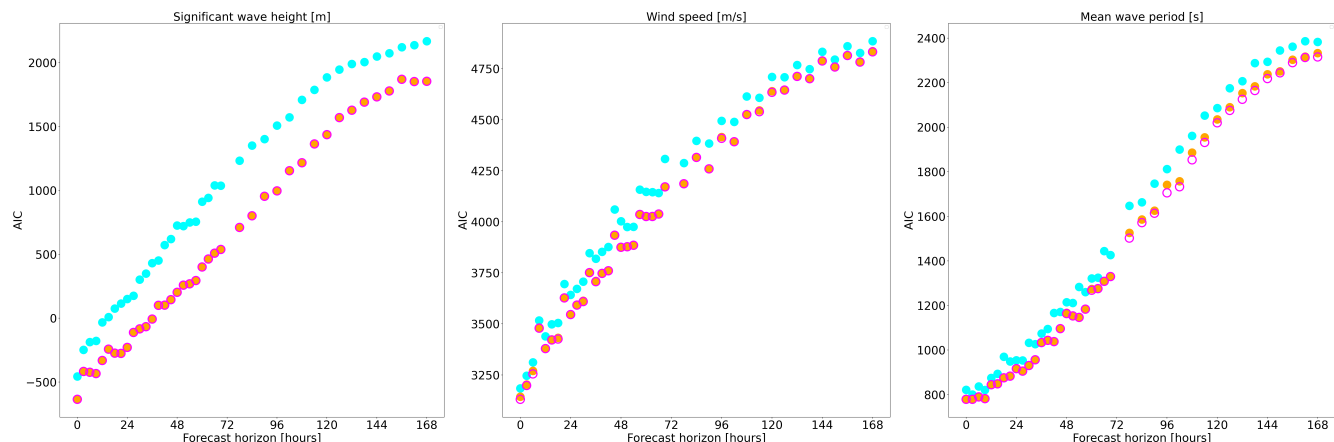
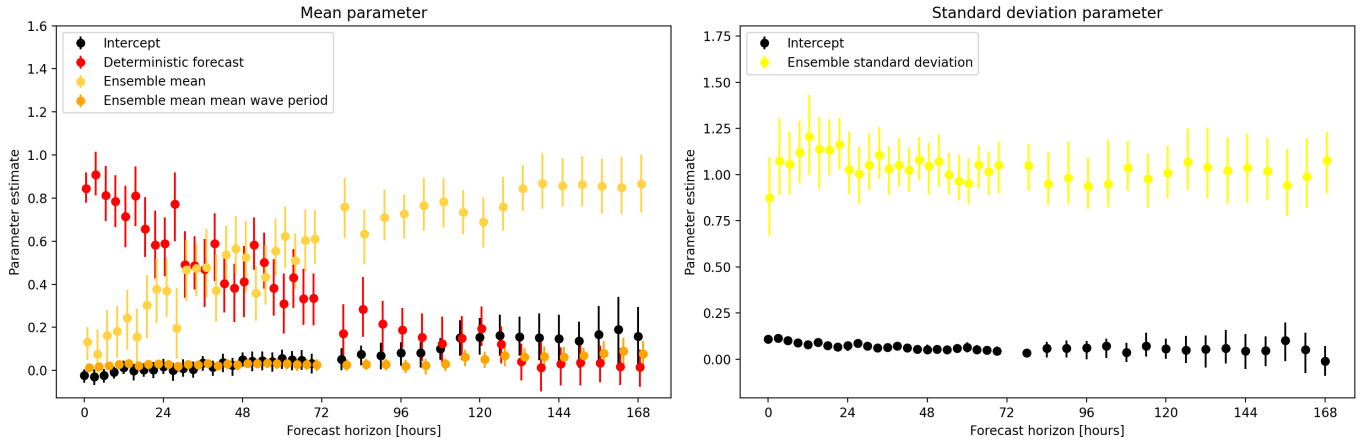


Figure 8: Variation of AIC with forecast horizon for significant wave height ( $H_S$ , left), wind speed ( $W$ , centre) and mean wave period ( $T_m$ , right). Each panel shows the growth of AIC with forecast horizon, for three calibration models: based on the consistent model for linear regression (cyan disc), based on the consistent model for NHGR (orange disc) and based on the optimal NHGR model choice (red circle). The consistent NHGR model form for each of  $H_S$ ,  $W$  and  $T_m$  is identified in Table 2.

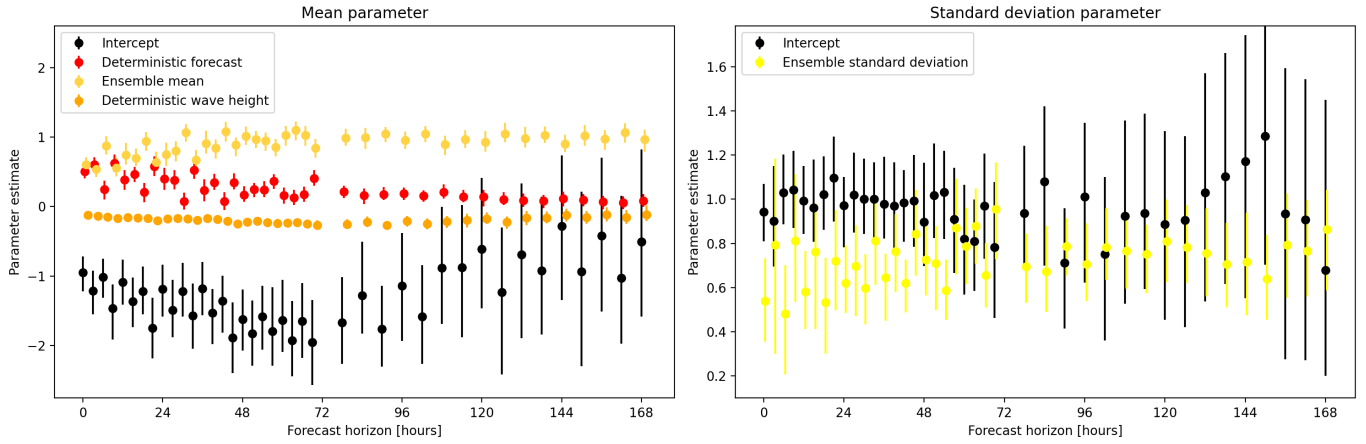
Parameter estimates for NHGR calibration, again using the standardisation scheme for covariates of the distributional mean described in Section 3.5, are illustrated in Figure 9. For the distributional mean parameters in the left hand panels, estimates are similar to those found using LR calibration for all of  $H_S$ ,  $W$  and  $T_m$ . For the distributional standard deviation parameters, it is striking that the ensemble standard deviation is very informative for  $H_S$ , and for  $T_m$  at short forecast horizons. Central 95% uncertainty bands for parameter estimates are calculated using bootstrap resampling. Figure SM3 quantifies benefit of incorporating ensemble variability (for the same physical quantity) to explain forecast variability in terms of percentage covariate contributions; covariate contributions are similar to those for LR. The importance of the ensemble standard deviation in estimating the distributional standard deviation is clear, especially for  $H_S$ .

Figure 10 compares forecasts and their variability from LR (cyan) and NHGR (magenta) models. Given observations already made, LR and NHGR are similar in terms of mean forecasts (discs in figure). The difference between

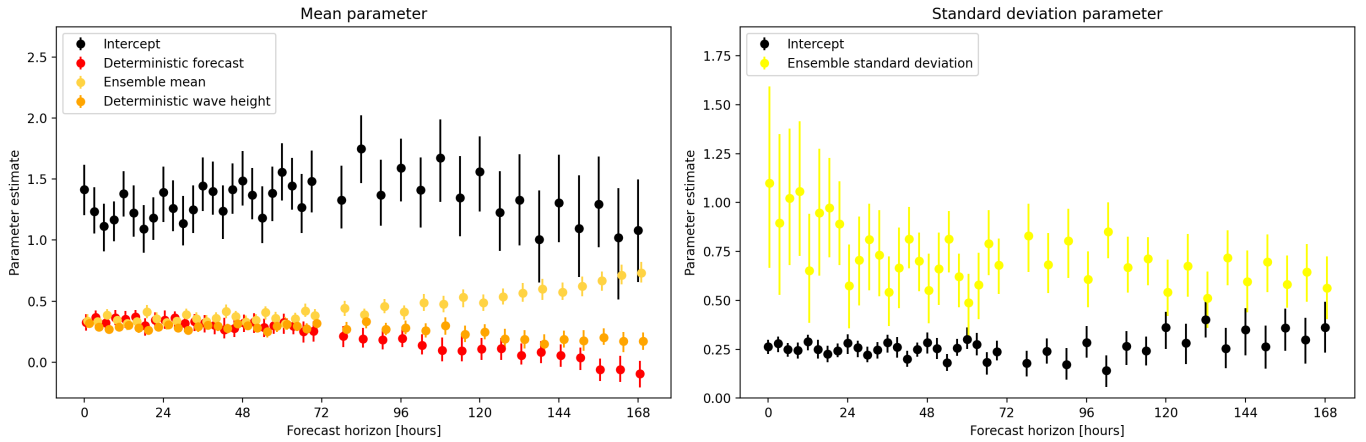




(a) Significant wave height  $H_S$ .



(b) Wind speed  $W$ .



(c) Mean wave period  $T_M$ .

Figure 9: Variation of estimated parameters with forecast horizon from consistent NHGR models for variables (a)  $H_S$ , (b)  $W$  and (c)  $T_M$ . Form of consistent NHGR calibration model given in Table 2. Central 95% uncertainty bands for parameter estimates are calculated using bootstrap resampling. Other details are given in Figure 6.

model performance is noticeable in the width of Gaussian forecast uncertainty bands (calculated using the estimated parameter values and error standard deviations, see Equations 1 and 2), particularly for  $H_S$ .

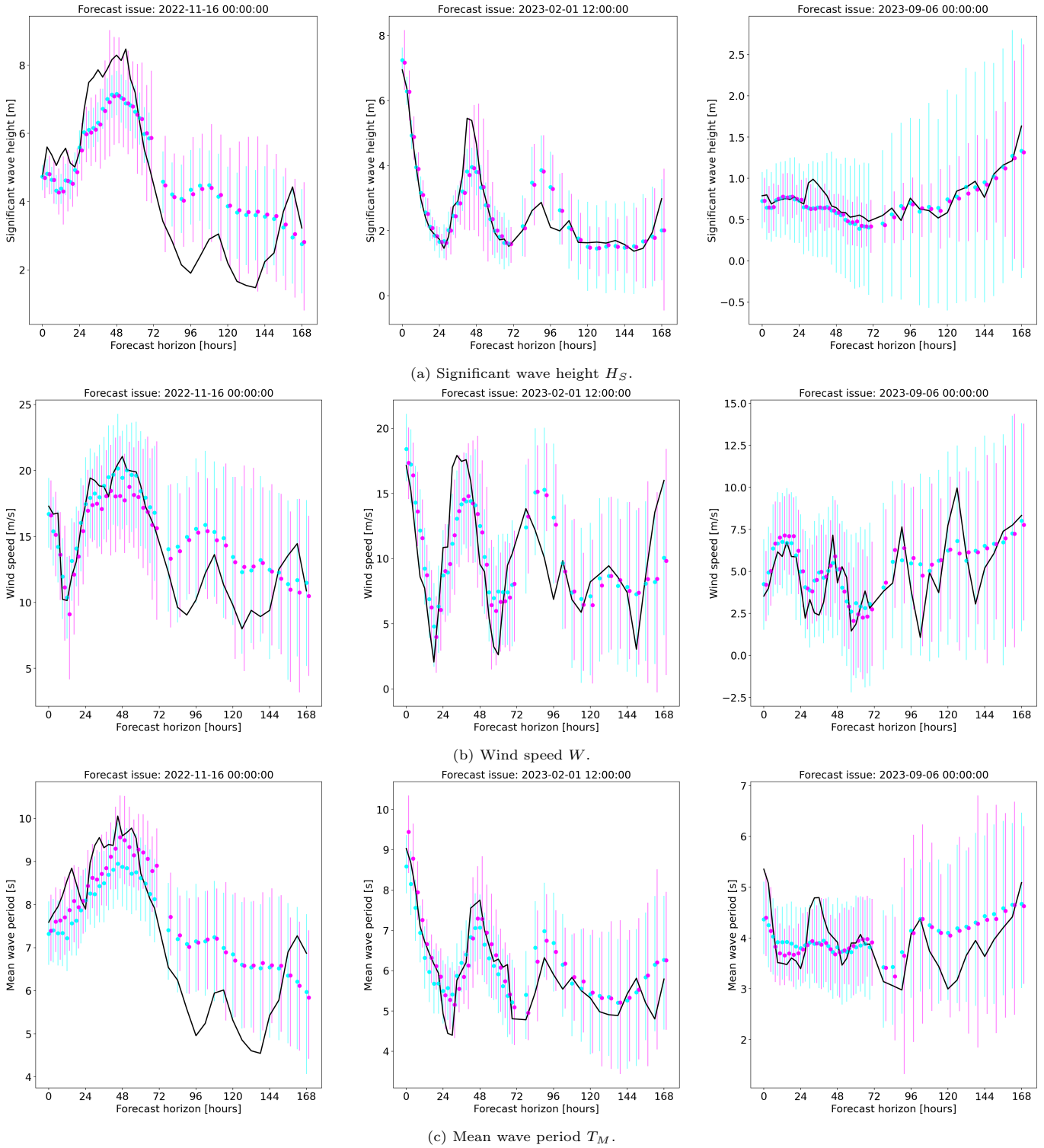
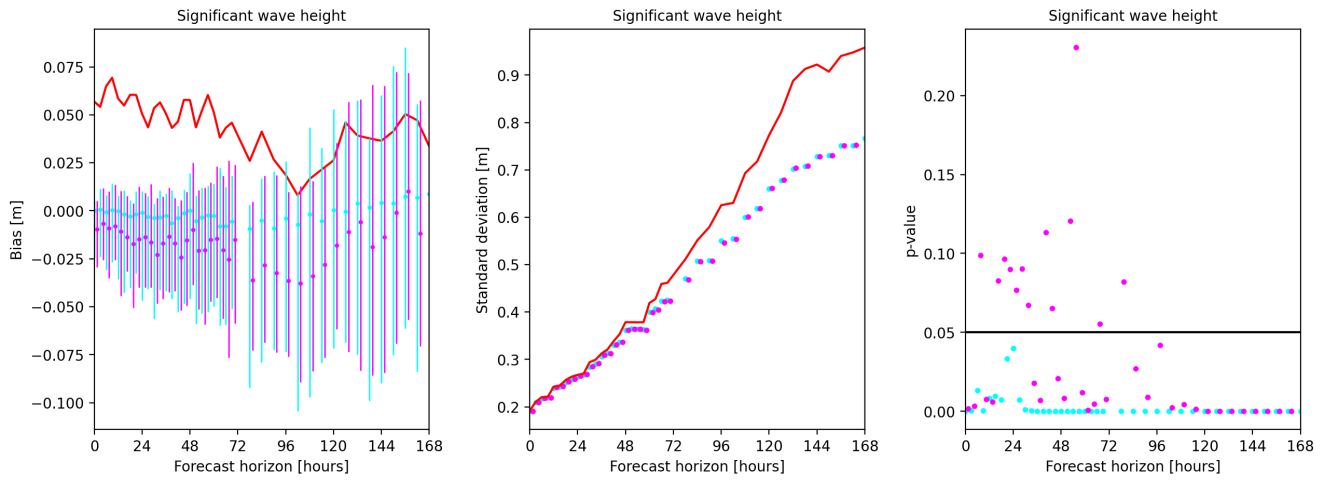
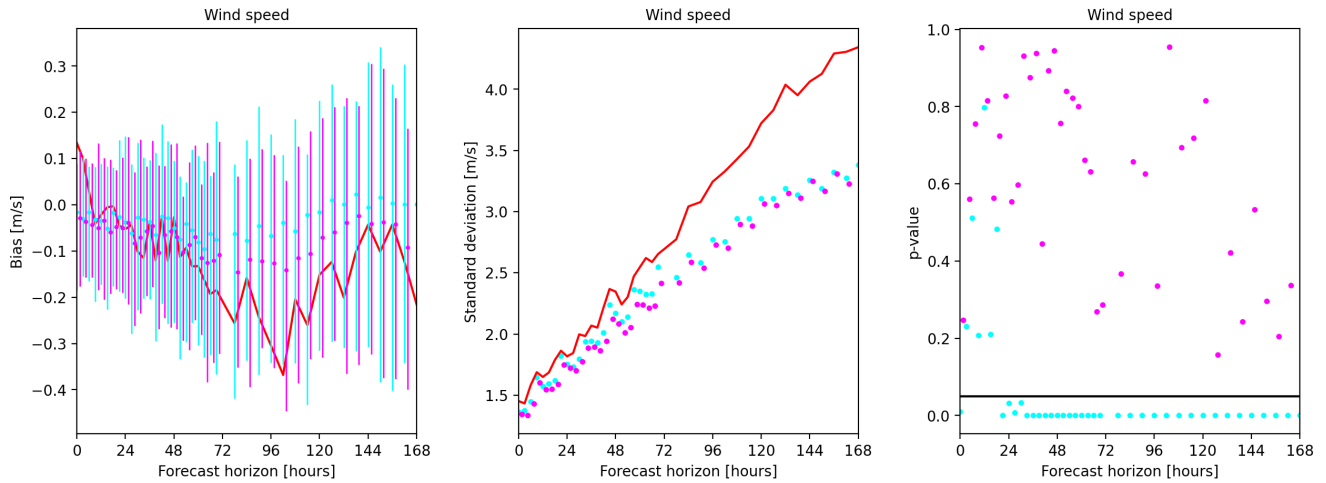


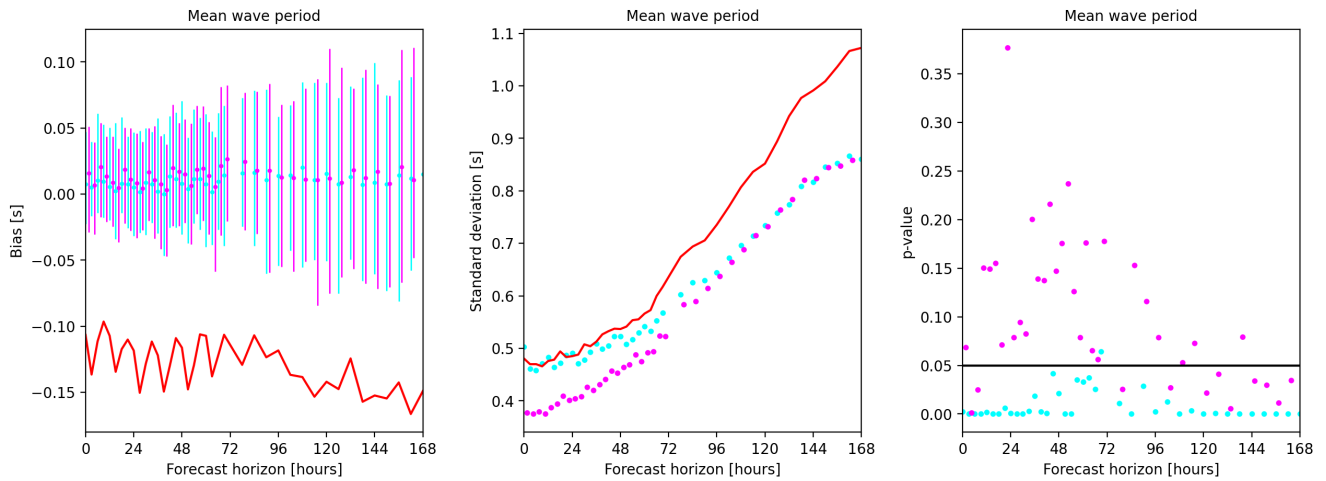
Figure 10: NHGR-calibrated forecasts and (future) reality for given variable at given forecast time. Three examples (columns) of forecasts on horizons  $\in [1, 168]$  for significant wave height (top), wind speed (middle) and mean wave period (bottom). Title of each column gives the time of forecast issue. (Future) reality illustrated using black line. Optimal calibrated forecast given in magenta, in terms of the mean (disc) and 95% Gaussian forecast uncertainty band (vertical line). Also shown for comparison are the corresponding optimal calibrated LR-forecasts (cyan) from Figure 7. Form of consistent LR and NHGR calibration models given in Tables 1 and 2. Box-whiskers have been translated horizontally by a small amount for clarity to avoid them being superimposed.



(a) Significant wave height  $H_S$ .



(b) Wind speed  $W$ .



(c) Mean wave period  $T_M$ .

Figure 11: Estimated mean (left) and standard deviation (centre) of forecast errors as a function of horizon, for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ . The right column gives p-values of a KS test, with null hypothesis that standardised forecast residuals are standard Gaussian distributed, as a function of horizon. Red lines refer to the (uncalibrated) deterministic forecast; cyan (magenta) to LR (NHGR) calibration models. Box-whiskers (left panels) and discs (centre and right panels) have been translated horizontally by a small amount for clarity to avoid them being superimposed. 95% uncertainty bands for bias (left) are calculated using bootstrap resampling.

#### 4.3. Comparing deterministic forecast, LR and NHGR forecast calibration model performance: within-sample assessment

Here we summarise the difference in performance of LR and NHGR calibration models by characterising the distribution of their forecast error as a function of horizon, specifically by estimating the bias and standard deviation of mean forecasts for the full period of data available. We also assess the similarity of the distribution of standardised residuals from the fitted models to the assumed Gaussian distribution, using a Kolmogorov-Smirnov (KS) test.

With  $Y(t+\tau)$  representing the in-situ measured metocean variable  $Y$  at time  $t+\tau$ , we seek to assess the performance of its calibrated forecasts  $\widehat{Y}_S(\tau|t)$  at forecast issue time  $t$  and forecast horizon  $\tau \geq 0$ , where  $S$  indicates the source of the forecast. Specifically, we will consider three cases  $S \in \{D, LR, NHGR\}$ , corresponding to the (uncalibrated) deterministic forecast, the LR calibration model and the NHGR calibration model. The deterministic forecast is included as a benchmark which with to assess the improvement in forecast performance offered by the LR and NHGR models. Each of the calibrated forecasts can be written in the form

$$\widehat{Y}_S(\tau|t) = \mu_S(\tau|t) + \epsilon \sigma_S(\tau|t), \text{ for } S \in \{D, LR, NHGR\} \quad (4)$$

since both LR and NHGR have a Gaussian error structure, where  $\mu_S(\tau|t)$  is the mean forecast and  $\sigma_S(\tau|t)$  the corresponding error standard deviation. See Equations 1 and 2 for the functional forms of these parameters for LR and NHGR; for deterministic forecast D,  $\mu_D(\tau|t)$  is simply the deterministic forecast itself, direct from the forecast provider, and  $\sigma_D(\tau|t) = 0$ . Further  $\epsilon$  a standard Gaussian random variate.

We then characterise model performance for D, LR and NHGR in terms of three summary statistics per forecast horizon, including sample estimates for the mean  $\mathbb{E}(Y(t+\tau) - \mu_S(\tau|t))$  and standard deviation  $\text{sd}(Y(t+\tau) - \mu_S(\tau|t))$  of the difference between reality and mean forecast. For LR and NHGR models, we also estimate the p-value associated with the KS test with null hypothesis that the standardised residuals  $(Y(t+\tau) - \mu_S(\tau|t))/\sigma_S(\tau|t)$  are drawn from a standard Gaussian distribution; p-values less than 0.05 indicate evidence that the null hypothesis can be rejected.

Results are given in Figure 11. In each row of the figure, the left hand panel illustrates the estimated bias  $\mathbb{E}(Y(t+\tau) - \mu_S(\tau|t))$  for LR (cyan) and NHGR (magenta) models, with the (uncalibrated) deterministic forecasts used as a benchmark. We also provide central 95% uncertainty intervals for bias estimates using a non-parametric bootstrap analysis. The centre panel illustrates the corresponding standard deviation  $\text{sd}(Y(t+\tau) - \mu_S(\tau|t))$ , again using the (uncalibrated) deterministic forecast standard deviation as a benchmark. The right hand panel illustrates p-values for the KS test for standardised residuals.

From the figure we see that the LR and NHGR calibration models provide reduced bias, especially for forecasting of  $T_m$ . We see considerable asymmetry in the bootstrap uncertainty intervals for bias of LR models for  $H_S$ , but note that the uncertainty intervals for bias of both LR and NHGR are very narrow relative to the typical range of values of  $H_S$ . We also see that forecast error standard deviation is smaller for LR and NHGR than for the uncalibrated deterministic forecast, particularly for longer forecast horizons. There is also evidence that the NHGR model outperforms LR for  $T_m$  for shorter horizons. Results of the KS test suggest that there is evidence to reject the null hypothesis of standard Gaussian standardised residuals from LR models (cyan) for all variables, and that standardised residuals from NHGR model fits (magenta) are more consistent with modelling assumptions, particularly for shorter horizons. Biases in forecasts from Figure 11 are unlikely to be of material concern to the metocean engineer, and may well be at the level of measurement error in practice. However, uncertainties in forecasts grow to levels which are likely to be of practical concern for horizons of 3 days and longer, for which 95% uncertainty bands are at least  $\pm 1$  m for  $H_S$ ,  $\pm 5$  m/s for  $W$  and  $\pm 1$  s for  $T_M$ .

#### 4.4. Comparing deterministic forecast, LR and NHGR forecast calibration model performance: out-of-sample assessment

Provided that data for the calibration model training period is representative of future observations of the environment, we can be confident that future model performance will be similar to that reported in Figure 11. We can also directly evaluate forecast performance for a time period following that used to estimate the calibration models. Figures SM4 and SM5 in the Supplementary Material illustrate forecast performance of models over out-of-sample Periods 1 (1 September 2023 - 30 April 2024) and Period 2 (1 September 2024 - 31 December 2024). Regrettably, for Period 1, comparable wave period data were not available, and for Period 2, comparable wind speed data were not available. (More specifically, for Period 1, for reasons beyond the authors' control, the ensemble forecast data from the forecast provider were not retained and hence not available for analysis. Further, just prior to Period 2, a recalibration of the wind sensor was performed, making fair assessment of forecast performance unviable.) Nevertheless, for  $H_S$  and  $W$  (Period 1) and  $H_S$  and  $T_m$  (Period 2), we see that the general characteristics of Figures SM4 and SM5 are similar to those of Figure 11. Forecast bias is small relative to standard deviation for the uncalibrated deterministic, LR-calibrated and NHGR-calibrated forecasts. Reduction in forecast standard deviation at longer horizons is clear

for both LR and NHGR relative to the deterministic forecast. The distribution of residuals from the NHGR model is generally more similar to the assumed Gaussian form.

Bias and variance are fundamental quantities used to characterise the performance of an estimator, favoured by us in the current work. Similarly, the KS statistic is a generic measure of the dissimilarity between distributions. However there are additional performance scoring rules and diagnostics, particularly interesting when evaluating ensemble forecasts, which could also be used for the current work. These include the continuous ranked probability score (CRPS; e.g. Gneiting et al. 2005) and the probability integral transform (PIT) histogram (e.g. Dawid 1984). More generally, the work of Hernandez et al. (2018) provides a review of performance, skill and accuracy assessment in operational oceanography, and Messner et al. (2020) reviews forecast verification tools, with a focus on wind power applications.

As an illustration, Figure 12 shows the variation of CRPS with forecast horizon for the two out-of-sample test Periods 1 and 2. For value  $y$  of the response at forecast horizon  $\tau$ , CRPS is calculated from the distribution  $F_{\hat{Y}}(\hat{y}; \tau)$  of the

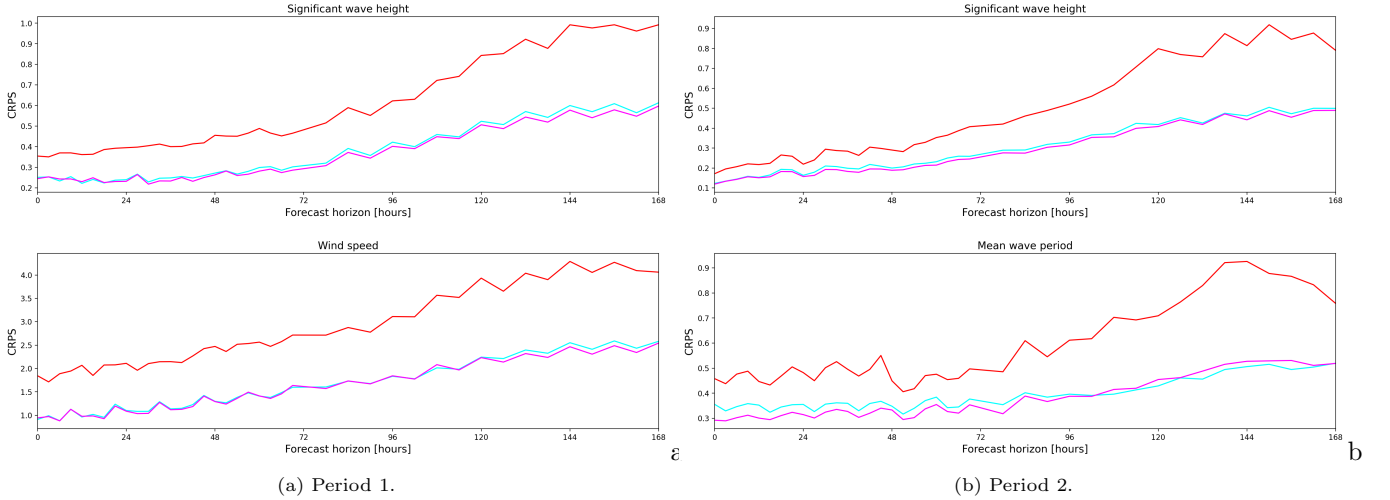


Figure 12: Variation of CRPS with forecast horizon for (a) out-of-sample Period 1 (1 September 2023 - 30 April 2024) and (b) out-of-sample Period 2 (1 September 2024 - 31 December 2024) for the uncalibrated deterministic forecast (red), the forecast calibrated using LR (cyan) and the forecast calibrated using NHGR (magenta).

forecast response  $\hat{Y}$ , using the expression

$$\text{CRPS}(y, \tau) = \int_{\hat{y}} (F_{\hat{Y}}(\hat{y}; \tau) - I(\hat{y} \geq y))^2 d\hat{y}$$

for indicator function  $I$  (with  $I(x) = 1$  if  $x$  is true and  $= 0$  otherwise), which can then be averaged over  $y$  to provide a function of  $\tau$  only (e.g. Gneiting et al. 2005), or inspected as a function of  $y$  for given  $\tau$ .  $\text{CRPS} = 0$  indicates perfect agreement between forecast and truth. Figure 12 shows that the average CRPS performance of calibrated forecasts better than that of the uncalibrated deterministic forecast, with an approximate reduction of 50% in value. Moreover, the NHGR calibrated forecast (magenta) is generally (but not always) slightly better than that calibrated using LR (cyan). Further, Figures SM5 and SM6 of the Supplementary Material show variation of CRPS with the value of response for four representative forecast horizons, for Periods 1 and 2.

Further, Figure 13 shows corresponding rank histograms which provide a visual interpretation of the distribution of observed values of response relative to the distribution of probabilistic forecast. As shown by Equation 4, both LR- and NHGR-calibration models yield a probabilistic forecast with Gaussian error structure. Therefore, the standardised residual  $r(\tau|t)$  for observation  $y(t + \tau)$  given the forecast at time  $t$  and forecast horizon  $\tau$ , defined by

$$r(\tau|t) = \frac{y(t + \tau) - \mu_S(\tau|t)}{\sigma_S(\tau|t)}$$

is expected to follow the standard Gaussian distribution with zero mean and unit variance (and cumulative distribution function  $\Phi(\bullet)$ ), with  $\mu_S$  and  $\sigma_S$  defined in Equation 4. Therefore, we expect the cumulative distribution function  $\Phi(r(\tau|t))$  evaluated at  $r(\tau|t)$  to be a uniform random number on  $[0, 1]$ . We can hence inspect the empirical density of  $\Phi(r(\tau|t))$  over all  $t$  and  $\tau$ , as shown in Figure 13, to assess its correspondence to a constant uniform density on  $[0, 1]$ . The figure suggests that there is reasonable overall agreement between residuals from the LR- and NHGR-calibrated

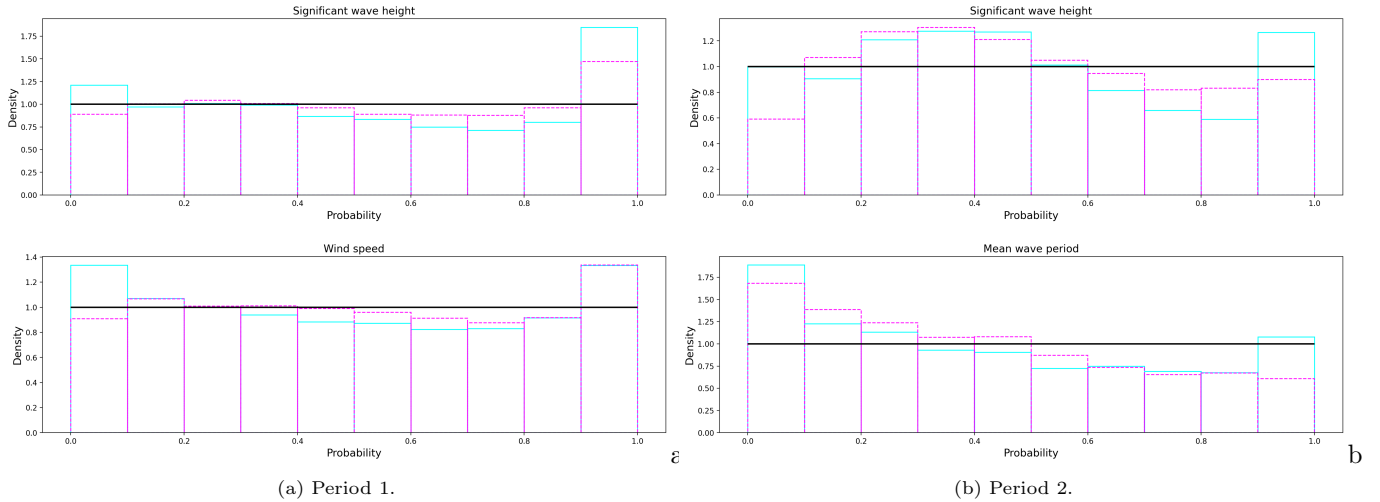


Figure 13: Rank histogram plots over all forecast horizons for (a) out-of-sample Period 1 (1 September 2023 - 30 April 2024) and (b) out-of-sample Period 2 (1 September 2024 - 31 December 2024) for the the forecast calibrated using LR (cyan) and the forecast calibrated using NHGR (magenta). The horizontal line (black) indicates the model-assumed histogram in each case.

forecasts and the model-assumed Gaussian distributional forms; departures from the horizontal black line in each panel indicate lack of agreement between model and reality. It appears that the NHGR-calibrated forecast provides somewhat better general agreement. We note that results of the KS testing visualised in Figures SM4 and SM5 for the out-of-sample periods, provide rather similar diagnostic information.

## 5. Discussion and conclusions

In this article, we assess the value of calibrating forecast models for significant wave height  $H_S$ , wind speed  $W$  and mean spectral wave period  $T_m$  for forecast horizons between zero and 168 hours from a commercial forecast provider, to improve forecast performance for a location in the central North Sea. We consider two straightforward calibration models, linear regression (LR) and non-homogeneous Gaussian regression (NHGR), incorporating deterministic, control and ensemble mean forecast covariates. We show that relatively simple calibration models (with at most three covariates) provide good calibration; addition of further covariates cannot be justified. Optimal calibration models (for the forecast mean of a physical quantity) always make use of the deterministic forecast and ensemble mean forecast for the same quantity, together with a covariate associated with a different physical quantity. The selection of optimal covariates is performed independently per forecast horizon, and the set of optimal covariates shows a large degree of consistency across forecast horizons. As a result, it is possible to specify a consistent model to calibrate a given physical quantity, incorporating a common set of three covariates for all horizons. For NHGR models of a given physical quantity, the ensemble forecast standard deviation for that quantity is skilful in predicting forecast error standard deviation, especially for  $H_S$ . We show that the consistent LR and NHGR calibration models facilitate reduction in forecast bias to near zero for all of  $H_S$ ,  $W$  and  $T_m$ , and that there is little difference between LR and NHGR calibration for the mean. Both LR and NHGR models facilitate reduction in forecast error standard deviation relative to naive adoption of the (uncalibrated) deterministic forecast, with NHGR providing somewhat better performance. Distributions of standardised residuals from NHGR models are generally considerably more like standard Gaussians.

Direct adoption of (uncalibrated) ensemble forecasts is not recommended, since there is evidence of large bias for  $T_m$ . For short horizons, the contribution of the deterministic forecast to the calibration model is highest, decaying with increasing horizon. In contrast, the importance of the ensemble mean forecast (for the same metocean quantity as that being forecast) increases. The relative contributions of deterministic and ensemble mean forecasts to the calibration varies across metocean responses. During the course of the study, we also considered different summaries of the ensemble distribution to the mean and standard deviation. We found e.g. that use of the ensemble median offered no improvement in general over the ensemble mean.

The choices of calibration models used here represent the simplest approaches that might reasonably be adopted in practice. Consequently there are many opportunities to extend the analysis. We noted evidence in the exploratory analysis that the forecast model generally tends to underestimate the very largest values. This might be an opportunity e.g. to include quadratic and higher order terms in covariates in the parametric form for the forecast mean (and, within the NHGR framework, for the forecast standard deviation). Alternatively, given that joint largest values of measured

and forecast variables can be considered extreme, it might be more appropriate to adopt extreme value models (e.g. Davison and Smith 1990, Heffernan and Tawn 2004, Jonathan et al. 2014, Towe et al. 2023, Towe et al. 2024) to characterise these regions more correctly, or more generally to relax the assumption of Gaussianity made by both LR and NHGR. There are opportunities also to exploit the extreme time-series structure of predictors and responses, following the Markov extremal model and related frameworks in Winter and Tawn (2016), Tendijck et al. (2019) and Tendijck et al. (2024). We might also consider joint calibration of multiple metocean variables. Multivariate predictive modelling is a large field of research and applications, offering a wealth of modelling strategies to predict a multivariate response from multivariate predictors. The presence of correlated predictors can lead to inflation in estimated model parameters (which can be quantified using measures such as the variance inflation factor), and inflation of prediction uncertainty. Nevertheless, joint modelling of multiple metocean variables provides the potential for better calibrated forecasts, including extremes. In the context of weather forecasting, Allen et al. (2024) provides a discussion of methods for assessing the calibration of multivariate probabilistic forecasts. Extension to calibration for multiple locations is also possible; for some locations at least, the calibration model might be sensitive to additional covariates, including e.g. wave and wind direction and season.

NHGR models are relatively simple to estimate. Although closed form uncertainty quantification (e.g. for parameter estimates and predictions) is not available, we find that non-parametric bootstrapping can be used successfully. Based on these findings, we recommend the adoption of NHGR models incorporating ensemble forecasts for medium-term forecasting tasks of the type described here, including for unmanning and related operations (see e.g. Towe et al. 2021).

## Acknowledgements

The authors would like to thank StormGeo for provision of the forecast data, and for useful discussions during the course of the analysis and preparation of this article. They further acknowledge the support of Shell colleagues Graham Feld, David Randell and Stan Tendijck, and helpful comments from two reviewers.

## References

- Adnan, R.M., Sadeghifar, T., Alizamir, M., Azad, M.T., Makarynsky, O., Kisi, O., Barati, R., Ahmed, K.O., 2023. Short-term probabilistic prediction of significant wave height using Bayesian model averaging: Case study of Chabahar port, Iran. *Ocean Engineering* 272, 113887.
- Akaike, H., 1974. A new look at the statistical model identification. *IEEE Trans. Autom. Control* 19, 716–723.
- Allen, S., Ginsbourger, D., Ziegel, J., 2023. Evaluating forecasts for high-impact events using transformed kernel scores. *SIAM-ASA J. Uncertain.* 11, 906–940.
- Allen, S., Ziegel, J., Ginsbourger, D., 2024. Assessing the calibration of multivariate probabilistic forecasts. *Q. J. R. Meteorol. Soc.* 150, 1315–1335.
- Astfalck, L., Bertolacci, M., Cripps, E., 2023. Evaluating probabilistic forecasts for maritime engineering operations. *Data-Centric Eng.* 4, e15.
- Bessa, R.J., Moehrlen, C., Fundel, V., Siefert, M., Browell, J., Haglund El Gaidi, S., Hodge, B., Cali, U., Kariniotakis, G., 2017. Towards improved understanding of the applicability of uncertainty forecasts in the electric power industry. *Energies* 10.
- Bjerregard, M.B., Moeller, J.K., Madsen, H., 2021. An introduction to multivariate probabilistic forecast evaluation. *Energy and AI* 4, 100058.
- Cerqueira, V., Torgo, L., 2024. Exceedance probability forecasting via regression for significant wave height prediction. URL: <https://arxiv.org/abs/2206.09821>, arXiv:2206.09821.
- Davison, A., Smith, R.L., 1990. Models for exceedances over high thresholds. *J. Roy. Statist. Soc. B* 52, 393.
- Dawid, A.P., 1984. Present position and potential developments: Some personal views: Statistical theory: The prequential approach. *J. Roy. Statist. Soc. A* 147, 278–292.
- Emiliano, P.C., Vivanco, M.J., de Menezes, F.S., 2014. Information criteria: how do they behave in different models? *Comput. Stat. Data Anal.* 69, 141–153.

- Gao, P., Director, H., Bitz, C., Raftery, A., 2022. Probabilistic forecasts of Arctic Sea ice thickness. *JABES* 27, 280–302.
- Gilbert, C., Browell, J., McMillan, D., 2021. Probabilistic access forecasting for improved offshore operations. *Int. J. Forecast.* 37, 134–150.
- Gneiting, T., 2011. Making and evaluating point forecasts. *J. Am. Statist. Soc.* 106, 746–762.
- Gneiting, T., 2014. Calibration of medium-range weather forecasts. ECMWF Technical Memoranda doi:10.21957/8xna7g1ta.
- Gneiting, T., Balabdaoui, F., Raftery, A.E., 2007. Probabilistic forecasts, calibration and sharpness. *J. Roy. Statist. Soc. B* 69, 243–268.
- Gneiting, T., Katzfuss, M., 2014. Probabilistic forecasting. *Annu. Rev. Stat. Appl.* , 125–151.
- Gneiting, T., Raftery, A.E., Westveld, A.H., Goldman, T., 2005. Calibrated probabilistic forecasting using ensemble model output statistics and minimum CRPS estimation. *Mon. Weather Rev.* 133, 1098 – 1118.
- Heffernan, J.E., Tawn, J.A., 2004. A conditional approach for multivariate extreme values. *J. Roy. Statist. Soc. B* 66, 497–546.
- Heinrich, C., Hellton, K.H., Lenkoski, A., Thorarinsdottir, T.L., 2021. Multivariate postprocessing methods for high-dimensional seasonal weather forecasts. *J. Am. Statist. Soc.* 116, 1048–1059.
- Hernandez, F., Smith, G., Baetens, K., Cossarini, G., Garcia-Hermosa, I., Drevillon, M., Maksymczuk, J., Melet, A., Regnier, C., von Schuckmann, K., 2018. *New Frontiers in Operational Oceanography*. Editors E. Chassignet, A. Pascual, J. Tintore and J. Verron. GODAE OceanView. chapter 29: Measuring performances, skill and accuracy in operational oceanography: new challenges and approaches. pp. 759–796.
- Hoehlein, K., Schulz, B., Westermann, R., Lerch, S., 2024. Postprocessing of ensemble weather forecasts using permutation-invariant neural networks. *Artif. Intell. Earth Syst.* 3, e230070.
- Jonathan, P., Randell, D., Wu, Y., Ewans, K., 2014. Return level estimation from non-stationary spatial data exhibiting multidimensional covariate effects. *Ocean Eng.* 88, 520–532.
- Messner, J.W., Pinson, P., Browell, J., Bjerregard, M.B., Schicker, I., 2020. Evaluation of wind power forecasts: an up-to-date view. *Wind Energy* 23, 1461–1481.
- Pinson, P., Madsen, H., Nielsen, H.A., Papaefthymiou, G., Kleockl, B., 2009. From probabilistic forecasts to statistical scenarios of short-term wind power production. *Wind Energy* 12, 51–62.
- Schefzik, R., Thorarinsdottir, T.L., Gneiting, T., 2013. Uncertainty quantification in complex simulation models using ensemble copula coupling. *Stat. Sci.* 28, 616–640.
- Sweeney, C., Bessa, R.J., Browell, J., Pinson, P., 2020. The future of forecasting for renewable energy. *Wiley Interdiscip. Rev. Energy Environ.* 9, e365.
- Tendijck, S., Jonathan, P., Randell, D., Tawn, J.A., 2024. Temporal evolution of the extreme excursions of multivariate  $k$ th order Markov processes with application to oceanographic data. *Environmetrics* 35, e2834.
- Tendijck, S., Ross, E., Randell, D., Jonathan, P., 2019. A non-stationary statistical model for the evolution of extreme storm events. *Environmetrics* 30, e2541.
- Towe, R., Randell, D., Kensler, J., Feld, G., Jonathan, P., 2023. Estimation of associated values from conditional extreme value models. *Ocean Eng.* 272, 113808.
- Towe, R., Ross, E., Randell, D., Jonathan, P., 2024. covXtreme: MATLAB software for non-stationary penalised piecewise constant marginal and conditional extreme value models. *Environ. Model. Softw.* 177, 106035.
- Towe, R., Zanini, E., Randell, D., Feld, G., Jonathan, P., 2021. Efficient estimation of distributional properties of extreme seas from a hierarchical description applied to calculation of un-manning and other weather-related operational windows. *Ocean Eng.* 238, 109642.



- Tyralis, H., Papacharalampous, G., 2024. A review of predictive uncertainty estimation with machine learning. *Artif. Intell. Rev.* 57, 94.
- van der Meer, D., 2021. A benchmark for multivariate probabilistic solar irradiance forecasts. *Solar Energy* 225, 286–296.
- Winkler, R.L., 1969. Scoring rules and the evaluation of probability assessors. *J. Am. Statist. Soc.* 64, 1073–1078.
- Winter, H.C., Tawn, J.A., 2016. Modelling heatwaves in central France: a case-study in extremal dependence. *J. Roy. Statist. Soc. C* 65, 345–365.

# Supplementary Material

## Calibration of medium-range metocean forecasts for the North Sea

Conor Murphy, Ross Towe and Philip Jonathan

### Overview

This document provides supporting plots for the article “Calibration of medium-range metocean forecasts for the North Sea”.

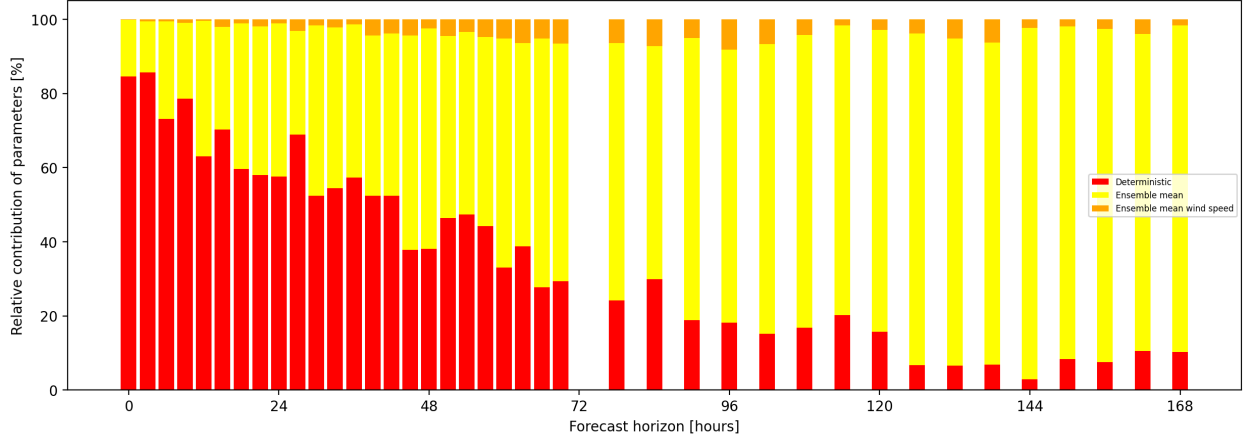
Figures SM1 and SM3 illustrate the contributions of different covariates to linear regression (LR, Figure SM1) and non-homogeneous Gaussian regression (NHGR, Figure SM3). These are calculated using estimates for LR and NHGR parameters  $\beta_j(\tau)$  (see Equations 1 and 2 of main text) from model fits where all covariates have been standardised to zero mean and unit variance, as

$$\frac{|\beta_j(\tau)|}{\sum_k |\beta_k(\tau)|}.$$

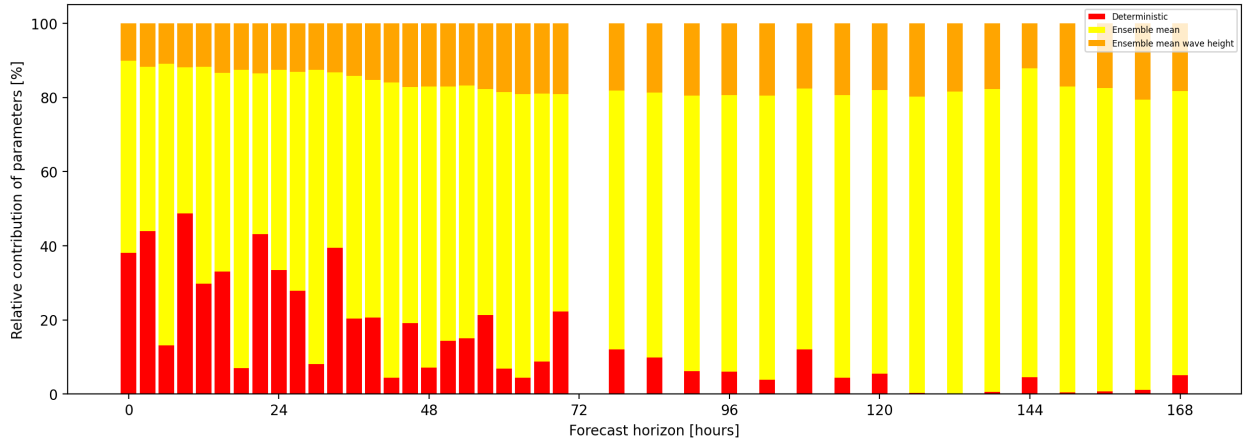
Here, the set of parameters  $\beta(\tau)$  corresponds to the set of the parameters  $b$  and  $c$  in the expressions LR or NHGR mean, and to the set of parameters  $d$  and  $e$  in the expression for the NHGR standard deviation. The intercept term  $a$  in LR and NHGR mean is ignored without loss of generality, but included for the NHGR standard deviation. Figures SM1 and SM3 provide illustrations to support the discussion on Figures 6 and 9 of the main text.

Figure SM2 provides the equivalent of Figure 4 of the main text for NHGR calibration.

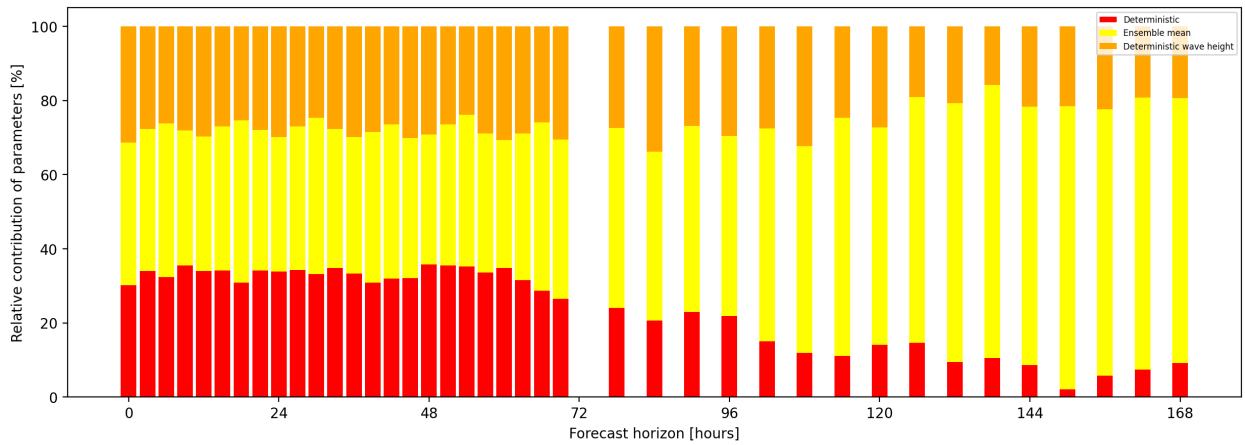
Figures SM4 and SM5 summarise out-of-sample forecast performance for the uncalibrated deterministic forecast, and the LR- and NHGR-calibrated forecasts, on two out-of-sample periods. Figures SM6 and SM7 illustrate corresponding continuous rank probability scores, and Figures SM8 and SM9 rank histograms.



(a) Significant wave height  $H_S$ .

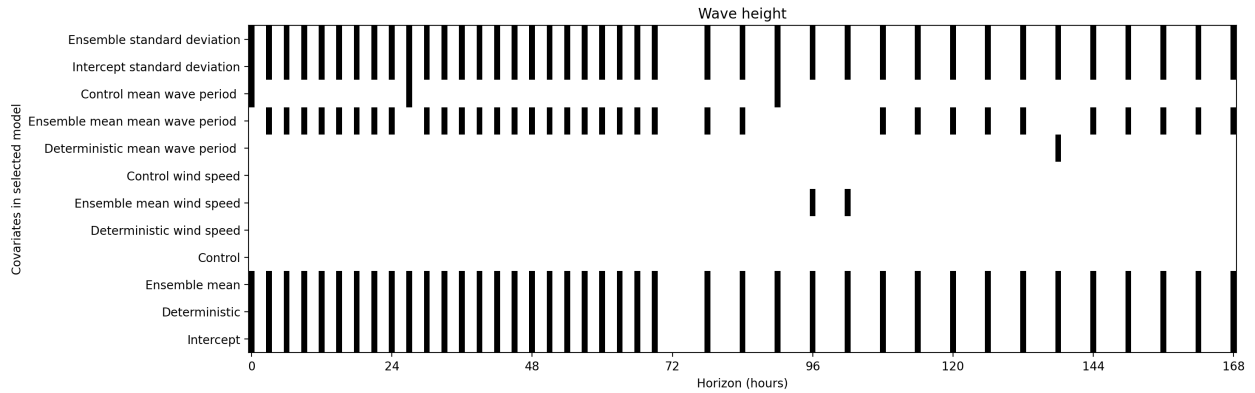


(b) Wind speed  $W$ .

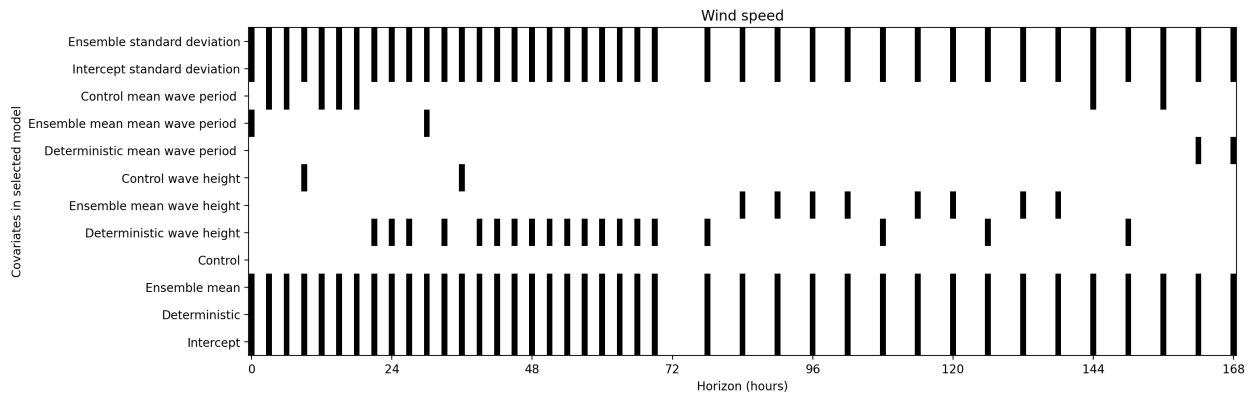


(c) Mean wave period  $T_M$ .

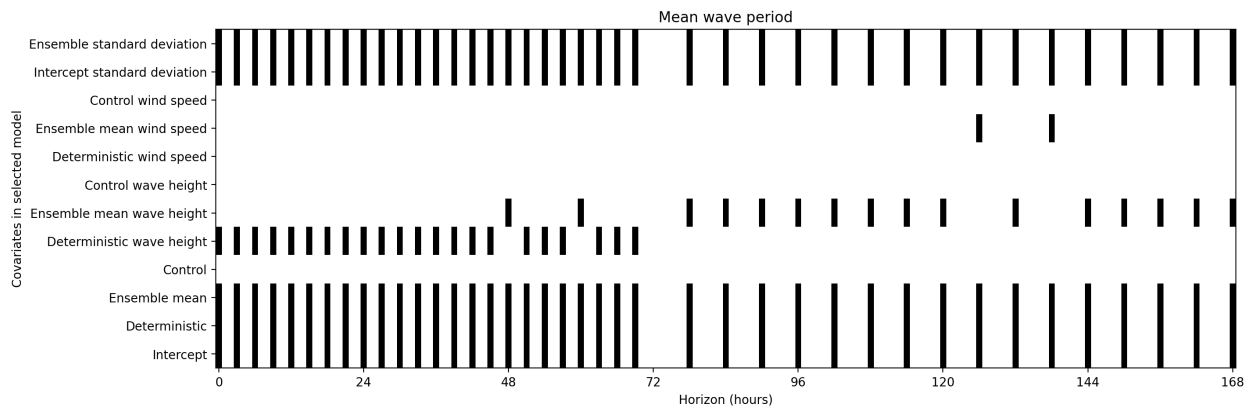
Figure SM1: Plot of linear regression covariate contributions to the calibration of forecasts for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ , as a function of forecast horizon. All covariates are standardised to zero mean and unit standard deviation prior to regression analysis. Contributions are reported in terms of the percentage absolute value of the parameter coefficient from the regression. The intercept effect is not shown. Consistent model forms are given in Table 1 of the main text.



(a) Significant wave height  $H_S$ .

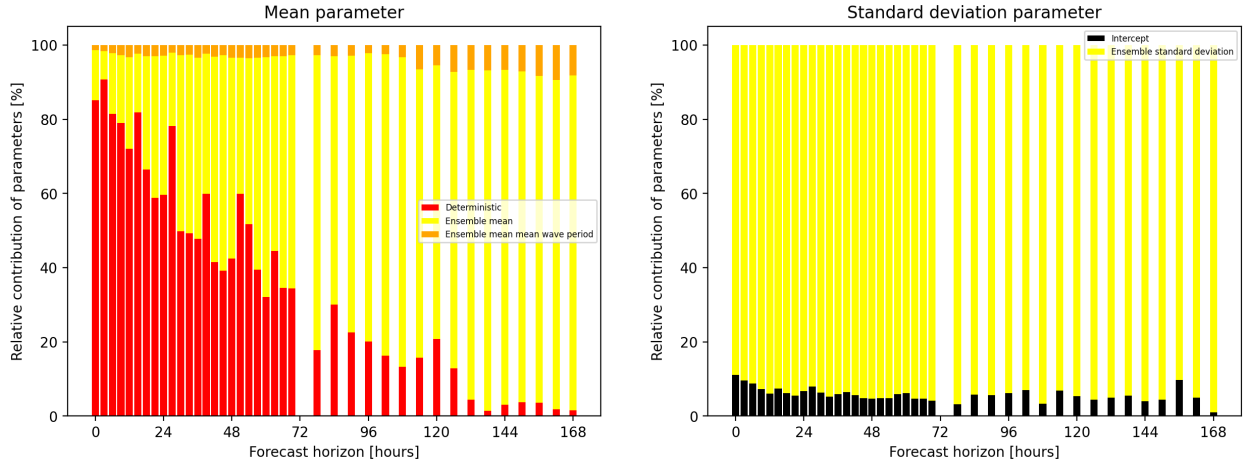


(b) Wind speed  $W$ .

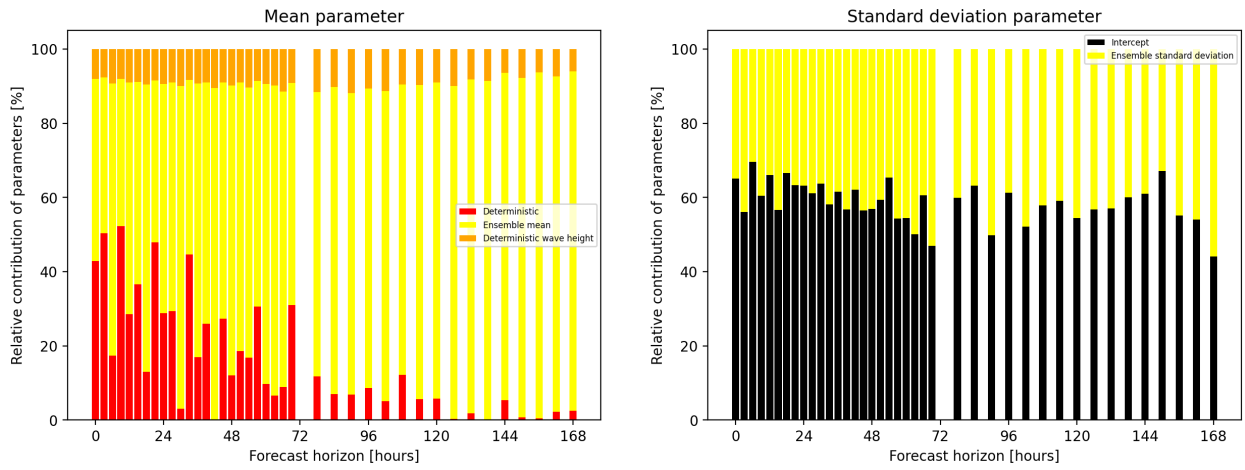


(c) Mean wave period  $T_M$ .

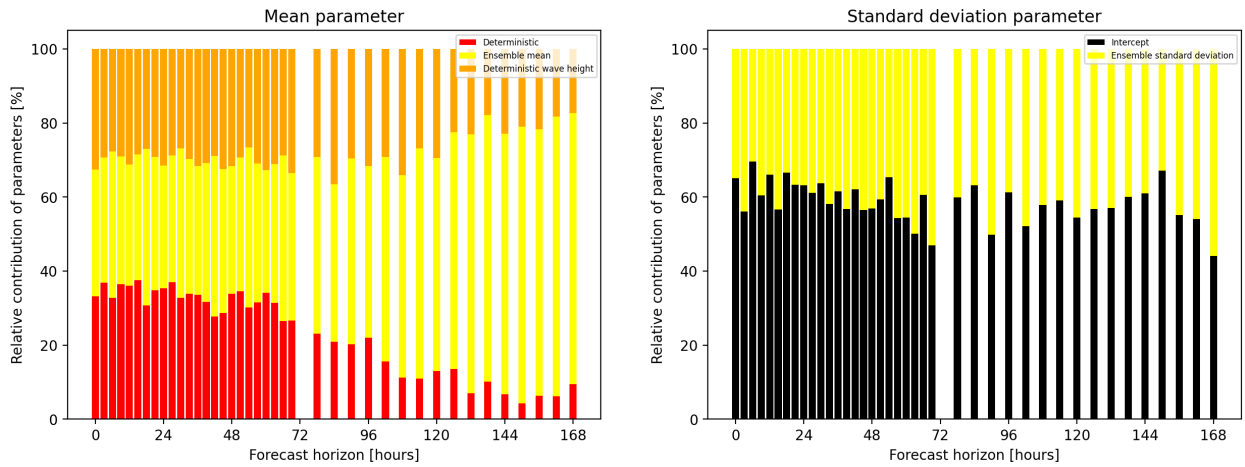
Figure SM2: Bar code plot of covariates included in optimal NHGR model for each forecast horizon. For discussion, see Section 4.2 of the main text.



(a) Significant wave height  $H_S$ .



(b) Wind speed  $W$ .



(c) Mean wave period  $T_M$ .

Figure SM3: Bar code plot of non-homogeneous Gaussian regression covariate contributions to the calibration of forecasts for (a) significant wave height  $H_S$ , (b) wind speed  $W$  and (c) mean wave period  $T_M$ , as a function of forecast horizon. The response and all covariates are standardised to zero mean and unit standard deviation prior to regression analysis. Contributions are reported in terms of the percentage absolute value of the parameter coefficient from the regression. Consistent model forms are given in Table 2 of the main text.

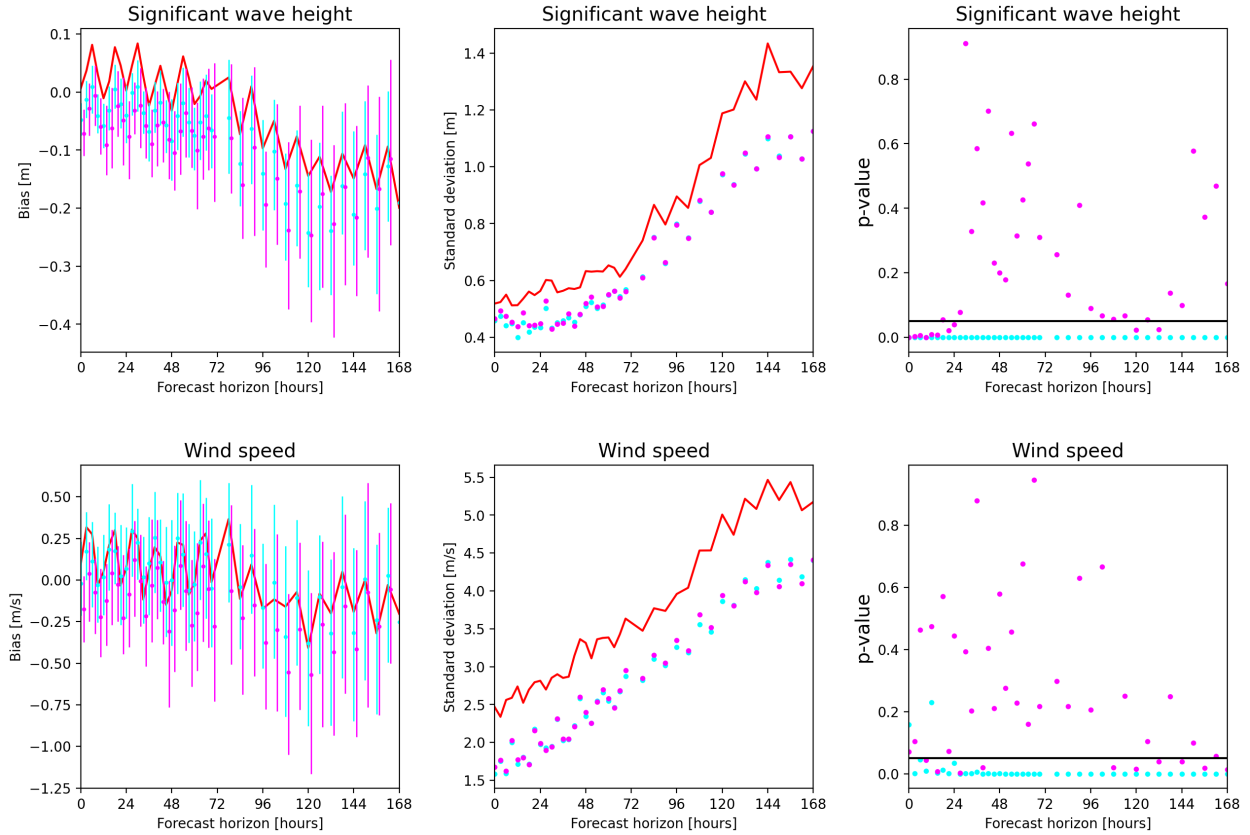


Figure SM4: Out-of-sample performance assessed for Period 1 (1 September 2023 - 30 April 2024), for calibration models estimated on data for period 17 May 2022 to 6 September 2023. Estimated mean (left) and standard deviation (centre) of forecast errors as a function of horizon, for (a) significant wave height  $H_S$  and (b) wind speed  $W$ . The right column gives p-values of a KS test, with null hypothesis that standardised forecast residuals are standard Gaussian distributed, as a function of horizon. Red lines refer to the (uncalibrated) deterministic forecast; cyan (magenta) to LR (NHGR) calibration models. Box-whiskers (left panels) and discs (centre and right panels) have been translated horizontally by a small amount for clarity to avoid them being superimposed. 95% uncertainty bands for bias (left) are calculated using bootstrap resampling. Results should be compared with those in Figure 11 of the main text. Bias is small relative to standard deviation. Reduction in forecast standard deviation at longer horizons is clear for both LR and NHGR relative to the deterministic forecast. The distribution of residuals from the NHGR model is generally more similar to the assumed Gaussian. Regrettably, comparable mean wave period data for the test period were not available.

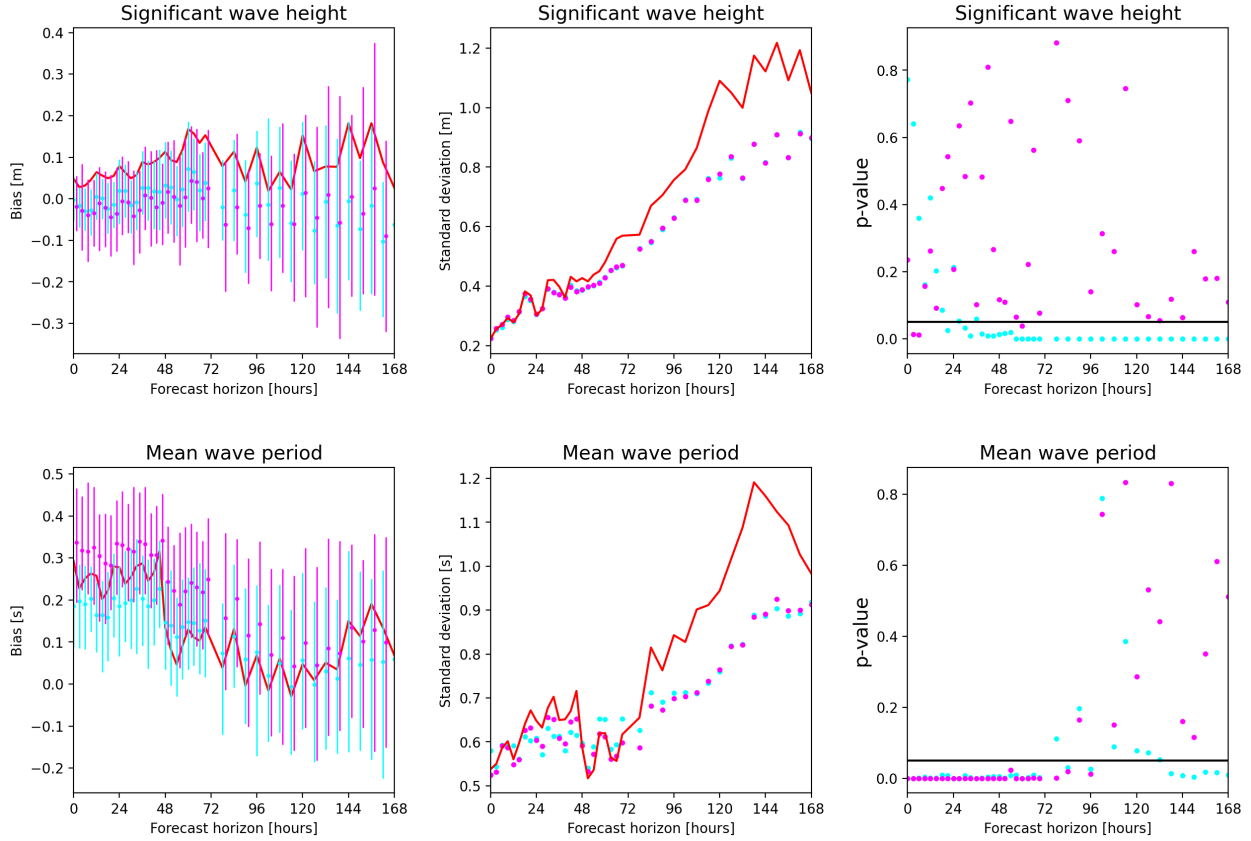
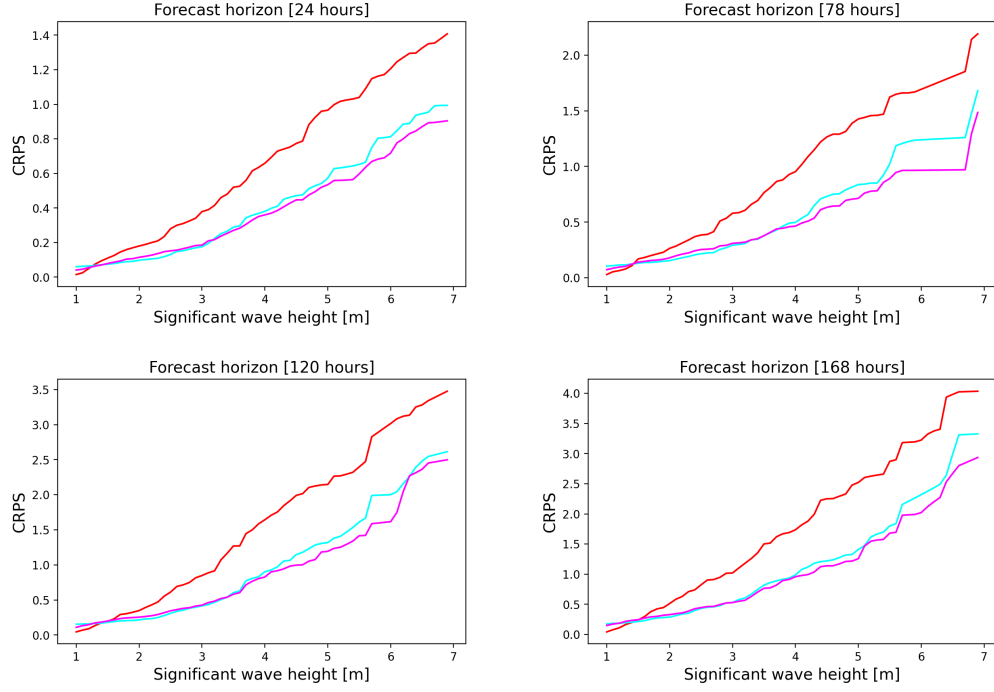
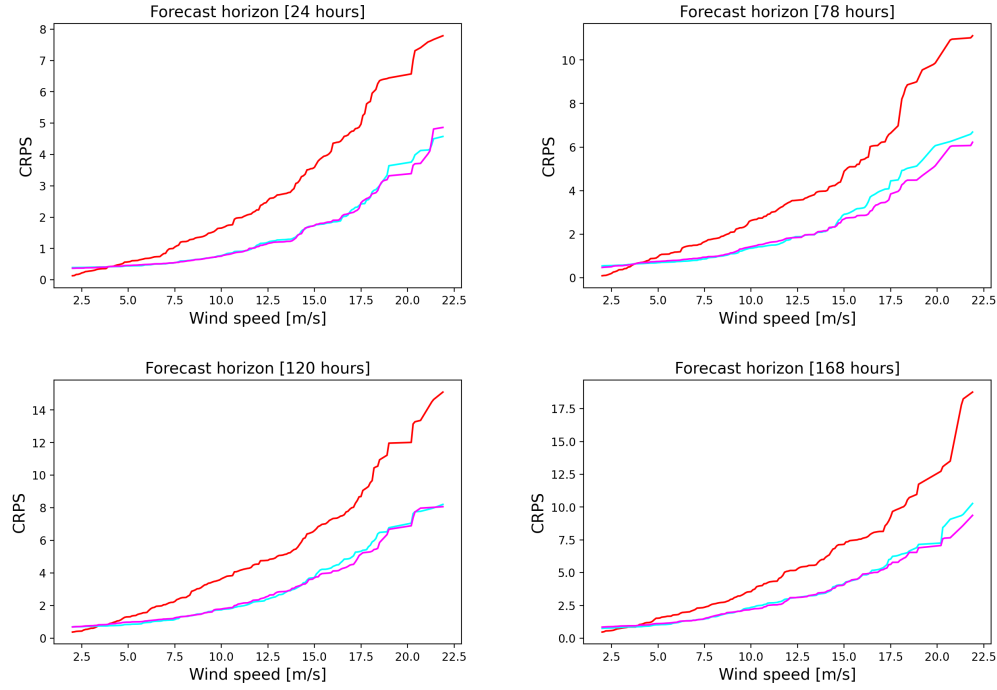


Figure SM5: Out-of-sample performance assessed for Period 2 (1 September 2024 – 31 December 2024, for calibration models estimated on data for period 17 May 2022 to 6 September 2023). Estimated mean (left) and standard deviation (centre) of forecast errors as a function of horizon, for (a) significant wave height  $H_S$  and (b) mean wave period  $T_M$ . The right column gives p-values of a KS test, with null hypothesis that standardised forecast residuals are standard Gaussian distributed, as a function of horizon. Red lines refer to the (uncalibrated) deterministic forecast; cyan (magenta) to LR (NHGR) calibration models. Box-whiskers (left panels) and discs (centre and right panels) have been translated horizontally by a small amount for clarity to avoid them being superimposed. 95% uncertainty bands for bias (left) are calculated using bootstrap resampling. Results should be compared with those in Figure 11 of the main text. Bias is small relative to standard deviation. Reduction in forecast standard deviation at longer horizons is clear for both LR and NHGR relative to the deterministic forecast. The distribution of residuals from the NHGR model is generally more similar to the assumed Gaussian for  $H_S$ . Regrettably, comparable wind speed data for the test period were not available.



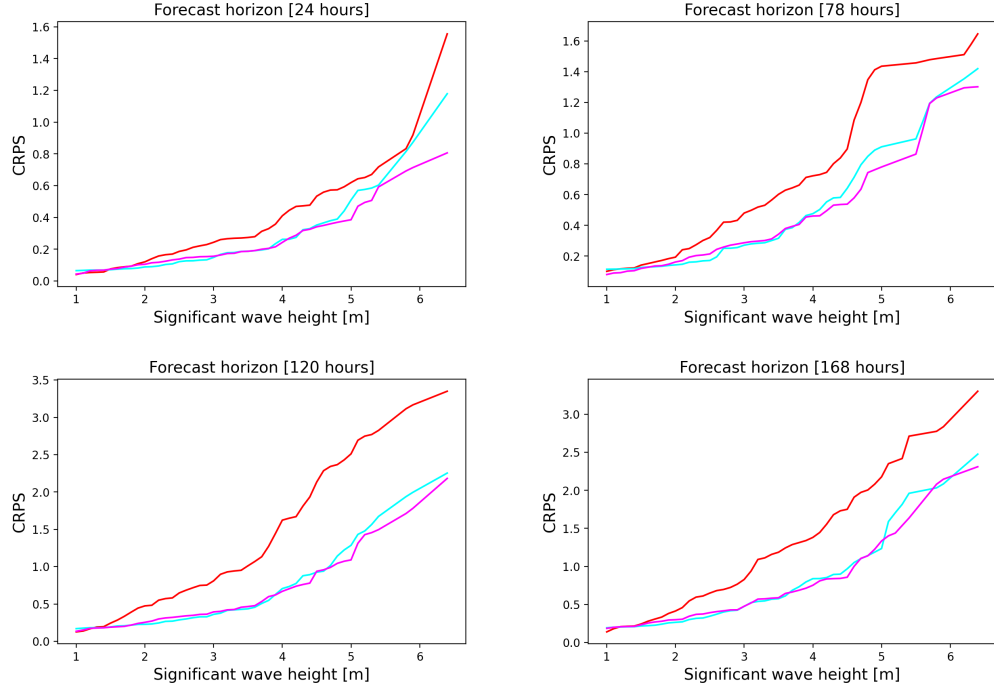
(a) Significant wave height  $H_S$ .



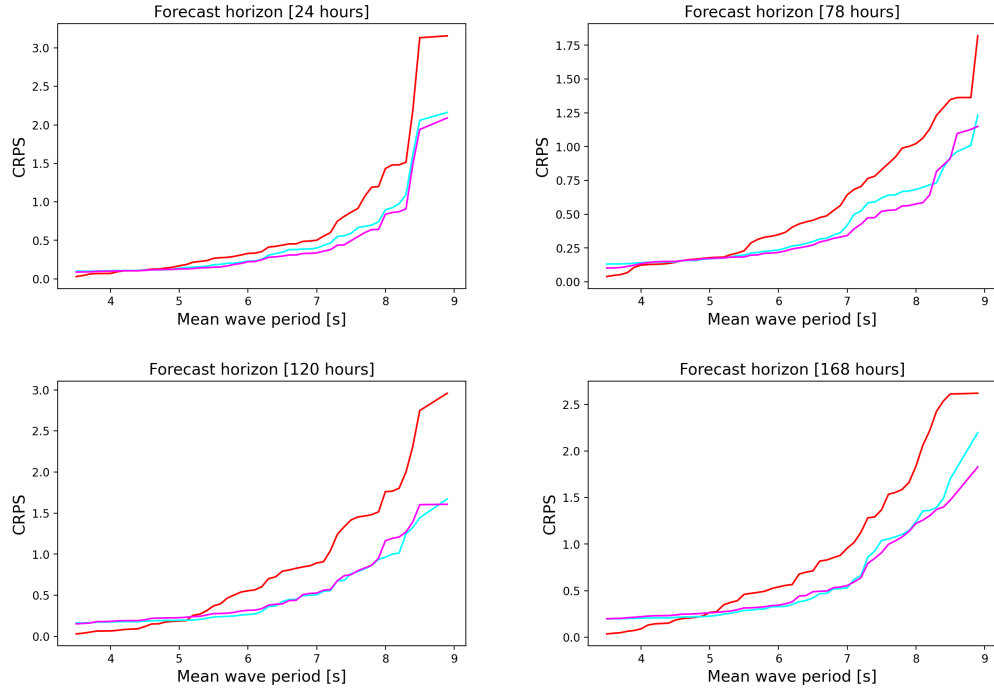
(b) Wind speed  $W$ .

Figure SM6: Variation of CRPS with value of response  $y$  for four forecast horizons  $\tau$  on out-of-sample Period 1 (1 September 2023 - 30 April 2024) for (a) significant wave height  $H_S$  and (b) wind speed  $W$ . CRPS shown for the uncalibrated deterministic forecast (red), the forecast calibrated using LR (cyan) and the forecast calibrated using NHGR (magenta). CRPS performance of calibrated forecasts better than that of the uncalibrated forecast, with an approximate reduction of 50% in value. Moreover, the NHGR calibrated forecast is generally slightly better than that calibrated using LR.



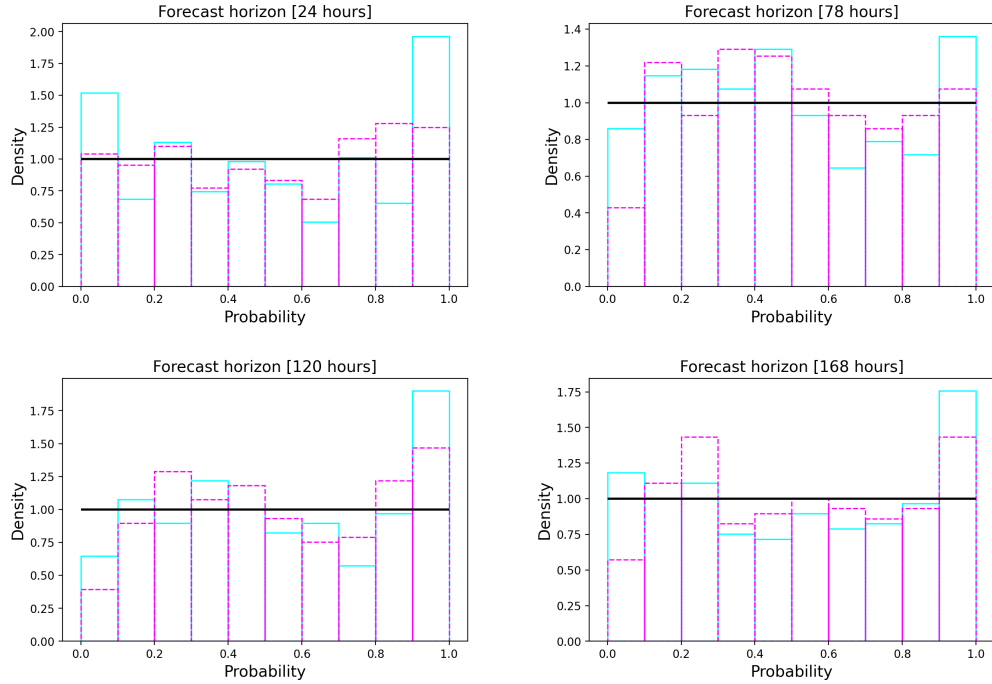


(a) Significant wave height  $H_S$ .

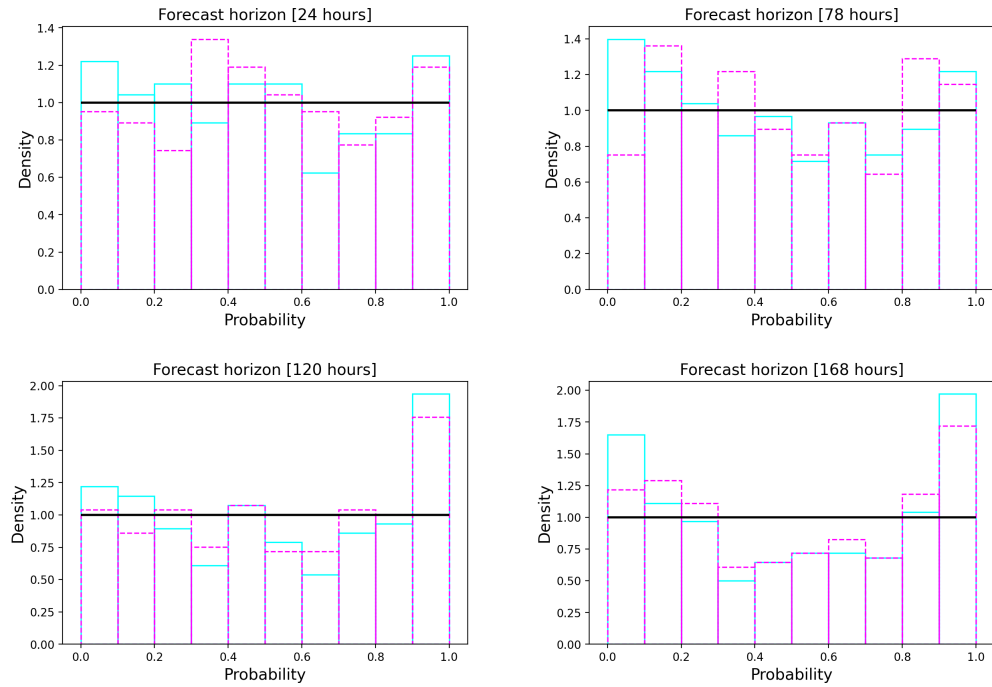


(b) Wind speed  $W$ .

Figure SM7: Variation of CRPS with value of response  $y$  for four forecast horizons  $\tau$  on out-of-sample Period 2 (1 September 2024 - 31 December 2024) for (a) significant wave height  $H_S$  and (b) mean wave period  $T_M$ . CRPS shown for the uncalibrated deterministic forecast (red), the forecast calibrated using LR (cyan) and the forecast calibrated using NHGR (magenta). CRPS performance of calibrated forecasts better than that of the uncalibrated forecast, with an approximate reduction of 50% in value. Moreover, the NHGR calibrated forecast is generally (but not always) slightly better than that calibrated using LR.

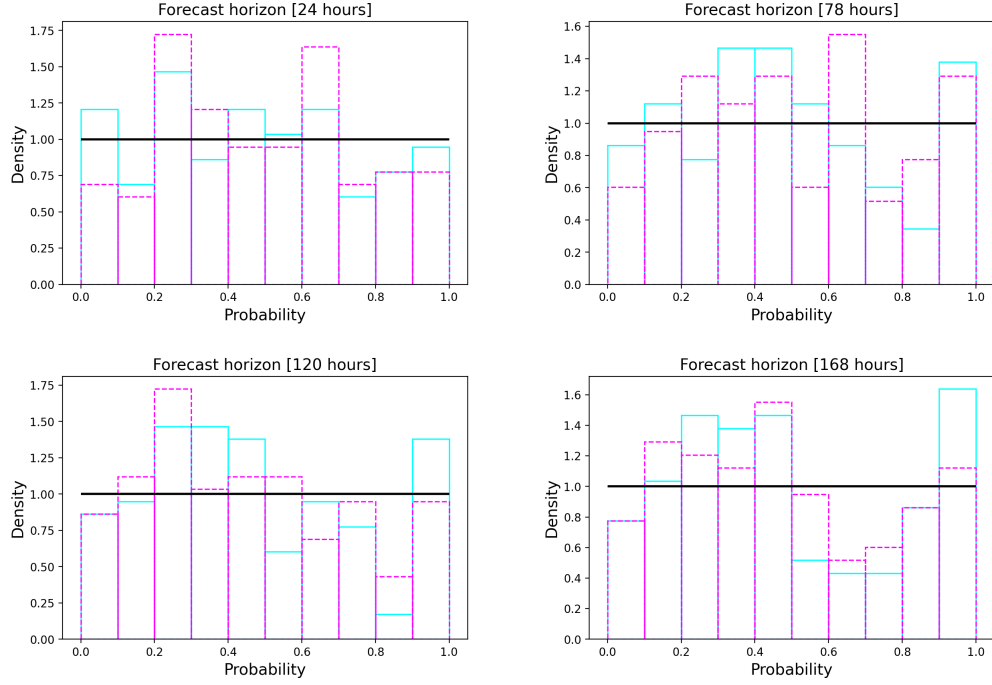


(a) Significant wave height  $H_S$ .

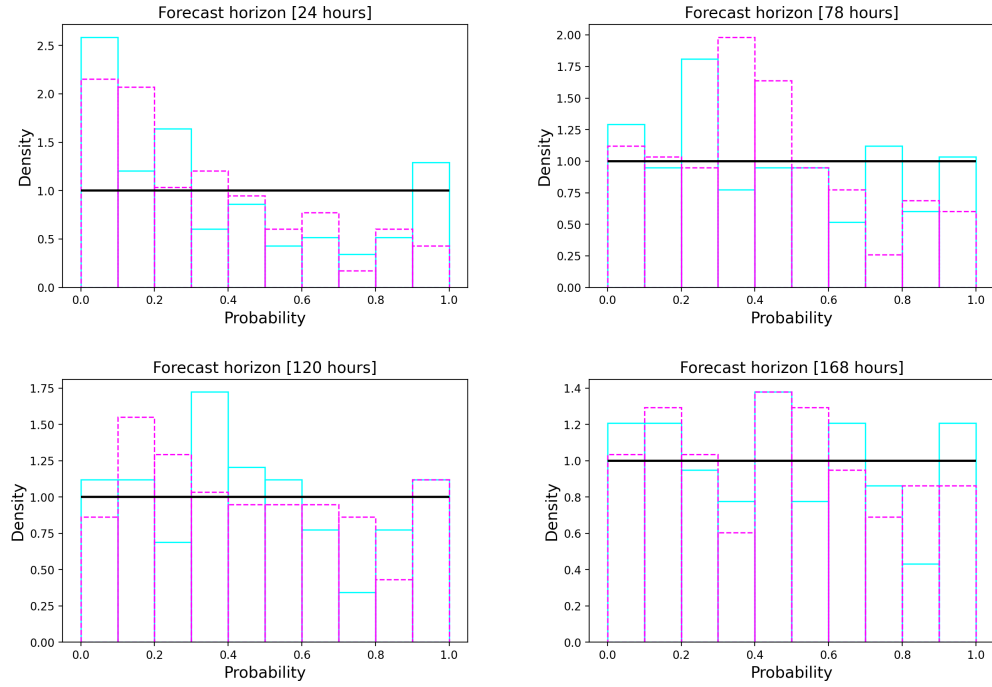


(b) Wind speed  $W$ .

Figure SM8: Rank histogram plots for four forecast horizons  $\tau$  on out-of-sample Period 1 (1 September 2023 - 30 April 2024) for (a) significant wave height  $H_S$  and (b) wind speed  $W$ . Histograms are shown for the forecast calibrated using LR (cyan) and the forecast calibrated using NHGR (magenta). The horizontal line (black) indicates the model-assumed histogram, and departures from it indicate lack of agreement between model and reality. Model performance for both LR- and NHGR-calibrated forecasts is seen to be reasonable, with NHGR calibration to be preferred.



(a) Significant wave height  $H_S$ .



(b) Mean wave period  $T_M$ .

Figure SM9: Rank histogram plots for four forecast horizons  $\tau$  on out-of-sample Period 2 (1 September 2024 - 31 December 2024) for (a) significant wave height  $H_S$  and (b) mean wave period  $T_M$ . Histograms are shown for the forecast calibrated using LR (cyan) and the forecast calibrated using NHGR (magenta). The horizontal line (black) indicates the model-assumed histogram, and departures from it indicate lack of agreement between model and reality. Model performance for both LR- and NHGR-calibrated forecasts is seen to be reasonable, with NHGR calibration to be preferred generally.