

Copenhagen, 26 October 2018

Dear editor,

Below, you find the referees' comments to the first revision of our manuscript with our responses in *italic*.

Response to Referee #1

The manuscript has improved a lot. I have some minor comments that I suggest to be taken account before publishing the manuscript.

The language of the manuscript should be checked by a native English speaker. There were also some typos.

Our english-born colleague corrected the english language in the manuscript.

Line 26: I suggest changing "severe surface waves" to "severe wave conditions"

OK - done

Lines 72-81: This should be in Discussion-section

OK - done

Lines 123-125: Is the North Atlantic grid forced by ECMWF-HRES also run 4 times a day? To my understanding the forcing is available only twice a day.

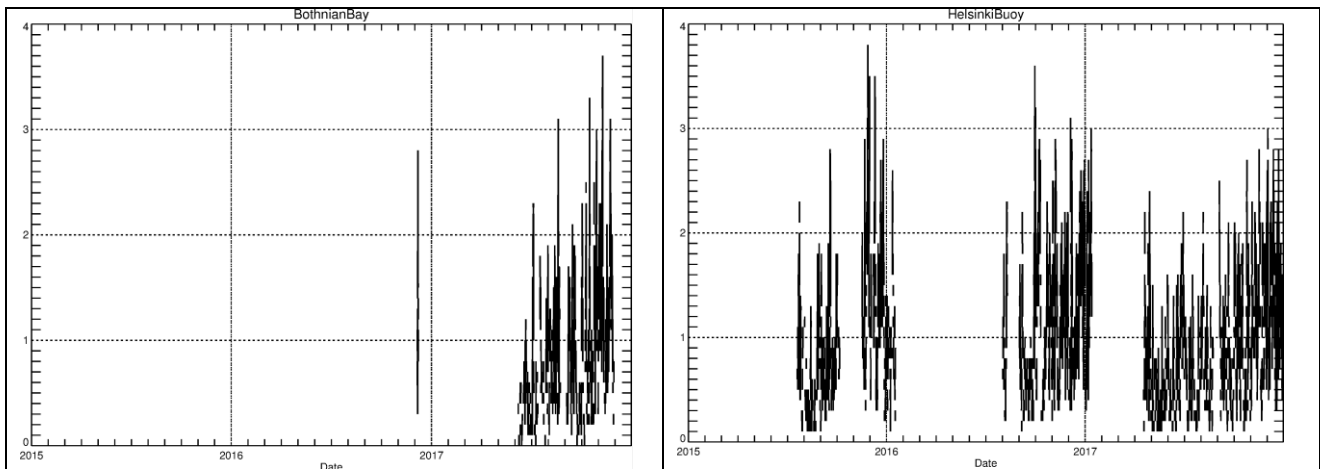
Yes, it is run four times/day, even if a new forcing is available only twice/day. This has been specified in the text.

Lines 153-154: If you would only increase grid resolution, not the spectral resolution. How would the computational efforts then compare. If you can not add an exact number, maybe some discussion about this would be appropriate. This experiment setup does not show, whether the possible benefits of the higher resolution comes from increasing horizontal or spectral resolution.

We have split the timing calculations into two parts in the manuscript, one for the spatial resolution, and one for the spectral resolution.

Lines 174-175: I again comment the selection of buoys used for comparison. Your requirement of more than 40% of temporal coverage basically leaves out all the buoys that are in the areas, where the seasonal ice cover typically ranges from Dec/Jan to May, such as the Bothnian Bay and Gulf of Finland.

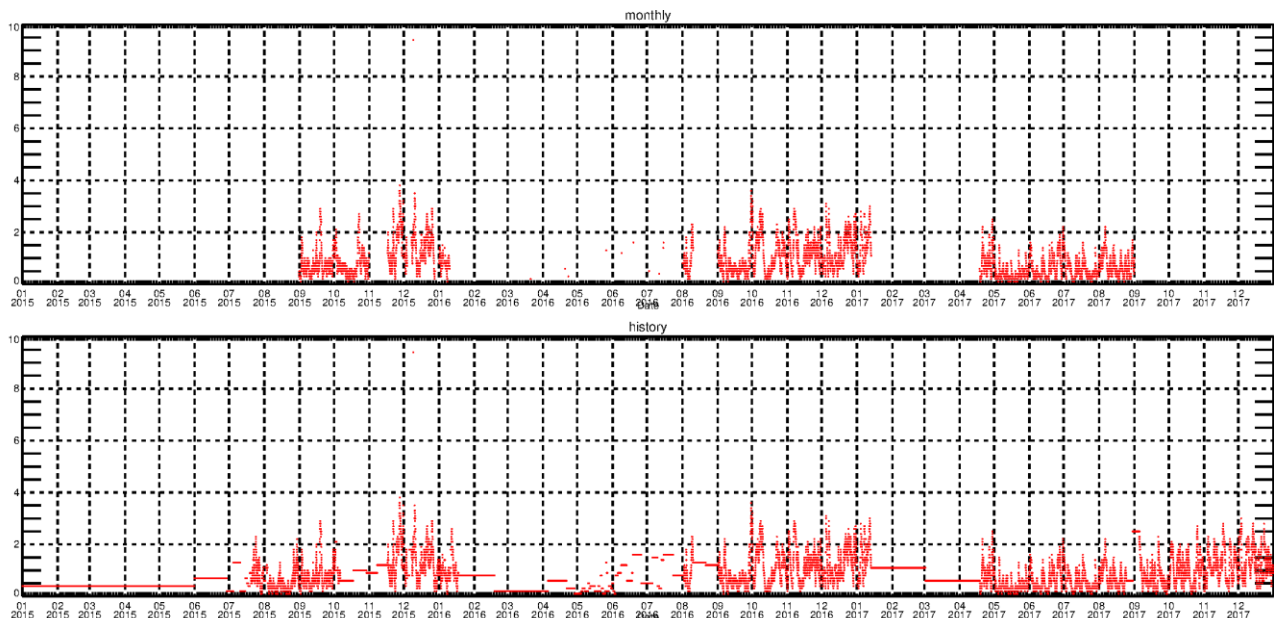
Maybe you should use this criteria in a way that accounts for the time at each buoy location the sea is ice-free.



We have considered the two suggested stations, Bothnian Bay and Helsinki Buoy/Gulf of Finland, see panels above. As for Bothnian Bay it has data for last half of 2017 only and that's the reason why we still exclude it from the study.

Helsinki Buoy has missing data; not only during the ice period (January-May), as expected, but in 2015 until July and with an additional break in October/November, and in 2016 until September. Only in 2017 operation is resumed in May. This not-so-good data coverage during summer may introduce biases in the verification measures, and therefore we continue to exclude station from the analysis. We have made a remark about this in the manuscript.

A closer inspection reveals that the reason for the not-so-good data coverage for Helsinki Buoy is corrupt data (constant value) during a part of the period, see panels below.



Line 176: Comment to the 'sites did not observe the full hour'. WaveRiders, which most of the buoys used in validation were, measure waves for 15-30 min period, and calculate the wave parameters based to that measurements. Did you check, whether the timestamp in the measured data files was the starting time, mid-time, or end time of this measurement period?

No, we did not go into details of the procedures of wave measurements. These may differ among instruments. We assume that the time given is representative for the measurement.

Table 4: What about higher values than 5m?

We have added a column '> 5 m' to the table.

Table 5: Why only Danish stations? I'm not sure that they are the best ones, at least as only ones, to describe the wind field accuracy in the Baltic. At least SMHI and FMI have coastal weather station data available in they open data portals. You should consider using them.

We have added Swedish and Finish coastal stations and now verify based on eight stations in total. The section on verification has been re-written, and the table replaced by a figure showing total rmse.

Line 310: According to Björkqvist et al the NB wave buoy measured 8.0 m , not 'almost 8m' during the storm

OK -corrected

Line 362: I would not state that the performance of LOWENSMEAN and LOWENS were "superior" to HIGH. For sure, it was better in many cases.

OK – we have changed the formulation

Lines 376-386: I would start with informing the readers that these numbers are not directly comparable. Maybe you could start with lines 384-386.

OK – we have re-arranged the text.

Line 398: I would not call Tuomi et al. 2017 a study, rather a standard product validation procedure performed in CMEMS.

OK – we have changed the text

Response to Referee #2

Better Baltic Sea Wave forecasts: Improving resolution or introducing ensembles?

Reviewed at iteration: Revised Submission

Recommendation: paper is accepted subject to minor/technical revisions.

Generic remarks: Thanks a lot for the authors for taking on board the previous review comments. In particular the improved documentation of the atmospheric forcing used by the wave models has made the paper and its conclusions a lot clearer. I would be happy to see this published without a further review subject to some technical (written English) corrections suggested below. However, I'd also like the authors to consider some potential minor revisions to the manuscript.

Potential minor revisions:

Section 2: Please confirm in the text whether any of the source term tuning parameters are different or the same in the various wave configurations. Previous review comments had also asked for the WAM version(s) to be documented.

We made a more precise formulation in the manuscript.

Section 2: Is it possible to add some justification as to why the authors might consider the 3km atmospheric system to offer a step change in wind forcing skill over the 5km configuration that underpins the ensemble – this is not an enormous resolution change, so which processes are believed to be improved? Ideally provide references for the atmosphere systems in question.

The difference in horizontal resolution between deterministic and ensemble forecasts are due to heuristics more than to science. The deterministic forecasts were introduced first and 3 km was what was affordable on the computer. Later, the ensemble forecasts were introduced, and here 5 km was regarded as the affordable resolution.

Section 5: Not suggesting that additional results/tables are presented – but is it possible to comment on the performance of the S05 atmosphere model's control member performance? Is it similar to, better or worse than S03? This would help put the effect of the ensemble into context.

The main idea of introducing the ensemble approach is to get an overall better forecast based on all ensemble members. The control member does not have a special role.

Section 6: As per section 5, if it is possible to comment on performance of S05's control member relative to the LOW wave model that will help contextualise the ensemble.

See above

Figure 9: Is it possible to show the range of LOWENSMEAN forecasts – it would be good to see if any of the ensemble members predicted the peak in SWH.

We have added all ensemble members to the figure and added an appropriate sentence in the text.

Section 7 (or 8): I'd still be keen to see a little more discussion on the systems, what has been verified and what that means in terms of application. Now that the source of wind forcing has been stated more explicitly, what we have are: LOW – a system with 3km winds integrated onto a 10km wave grid; HIGH – 3km winds integrated onto a 5km grid; and LOWENS – a system with 5km winds integrated on a 10km grid but run as an ensemble in order to sample uncertainty in the atmosphere model and evolution of the meteorological conditions. Validation is carried out on a site specific basis. So, the important points to consider are:

1. Whilst the validation is a very standard approach for wave models, it is actually documenting site-specific forecast performance for the model. This is an important consideration when testing high resolution systems, where 'double-counting errors' start to get introduced and site specific verification does not improve as a result. On the other hand, the majority of forecast use for these models are site-specific, so the results are absolutely valid when one considers the products that get generated from these models.
2. The comparison between LOW and HIGH suggest that integrating the 3km winds onto a coarser wave model grid has no impact in terms of the verification at offshore sites. Why might this be? Speculating, this is probably because wave development is a function of winds, but integrated over a longer fetch area – so the resolution of the atmosphere model is (within reason) perhaps less important than the model's ability to properly place major synoptic features. Figure 2 hints at this, since there appears to be little addition structure in the HIGH model field offshore compared to LOW.
3. The ensemble mean gives the best site-specific forecasts offshore, in terms of limiting the overall error. This is despite using a coarser resolved atmospheric model and the coarser wave model. This is consistent with the argument above that uncertainty large scale feature development is perhaps (generally) more important to wave forecasts in open waters than getting the momentum exchange associated with small scale atmospheric features correct. I'd suggest that argument is further supported if the LOWENS control member verification at the offshore locations is not significantly worse than for LOW and HIGH.

I appreciate it if points 2 and 3 feel too speculative for the authors to consider including in their discussion, but I would advocate some comment along the lines of point 1.

We have added a formulation addressing your point 1 in the beginning of the discussion section.

Section 8: The under-spread in the ensemble suggests that there is scope for improving that system. I'd suggest this is a valid conclusion.

We agree and have added a sentence at the end of the conclusion section.

Technical revisions:

Line 31-32: dissipation of the wave energy mainly occurs through internal processes, e.g. whitecapping.

OK, done.

Line 48-49: The equations of the NWP model are discretized on a horizontal grid with a certain spatial resolution, which influences the maximum spatial resolution of the wave model.

OK - done

Line 52: Over time, technical development has increased available computational resources, making it possible to increase...

OK - done

Line 75: ...modelled sea-surface temperatures (SSTs) by the NEMO...

OK - done

Line 76-77: Introducing such coupling may demand a high horizontal resolution, in atmosphere, wave and ocean models, in order to describe the fluxes most satisfactorily.

OK - done

Line 84: ...wind forecasts is in Section 5, whilst verification of a principle wave forecast variable, significant wave height (SWH), is presented in Section 6.

OK - done

Line 123: Each forecast run...

Thank you – done

Line 133: ...with characteristics identical to LOW, but using a parallel...

OK - done

Line 167: ...the area with SWH above 6m extends further southward...

Thank you - done

Line 170-171: Observed series of SWH from wave measurement sites in the Baltic Sea, obtained from the Copernicus Marine Environment Monitoring Service (CMEMS) database, are used.

OK - done

Line 306: ...HIGH forecast has a significantly smaller under-prediction bias than the other forecast classes.

OK - done

Line 364: The conclusions hold,...

OK - done

Line 410: ...field approaches an ice-covered area,...

OK - done

Lines 412-413: ...when dense enough, acts as a solid shield that effectively removes all local wave energy...

Thank you - done

Line 414: ...thick enough for this to be approximately correct.

Thank you - done

Line 445-446: ...there are no indications that a further increase of the WAM model will result in enhanced site-specific forecast performance.

OK - done

Better Baltic Sea wave forecasts: Improving resolution or introducing ensembles?

Torben Schmith, Jacob Woge Nielsen, Till Andreas Soya Rasmussen, Henrik Feddersen

Danish Meteorological Institute, Copenhagen, Denmark

Correspondence to: Torben Schmith (ts@dmi.dk)

Abstract. The performance of short-range operational forecasts of significant wave height in the Baltic Sea is evaluated. Forecasts produced by a base configuration are inter-compared with forecasts from two improved configurations: one with improved horizontal and spectral resolution and one with ensembles representing uncertainties in the physics of the forcing wind field and the initial conditions of this field. Both the improved forecast classes represent an almost equal increase in computational costs. The inter-comparison therefore addresses the question: would more computer resources most favorably be spent on enhancing the spatial and spectral resolution or, alternatively, on introducing ensembles? The inter-comparison is based on comparisons with hourly observations of significant wave height from seven observation sites in the Baltic Sea during the three-year period 2015-2017. We conclude that for most wave measurement sites, the introduction of ensembles enhances the overall performance of the forecasts, whereas increasing the horizontal and spectral resolution does not. These sites represent offshore conditions, well exposed from all directions with a large distance to the nearest coast and with a large water depth. Therefore, the detailed shoreline and bathymetry is also a priori not expected to have any impact. Only ~~for at~~ one site, do we find that increasing the horizontal and spectral resolution significantly improved the forecasts. This site is situated in nearshore conditions, close to land, with a nearby island and therefore shielded from many directions. This study therefore concludes that to improve wave forecasts in offshore areas, ensembles should be introduced. For near shore areas, the study suggests that additional computational resources should be used to increase the resolution.

1 Introduction

Severe ~~surface waves~~wave conditions affect ship navigation, offshore activities and risk management in coastal areas. Therefore, reliable forecasts of wave conditions are important for ship routing and planning purposes when constructing, maintaining and operating offshore facilities, such as wind farms and oil installations.

Waves are generated by energy transfer from surface winds that act on the sea. The energy transfer is determined by the *fetch* (the distance, over which the wind acts), and by the *duration* of the wind. For *deep water waves*, defined as the wave height being much smaller than the water depth, dissipation of the wave energy mainly occurs through internal processes, e.g. whitecapping~~dissipation of the wave energy occurs through internal dissipation mainly~~. For *shallow water waves*, defined as the wave height being comparable to the water depth, dissipation through bottom friction and through wave breaking over a shallow and sloping sea bed becomes important. Shallow water waves may also be refracted over a varying bathymetry

37 Therefore, a correct and detailed description of the bathymetry is important for correctly forecasting waves
38 in coastal areas and other shallow sea areas. Other factors with a potential effect on the development of
39 waves include nonlinear wave-wave interaction, ocean currents, time-varying water depth due to variations
40 in sea level, and sea ice coverage.

41 The Baltic Sea is connected to the world ocean through the Danish waters with shallow and narrow Straits
42 (see Figure 1), and this allows virtually no external wave energy to be propagated into the area. The Baltic
43 Sea consists of a number of basins with depths exceeding 100 m, separated by sills and water areas with
44 more moderate water depths. Between Finland and Sweden lies an archipelago with complicated
45 bathymetry on very small spatial scales. The wind is in general westerly over the area, and the most
46 prominent cause for severe wind and wave conditions is low pressure systems passing eastward over
47 central Scandinavia. Winter ice occurs in the northern and eastern parts of the Baltic Sea. There is no
48 noticeable tidal amplitude or permanent current systems.

49 Short-term forecasting of surface waves is done by a wave model, forced with forecasted wind from an
50 atmospheric numerical weather prediction (NWP) model. The equations of the NWP model are discretized
51 on a horizontal grid with a certain spatial resolution, which influences the maximum spatial resolution of
52 the wave model.~~The equations of the NWP model are discretized on a horizontal grid with a certain spatial~~
53 ~~resolution, which determines the maximum spatial resolution of the wave model.~~ The available computer
54 resources ~~put a limit~~s on the horizontal grid spacing, ~~that~~which can be afforded.

55 Over time, technical development has increased available computational resources, which traditionally
56 hasve been used to making it possible to increase~~Technical development has increased the computational~~
57 ~~resources, making possible to increase~~ the horizontal spatial resolution of the NWP and wave models. This
58 allows for an improved description and forecasting of the synoptic and mesoscale atmospheric systems,
59 including the details of the associated wind field. In addition, a more detailed description of the bathymetry
60 improves the correct description of dissipation and refraction of waves, as argued above. Additional
61 computer resources may also be used to improve the spectral resolution in the wave model. This includes
62 the directional resolution and the number of frequencies included.

63 Increasing computer resources have also made ensemble NWP possible. The purpose of ensemble
64 forecasts is to improve forecast skill by taking both the initial error of the forecast and the uncertainty of
65 the model physics into account. Furthermore, ensemble forecast allows for probabilistic forecasts,
66 identified as a priority for operational oceanography (She et al., 2016), and allows for quantifying forecast
67 uncertainty. Ensemble wave forecast systems have been implemented at global scale (Alves et al., 2013;
68 Cao et al., 2009; Saetra and Bidlot, 2002) and more regionally in the Norwegian Sea (Carrasco and Saetra,
69 2008), and in the German Bight and Western Baltic (Behrens, 2015).

70 From the above discussion it is evident that additional computer resources can be used in different ways to
71 change the wave forecast setup, in order to increase the forecast quality. The purpose of the present study
72 is to investigate the effect on the forecast quality of increasing the horizontal resolution and the spectral
73 resolution vs. introducing ensemble forecasts. This will be done by verifying the DMI operational
74 forecasting of wave conditions in the Baltic Sea in different configurations against available observations of
75 significant wave height.

~~It should be mentioned that improving wave forecasts is not the only driving factor in reducing the grid size of the wave model. Coupling the wave model with atmosphere or ocean circulation models may give a better description of vertical fluxes of heat and momentum (Cavaleri et al., 2012). For instance, Alari et al. (2016) documented a significant improvement of modelled SSTs by the NEMO circulation model in the Baltic Sea when a two-way coupling to the wave model WAM was introduced. Doing such couplings may demand a high horizontal resolution to describe the fluxes most satisfactorily.~~

~~Also if~~ Increasing the horizontal resolution of the NWP-system may also lead to improved wind forecasts, due to in particular better descriptions of processes in extratropical cyclones. In these cases, where the wind field is strong and varying on a small spatial scale, ~~also~~ wave forecasts may also be improved by running the wave model in a similarly high resolution.

This paper is arranged as follows. Section 2 describes the model and setup, Section 3 describes the observations used and the verification methodology is described in Section 4. Verification of DMI-HIRLAM wind forecasts is in Section 5, whilst verification of the significant wave height (SWH) is presented, ~~and the SWH forecast verification is~~ in Section 6. Results of the verification are discussed in Section 7 and conclusions made in Section 8.

2 Model and setup

The DMI operational wave forecasting system DMI-WAM uses the 3rd generation spectral wave model WAM Cycle4.5.1 (Günther et al., 1992), with one minor change of source term functions. To speed up wave growth from calm sea, the spectral energy has a lower limit corresponding to a wave height of 7 cm. It is forced by the regional NWP model DMI-HIRLAM and the global NWP model ECMWF-GLM. WAM Cycle4.5 solves the spectral wave equation, and calculates the wave energy as a function of position, time, wave period and direction. Derived variables, such as the significant wave height (SWH), are calculated as suitable integrals of the wave energy spectrum.

The DMI-WAM model system forecasts waves in a larger area than the Baltic Sea and therefore has a setup with two nested spatial domains of different geographical extent (see Figure 1): North Atlantic (NA) and North Sea/Baltic Sea (NSB), of which forecast results from the NSB-domain are analyzed in this study. The NA domain uses the JONSWAP wave spectrum for fully developed wind-sea (Hasselmann et al., 1973) along open model boundaries, while the NSB domain use modeled wave spectra from the NA domain at its open boundaries (one-way nesting).

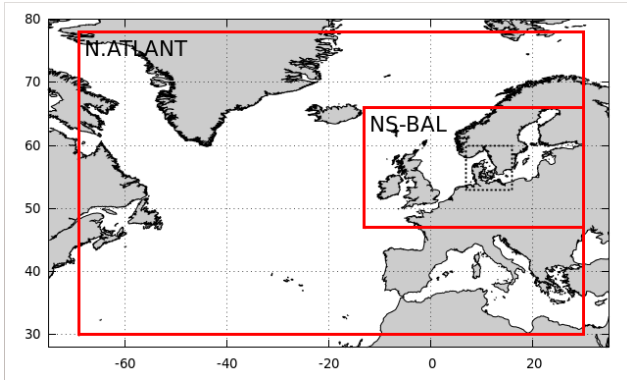


Figure 1 Nesting of domains in DMI-WAM. Outer frame is North Atlantic (NA) domain, inner frame is the North Sea/Baltic Sea(NSB)-domain. Dotted frame is the Transition Area. Only data from the NSB-domain are analyzed in this study.

The wave energy is discretized into a number of wave directions and frequencies. To facilitate wave growth from calm sea, a lower limit is applied to the spectral energy. The resulting surface roughness parameterizes the effect of capillary waves, and corresponds to a minimum significant wave height of 7 cm.

The energy source is the surface wind. The sink terms are wave energy dissipation through wave breaking (white capping), wave breaking in shallow areas, and friction against the sea bed. Depth-induced wave breaking (Battjes and Janssen, 1978) is used in the NSB domain only, since in the NA domain, the depth maps are not detailed enough for activation of this effect. The wave energy is redistributed spatially by wave propagation and depth refraction, and spectrally by non-linear wave-wave interaction. Interaction with ocean currents and effects due to varying sea level caused by tides or storms are not incorporated.

In addition to a land mask, we have a time-varying ice mask. Below ice 30% concentration, sea ice is assumed to have no effect. Above 30% ice concentration, no wave energy is generated or propagated, i.e. the effect is like that of land. The applied sea ice concentrations originate from OSISAF (<http://osisaf.met.no/p/ice/>) with a frequency of 24 hours and around 25 km true horizontal resolution, gridded to ~10 km horizontal resolution and interpolated to the WAM-grid. The ice cover is initialized every day at 00z, and kept constant throughout each forecast run.

The surface wind forcing is provided by different atmospheric models for the two domains. For the NA domain, wind is provided by the ECMWF-HRES global weather forecast every 3 hours. For the NSB domain, the surface wind is provided every hour by DMI-HIRLAM. Setup details are summarized in Table 1

Table 1 Specifications of DMI-WAM nested setup.

Domain	North Atlantic	North Sea/Baltic Sea
Longitude	69W-30E	13W-30E
Latitude	30N-78N	47N-66N
Atmospheric forcing	ECMWF-HRES	DMI-HIRLAM
Boundary condition	JONSWAP	One-way nested
Depth-induced wave breaking	No	Yes

Each forecast run is initialized using the sea state at analysis time, calculated by the previous run as a six hour forecast. The operational DMI-WAM suite is run four times a day to 48 h forecast range. This is also true for the North Atlantic domain, even when new forcing is available twice per day only. This is for practical reasons, since the North Atlantic domain is very cheap to run. Spatial fields of forecasted SWH and other variables are output in hourly time resolution.

Historically, three different configurations of the DMI-WAM setup have been used, and data from these for the period 2015-2017 is the basis for the present verification. In the old LOW configuration, the horizontal resolution is around 50 km in the NA domain and around 10 km in the NSB domain, ~~and the~~ wave energy is resolved in 24 directions and at 32 frequencies, corresponding to wave periods ~~of between~~ 1.25-23.94 s and wave lengths ~~between of~~ 2.4-895 m (in deep water). Bathymetry is ETOPO (Amante and Eakins, 2009) in the NA domain, and the Baltic bathymetry from IOW (<https://www.io-warnemuende.de/topography-of-the-baltic-sea.html>) supplemented by depth data from the Danish Geodata Agency (DGA) in the NSB domain. More recently, an ensemble configuration (LOWENS) has been introduced with characteristics identical to LOW, but ~~with using a~~ parallel run of 11 ensemble members forced with perturbed atmospheric fields (initial conditions and physics). Finally, in the ~~also~~ recently introduced HIGH configuration, the horizontal resolution is around 25 km in the NA domain and around 5 km in the NSB domain, ~~and the~~ wave energy ~~is~~ resolved in 36 directions and 35 frequencies, corresponding to wave periods ~~between of~~ 0.94-23.94 s, and wave lengths ~~between of~~ 1.37-895 m (in deep water). Bathymetry is RTopo (Schaffer et al., 2016).

All configurations are forced by winds from ECMWF-HRES in the NA domain and DMI-HIRLAM in the NSB domain. In the NSB domain, the LOW and HIGH are forced by the S03 version (3 km horizontal resolution), while LOWENS is forced by the S05 version (5 km horizontal resolution). The S03 and S05 versions of DMI-HIRLAM were used operationally by DMI as deterministic and ensemble weather forecast models in the 2015-17 period. While the better resolution of S03 might have an impact on forecasts where orographic effects matter, the impact on wind forecasts over sea is expected to be insignificant. The DMI-HIRLAM winds are interpolated to the WAM grids by bilinear interpolation. To diminish coastal effects, DMI-HIRLAM delivers a special *water-wind* to DMI-WAM, in which the surface roughness everywhere is assumed to be that of water. This enhances the wind speed in the coastal zone, most important in semi-enclosed areas (bays, fjords, etc.). It is basically a way to sharpen the land/sea boundary, reducing influence of land roughness on near-shore winds. An overview of the DMI-WAM configurations is provided in Table 2 ~~Table 2~~.

Table 2 Details of DMI-WAM configuration used in this study.

	DMI-WAM Horizontal resolution [km]		# wave directions	# wave spectral frequencies	Bathymetry		Atmospheric horizontal resolution [km]		Ensemble members	
	North Atlantic	NSB			North Atlantic	NSB	North Atlantic (ECMWF)	NSB (DMI- HIRLAM)	North Atlantic	NSB
LOW	50	10	24	32	ETOPO	IOW/DGA	16	3	-	-
LOWENS	50	10	24	32	ETOPO	IOW/DGA	16	5	-	11
HIGH	25	5	36	35	RTopo	RTopo	16	3	-	-

162 When replacing the LOW forecast configuration with the HIGH configuration, the required computational
163 resources for running DMI-WAM are increased by a factor of 2^2 (increase in horizontal resolution) \times 1.75
164 (effective decrease in time step) = 7 due to higher spatial resolution, and by a factor of \times 1.5 (increase of
165 number of directions) \times 35/32 (increase of number of spectral frequencies) = 1.6. This gives a total factor of
166 $7 \times 1.6 \approx 11.5$. From the LOW to the LOWENS configuration, it is increased by a factor of 11 (number of
167 ensemble members). Since these increases in computational effort are very similar, an inter-comparison
168 can contribute to answering the question: should additional computer resources be used for increasing the
169 spatial and spectral resolution, or for sampling the uncertainty in meteorological conditions using
170 ensembles.

171 The LOW and HIGH configurations both produce a class of deterministic forecast, which are also named
172 LOW and HIGH, respectively. The LOWENS configuration produces a class of probabilistic forecast, called
173 LOWENS. In addition, the ensemble mean defines a class of deterministic forecasts, called LOWENSMEAN.

174 To illustrate differences to be expected among the deterministic forecasts, we show 48 h forecasts of SWH
175 valid at the peak of the 'Toini' storm on 10 January 2017.

176

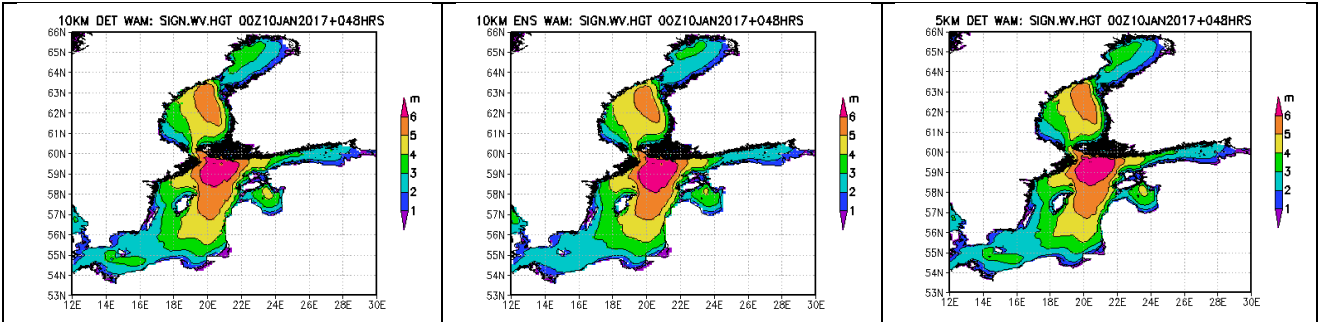


Figure 2 Forecasted (48h) SWH at the peak of ‘Toini’ storm 10 january 2017 00z for LOW (left), LOWENSMEAN (middle) and HIGH (right) forecasts.

All three forecasts agree in the gross features of the forecasted SWH field. However, there are differences, e.g., northeast of the island of Gotland, the area with SWH above 6 m extends further southward in the LOWENSMEAN forecast, than in the LOW and HIGH forecasts.

3 Observations

Observed series of SWH from wave measurement sites in the Baltic Sea, obtained from the Copernicus Marine Environmental Monitoring System (CMEMS) database, are used. None of the series has a continuous record over the three-year period 2015 – 2017. Data gaps may be due to malfunction, maintenance or withdrawal of the instrument. The latter occur during winter due to the possibility of ice. We selected sites with valid observations that covered more than 40% and were distributed reasonably throughout the study period. To avoid biases in the verification measures due to under- or overrepresentation of particular seasons, we also aimed at having an approximately even coverage throughout the year.

-Figure 3 and Table 3 show the positions and water depths of the wave measurement sites together with the bathymetry of the Baltic Sea. Some sites did not observe at the full hour. Observations from these sites were ascribed to the nearest full hour, if the time distance between the observation time and the full hour was less than 15 min, otherwise not used. All observation series used are shown in Figure S1. The frequency of observed SWH in different intervals for each site is given in Table 4

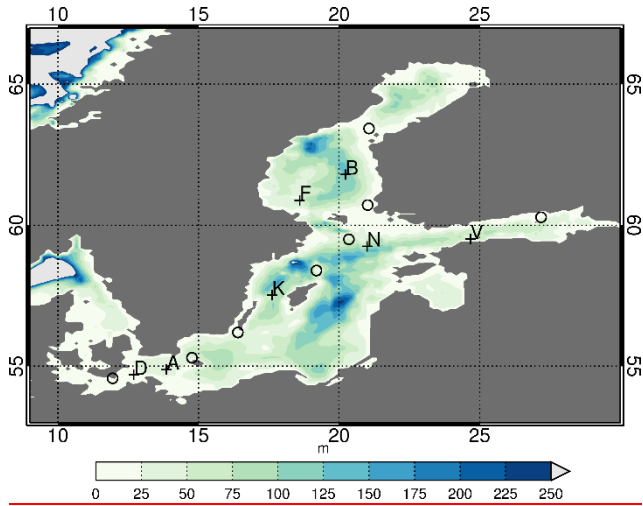


Figure 3 Map of the Baltic Sea with bathymetry and positions of wave measurement sites marked with crosses. For details about sites, see Table 3. Meteorological stations used in the wind verification of DMI-HIRLAM are marked with circles.

Table 3 Details of wave measurement sites.

Observation site	Lon	Lat	Depth [m]	
			Model	Actual
A Arkona WR	13.9	54.9	46	45
B Bothnian Sea	20.2	61.8	118	~120
D Darsser Sill WR	12.7	54.7	20	21
F Finngrundet WR	18.6	60.9	56	67
K Knolls Grund	17.6	57.5	63	90
N Northern Baltic	21.0	59.2	68	~100
V Vahemadal	24.7	59.5	18	5

Table 4 Observed frequency of SWH in different bins for wave measurement sites.

SWH [m]	0-1	1-2	2-3	3-4	4-5	>5
Arkona WR	0.47	0.39	0.12	0.01	<0.01	<0.01
Bothnian Sea	0.46	0.38	0.12	0.02	0.01	<0.01
Darsser Sill WR	0.67	0.31	0.02	<0.01	<0.01	<0.01
Finngrundet WR	0.69	0.27	0.04	0.01	<0.01	<0.01
Knolls Grund	0.62	0.31	0.06	0.01	<0.01	<0.01
Northern Baltic	0.39	0.37	0.18	0.05	0.01	<0.01
Vahemadal	0.78	0.20	0.02	<0.01	<0.01	<0.01

4 Verification methodology

In this section, a short overview of the verification procedure will be given. For background and more details regarding the verification measures, we refer to (Jolliffe and Stephenson, 2003)

208 For each measurement series of SWH, the corresponding forecast series for all forecast classes and for
 209 forecast range zero to 48 h for the grid point nearest to the position of the wave measurement site was
 210 extracted from the model output.

211 For the deterministic and continuous forecast classes (LOW, LOWENSMEAN and HIGH), we use the
 212 conventional performance measures *root mean square error* (RMSE), defined as the square root of the time
 213 average of the sum of squared differences between forecast and observation:

$$RMSE(\tau) = \langle (h_{s,fcst}^{\tau} - h_{s,obs})^2 \rangle$$

214 the bias

$$BIAS(\tau) = \langle h_{s,fcst}^{\tau} - h_{s,obs} \rangle,$$

216 and the correlation coefficient

$$CC = \frac{\langle (h_{s,fcst}^{\tau} - \langle h_{s,fcst}^{\tau} \rangle)(h_{s,obs} - \langle h_{s,obs} \rangle) \rangle}{\sqrt{\langle (h_{s,fcst}^{\tau} - \langle h_{s,fcst}^{\tau} \rangle)^2 \rangle \langle (h_{s,obs} - \langle h_{s,obs} \rangle)^2 \rangle}}$$

217 where $h_{s,obs}$ is the observed SWH and $h_{s,fcst}^{\tau}$ is a corresponding forecast with forecast range τ .

218 | The RMSE is a positive definite quantitative measure, and smaller values mean a better forecast. The bias
 219 | can take positive and negative values, and a good forecast has a numerically small value. The averaging,
 220 | indicated by $\langle \cdot \rangle$, ~~can be~~ found based on all available values during the three-year period. Also, the RMSE
 221 | and BIAS as function of $h_{s,obs}$ will be considered.

222 A framework for verifying probabilistic forecasts is the *continuous ranked probability score* (CRPS), defined
 223 as

$$CRPS(\tau) = \langle \int [F^{\tau}(h_s) - H(h_s - h_{s,obs})]^2 dh_s \rangle,$$

225 where $F^{\tau}(h_s)$ is the forecasted probability distribution, $h_{s,obs}$ is the observed value, and $H(\cdot)$ is the
 226 Heaviside step function. A small CRPS occurs when the median of the probabilistic forecasts are close to the
 227 observed values. Also a sharp probabilistic forecast with a small spread favors a small CRPS. This means that
 228 the best forecast is achieved when CRPS is small. CRPS can be applied to both the probabilistic forecast
 229 class LOWENS, as well as the deterministic forecast classes, LOW, LOWENSMEAN and HIGH, since these
 230 can be regarded as probabilistic forecasts with a step probability distribution. For the deterministic forecast
 231 classes, the CPRS equals the *mean absolute error*.

232 Besides the continuous and probabilistic forecasts, also the binary forecast of the SWH exceeding a
 233 specified threshold is considered. The performance measure used is the Brier Score, defined as

$$BS(\tau) = \langle (p - x)^2 \rangle,$$

235 where p is the forecasted probability with forecast range τ of exceeding the threshold and x takes the
 236 value of 1 or 0 dependent on whether the threshold actually was exceeded or not. The Brier Score is thus a

237 positively definite measure, where values are between zero and one, and the lower the value, the better
238 the forecast.

239 **4.1 Calculation of confidence bands**

240 All the measures described above are subject to sampling uncertainty; if they had been calculated on data
241 from another time period than 2015-2017, they would have had different values. To estimate this sampling
242 uncertainty and thereby obtain confidence bands, we applied a block bootstrapping procedure, where a
243 large number of resampled series with the same length as the original series (three years) were created. A
244 blocking length of one month was chosen. This choice takes the atmospheric decorrelation time scale of a
245 few weeks into account and it allows a large number of different resampled series to be made.

246 Each resampled series is constructed as follows: The resampled series will contain three ~~January's~~January
247 months, and each of these is randomly chosen, with replacement, of the three ~~January's~~January months
248 from the original series. A similar procedure applies for February, etc. In this way, the resampled series are
249 most likely different but the annual cycle is preserved. Both the observed series and the forecast series are
250 resampled. For each pair of resampled series bootstrapped value of the performance measures are
251 calculated. Repeating the resampling procedure, we obtain 1000 resampled values of the measures, from
252 which their approximate statistical distribution and confidence bands can be calculated. As a standard,
253 confidence bands (5/95%) are calculated by the bootstrap procedure described above and this allows for a
254 quantitative inter-comparison of the performance measures for the different forecast classes: if the
255 confidence bands do not overlap then there is a significance difference.

256 **5 Verification of the wind forecasts**

257 In order to illustrate the benefit of the meteorological ensemble on wind forecasts the S03 deterministic
258 and S05 ensemble mean have been verified against available wind observations for eight coastal
259 meteorological stations around the Baltic Sea (Figure 3). The RMSE of all stations for the period 1 Jan 2015 -
260 31 Dec 2017 is shown in Figure 4 as a function of forecast range. This reveals that the S05 ensemble mean is
261 more accurate than S03, especially at the longer forecast ranges. Similar results are found for other
262 verification scores, such as correlation and hit rate (not shown).

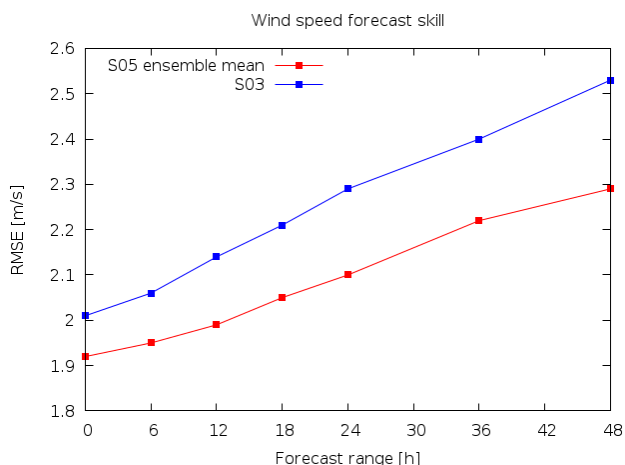


Figure 4 Verification of wind speed. Average RMSE between model and observations for eight coastal meteorological stations in the Baltic Sea area.

The two configurations of DMI-HIRLAM used (see Table 2) have been verified against available wind observations from Danish coastal stations, i.e. covering the western part, of the Baltic Sea, for the period 1 January 2015–31 December 2017. For the S05 configuration, the ensemble mean is verified.

Table 5 Verification results for DMI-HIRLAM against Danish coastal stations for the period 1 January 2015–31 December 2017. Positions of stations are marked on Figure 3.

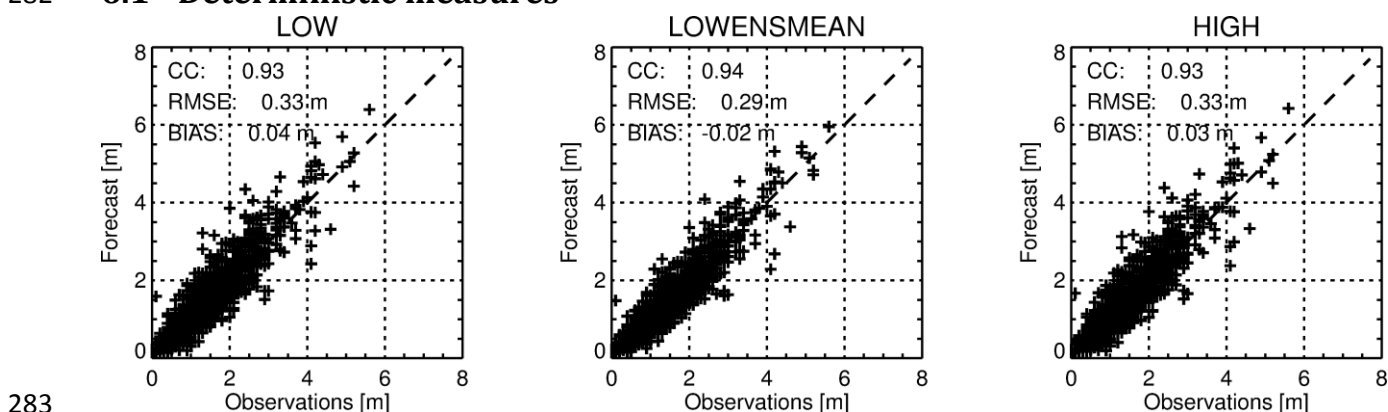
FCST	BIAS [ms^{-1}]		RMSE [ms^{-1}]		CC		Hit rate, error $\leq 2 \text{ ms}^{-1}$	
	S05 (EM)	S03	S05 (EM)	S03	S05 (EM)	S03	S05 (EM)	S03
Gedser (WMO 06149) :								
0	0.48	0.46	1.57	1.56	0.90	0.91	0.82	0.82
6	0.43	0.46	1.58	1.62	0.90	0.90	0.81	0.80
12	0.45	0.49	1.66	1.72	0.89	0.89	0.79	0.78
18	0.45	0.49	1.73	1.83	0.88	0.87	0.77	0.76
24	0.45	0.51	1.77	1.91	0.87	0.86	0.76	0.74
36	0.42	0.51	1.92	2.03	0.85	0.84	0.73	0.70
48	0.44	0.46	2.03	2.16	0.82	0.81	0.70	0.67
Hammer Odde Lighthouse (WMO 06197) :								
0	0.31	0.19	1.24	1.24	0.92	0.91	0.90	0.90
6	0.22	0.12	1.28	1.33	0.91	0.90	0.88	0.87
12	0.24	0.13	1.34	1.42	0.90	0.88	0.87	0.85
18	0.25	0.15	1.38	1.48	0.89	0.87	0.86	0.84
24	0.26	0.14	1.43	1.57	0.88	0.86	0.84	0.82
36	0.24	0.11	1.53	1.67	0.86	0.84	0.82	0.79
48	0.23	0.10	1.62	1.80	0.85	0.81	0.80	0.77

Table 5 summarizes verification results for 10m wind forecasts for the 3km-resolution S03 model and the ensemble mean of the 5km-resolution S05 model for two Danish coastal stations in the western part of the Baltic Sea. A comparison between the two model forecasts shows a small positive bias and RMS errors increasing with forecast range from approx. 1 ms^{-1} to approximately 2 ms^{-1} for 48h forecasts. The error of

276 the ensemble mean forecasts generally increases less with forecast range than the error of the high-
 277 resolution forecasts. Similarly, the correlation and the hit rate ($\text{error} \leq 2 \text{ ms}^{-1}$) decrease with forecast
 278 range, but less so for the ensemble mean forecasts. That is, in terms of wind forcing the ensemble mean of
 279 the S05 model provides slightly more accurate forecasts than the higher resolution, deterministic S03
 280 model, especially for the longer forecast ranges.

281 6 Verification of forecasted SWH against observations

282 6.1 Deterministic measures



284 Figure 5 Scatter plot of 24 h forecasts and corresponding observations of significant wave height at site
 285 Bothnian Sea for the LOW, LOWENSMEAN and HIGH forecast classes. Dotted line is the diagonal,
 286 representing a 1:1 agreement between observations and model.

288 To get an idea of the overall quality of the forecasts, [Figure 5](#) shows scatter plots between 24 h
 289 forecasted and observed SWH for station Bothnian Sea. The points are distributed along the diagonal in all
 290 three configurations with correlation coefficients above 0.9. The RMSE is 0.33 m for both LOW and HIGH
 291 but is lower at 0.29 m for the LOWENSMEAN forecasts, which also have the numerically lowest bias. Also
 292 for other sites, such as Arkona WR (see [Figure 6](#)), the RMSE for LOWENSMEAN forecasts is
 293 lower than for the LOW and HIGH forecasts, and similarly for the bias. However, the scatter plot appears
 294 differently for this station, because there is a tendency for over-predicting high waves for all three forecast
 295 classes.

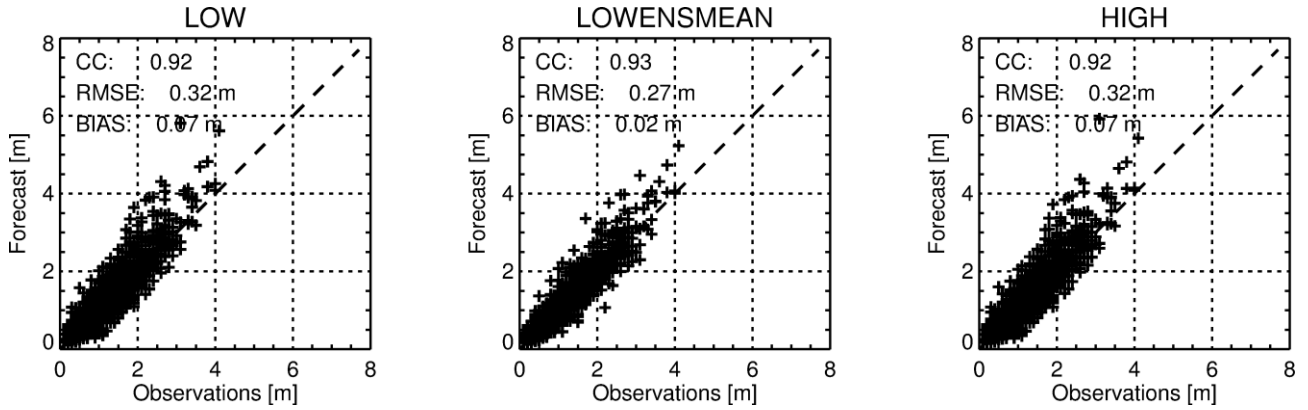


Figure 6 As Figure 5 but for site Arkona WR. Scatter plot of 24 h forecasts and corresponding observations of significant wave height at site Arkona WR for the LOW, LOWENSMEAN and HIGH forecast classes. Dotted line is the diagonal, representing a 1:1 agreement between observations and model.

We now turn to the RMSE as function of forecast range, of which plots for all sites can be found in Figure S2. For all sites, the RMSE increases slightly as function of forecast range. All sites except Vahemadal exhibit qualitatively similar **behaviour**: the RMSE for the LOW and HIGH forecasts are almost similar, while it is lower for the LOWENSMEAN forecasts. Thus, for Arkona WR (shown in Figure 7), Bothnian Sea and Darss Sill WR, the RMSE of the LOW and the HIGH forecasts have overlapping confidence bands. The RMSE for LOWENSMEAN gradually diverges to a lower value (around 5 cm) and for large forecast ranges, the confidence bands do not overlap with those for the LOW and HIGH forecast classes. The remaining sites except Vahemadal behave **similarly**, but with overlapping confidence bands even for the largest forecast ranges.

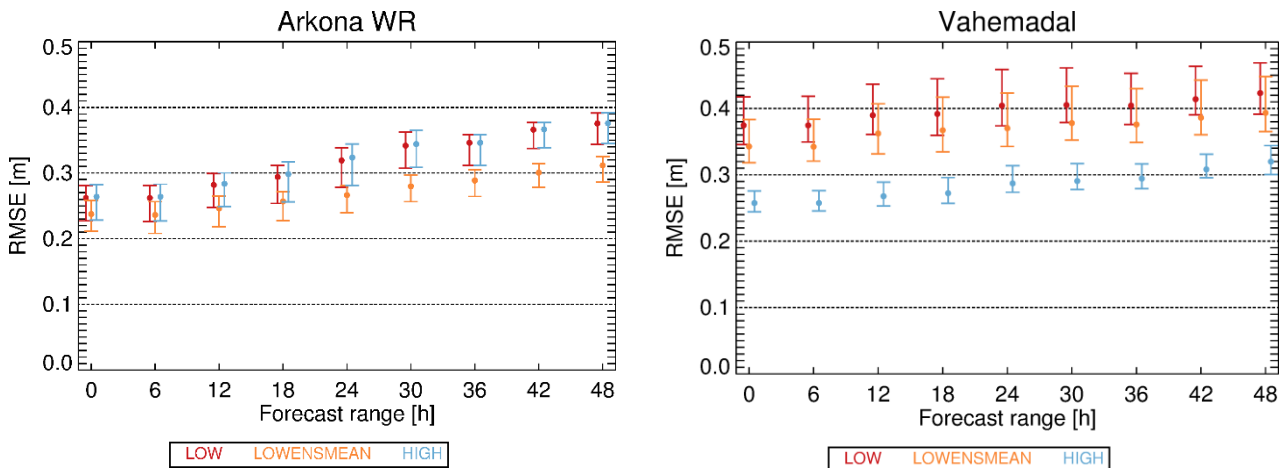
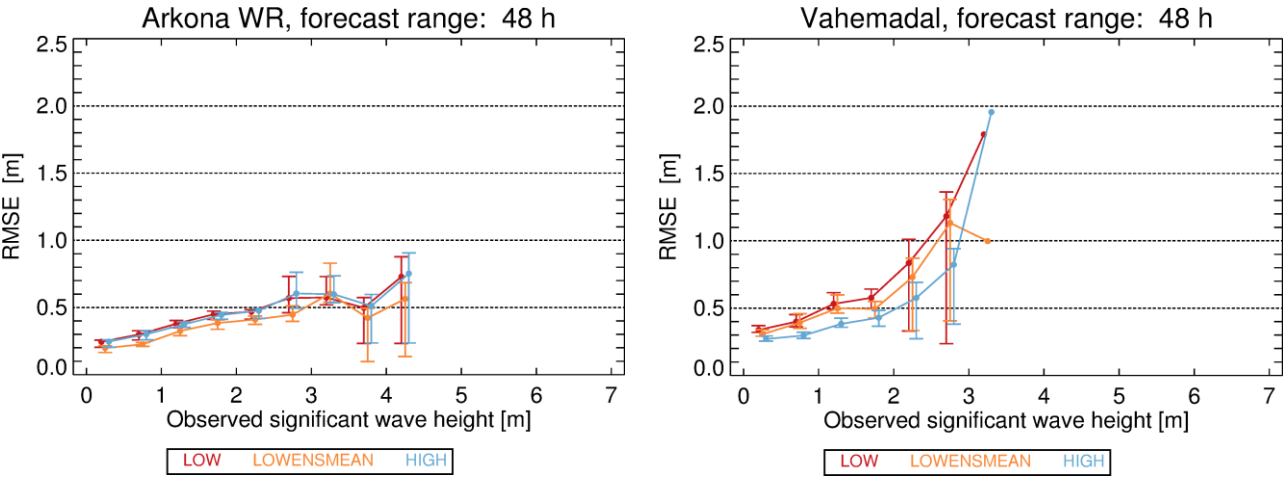


Figure 7 RMSE for selected forecast ranges for Arkona WR (left panel) and Vahemadal (right panel) for LOW, LOWENSMEAN and HIGH forecasts. Error bars show 5/95% confidence bands calculated by bootstrapping.

The site Vahemadal (Figure 7) has a different behavior. For this site, the HIGH forecast class has a significantly smaller RMSE and with non-overlapping confidence bands with the RMSE of the LOW and LOWENSMEAN forecasts. This site also has a non-negligible bias of around 12 cm for the HIGH and around 20 cm for the LOW and LOWENSMEAN forecasts; this bias is independent of forecast range (not shown).

318 **6.1.1 Performance depending on observed SWH**
 319



320 **Figure 8 RMSE as function of SWH for Arkona WR (left panel) and Vahemadal (right panel) for LOW, LOWENSMEAN and HIGH**
 321 **forecasts and forecast range 48 h. Error bars show 5/95% confidence bands calculated by bootstrapping.**

322 The RMSE of the forecasts depends on the magnitude of the SWH. Plots for all sites for [the 24 h](#) and 48 h
 323 forecast ranges of RMSE as function of the SWH can be found in Figures S3 and S4. The RMSE for Arkona
 324 WR and Vahemadal as [a](#) function of the SWH for [the](#) forecast range 48 h is shown in [Figure 8Figure 8Figure](#)
 325 [7](#). The RMSE increases as [a](#) function of the observed SWH for both sites. For Arkona WR, the LOWENSMEAN
 326 forecast class has the lowest RMSE, although with confidence bands overlapping with the other forecast
 327 classes. This behavior is seen at all sites, except Vahemadal. For Vahemadal, the HIGH forecast class has the
 328 lowest RMSE, and up to a SWH of 2 m, the confidence band is well separated from the confidence bands of
 329 the other forecast classes.

330 Also the bias depends on the SWH. Plots for all sites for 24 and 48 h forecast range of the bias as function of
 331 the SWH are displayed in Figures S5 and S6. For small SWH, the bias is close to zero for most sites. For some
 332 sites, the bias remains close to zero for increasing SWH, as shown for Arkona WR in left panel of [Figure](#)
 333 [9Figure 9Figure 8](#), while for others it becomes different from zero for large values of SWH. There [re](#) is no
 334 noticeable different in the bias of the different forecast classes, except for Vahemadal, shown in right panel
 335 of [Figure 9Figure 9Figure 8](#), where the HIGH forecast class has a significantly smaller [under-prediction](#) bias
 336 than the other forecast classes.

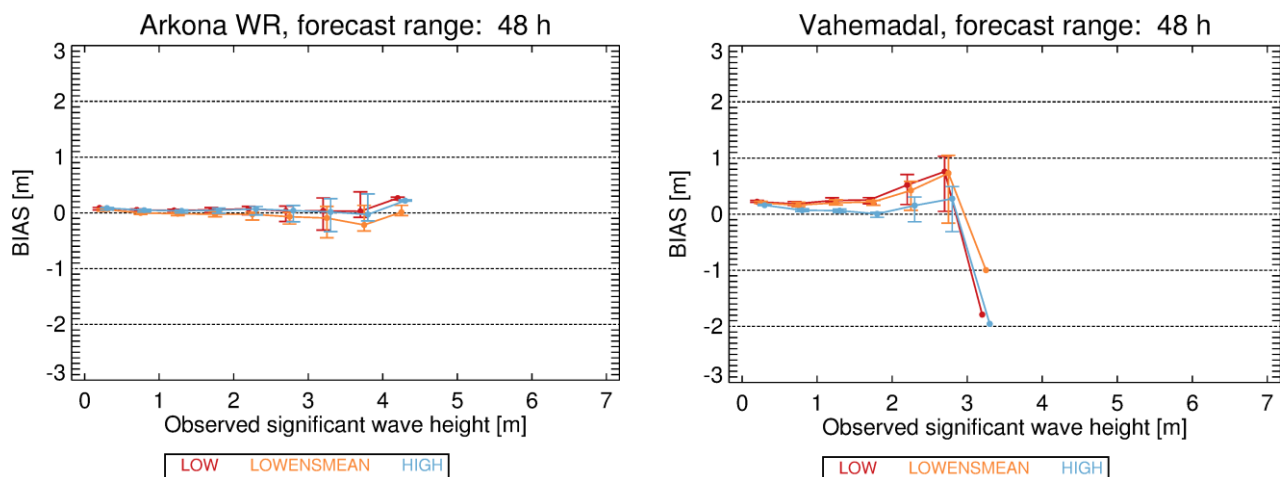


Figure 9 Bias as function of SWH for Arkona WR (left panel) and Vahemadal (right panel) for LOW, LOWENSMEAN and HIGH forecasts and forecast range 24 h. Error bars show 5/95% confidence bands calculated by bootstrapping.

6.1.2 Forecasts during 'Toini' storm

The Toini storm on 11. January 2017, where a SWH of almost 8.0 m was recorded on at Northern Baltic (Björkqvist et al., 2017a), is within our verification period.

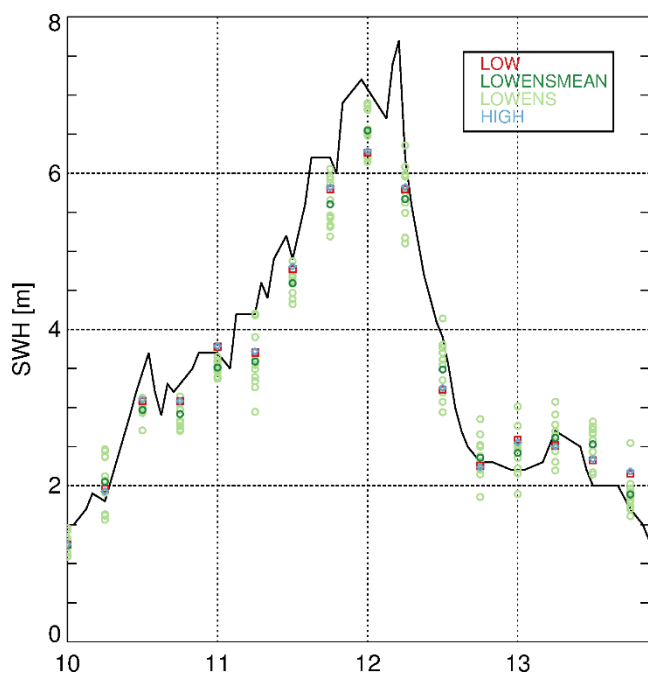


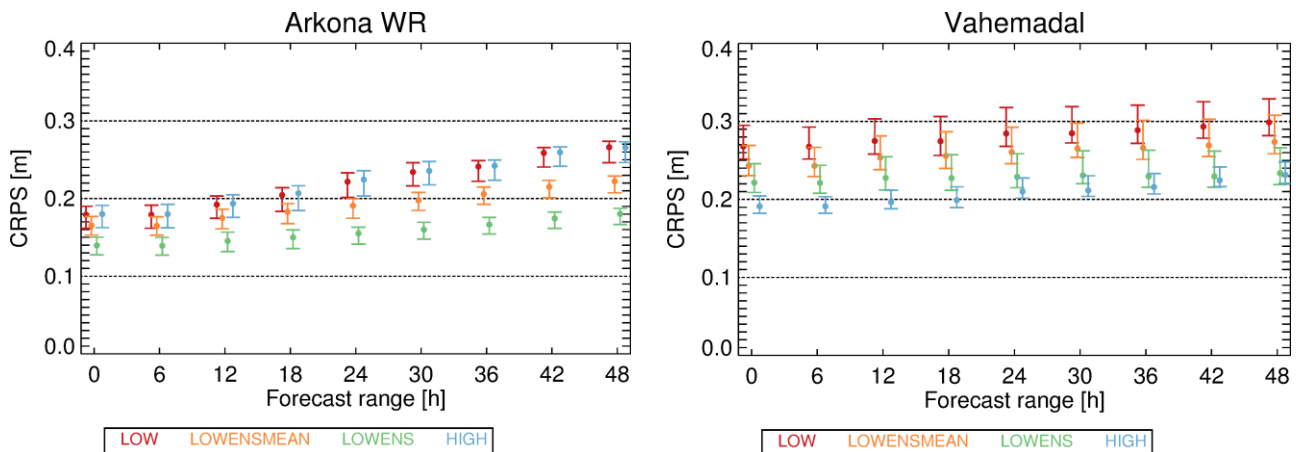
Figure 10 Observed SWH for Northern Baltic during, 10-13 January 2017, including the Toini storm. Open circles are 48 h forecasts.

Figure 10Figure 10Figure 9 shows the observed SWH at Northern Baltic during 10-13 January 2017, i.e. including the Toini storm, peaking in the early hours of 12 January, together with 48 h forecasts. In this case there is no apparent 'best' forecast. Near the peak, LOWENSMEAN performs best, but both before and after, the HIGH/LOW performs better. Furthermore, that in most cases, the LOW and HIGH forecasts are very similar in most cases, indicating that the higher resolution does not improve the forecasts. Finally, we note that the observations generally are within or just a little outside the range of the ensemble forecast.

352 6.2 Probabilistic metrics

353 The 11 ensemble members of the LOWENS forecast class defines a statistical distribution function, which is
 354 a probabilistic forecast of the wave conditions. Besides, the deterministic forecast classes LOW,
 355 LOWENSMEAN and HIGH may be regarded as probabilistic forecasts with probability one for the
 356 deterministically forecasted future state and probability zero for all other states.

357 As described in Section 4, we use CRPS to describe performance of probabilistic forecasts. CRPS for all sites
 358 for selected forecast ranges can be found in Figure S7. As typical examples, [Figure 11](#)[Figure 11](#)[Figure 10](#)
 359 displays this plot for Arkona WR and Vahemadal.



360 **Figure 11** CRPS for selected forecast ranges for Arkona WR (left panel) and Vahemadal (right panel) for LOW, LOWENSMEAN,
 361 LOWENS and HIGH forecasts. Error bars show 5/95% confidence bands calculated by bootstrapping.

362 All sites except Vahemadal behave qualitatively as Arkona WR: the LOWENSMEAN forecast class has a
 363 lower CRPS compared to both the HIGH and LOW classes, although the difference is significant (non-
 364 overlapping confidence bands) for Arkona WR, Bothnian Sea and Darsser Sill WR only, and only for the
 365 largest forecast ranges. Furthermore, for all these sites, the LOWENS forecast class has an even lower CRPS,
 366 with confidence bands separated from those of all other forecasts classes. Again, Vahemadal behaves
 367 differently; here the HIGH forecast class has the best performance in terms of CRPS. However, for large
 368 forecast ranges, the LOWENS forecast class tends to perform equally well.

369 6.3 Binary forecasts

370 For the probabilistic LOWENS forecast class, a binary forecast can be derived as the probability of exceeding
 371 a defined threshold of SWH. For the deterministic forecast classes: LOW, LOWENSMEAN and HIGH, this
 372 probability of exceedance is either zero or one. As described in Section 4, the Brier Score is used as
 373 performance measure for probabilistic, binary forecasts.

374 The Brier Score as a function of threshold is shown for all sites in Figures S7 and S8. [Figure 12](#)[Figure](#)
 375 [12](#)[Figure 11](#) shows the Brier Score as a function of threshold for Arkona WR and Vahemadal for 48 h
 376 forecast range. For Arkona WR, the Brier Score for the LOWENS forecast class is the smallest, however the
 377 confidence intervals overlap with confidence intervals from the other forecasts above the 2 m threshold.
 378 Also the LOWENSMEAN forecast class has a low Brier Score. This behavior is common to all sites except
 379 Vahemadal. For Vahemadal, the Brier Score is smallest for the HIGH forecasts for thresholds above 1 m.

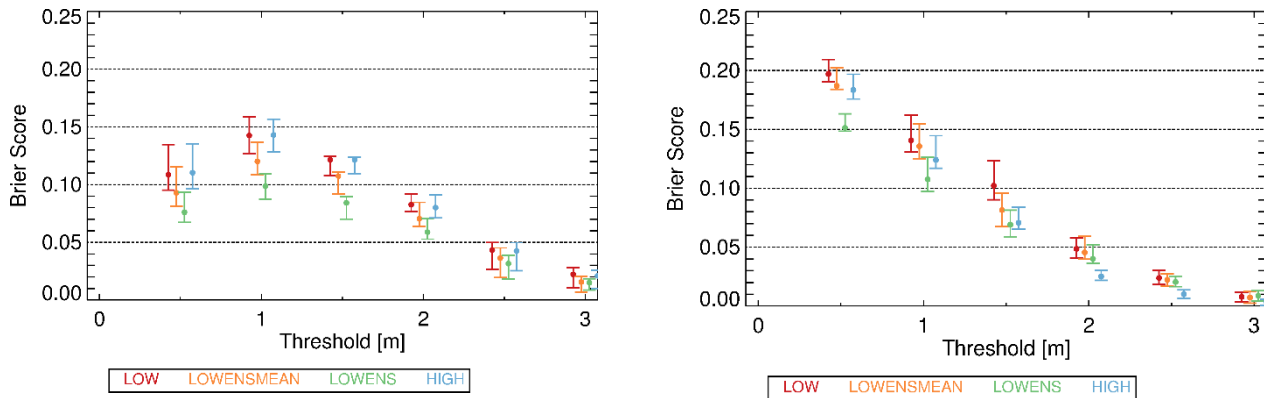


Figure 12 Brier score for Arkona WR (left panel) and Vahemadal (right panel) for binary forecast for forecast range 48 h.

6.4 Rank histogram

Rank histograms serve the purpose of illustrating the reliability of probabilistic ensemble forecasts. It is a histogram of the rank of the observation, when the observation and all ensemble members of the corresponding forecast are pooled together. If the observations and the ensemble members belong to the same distribution, then the rank histogram will be flat, while a U-shaped histogram indicates too small variance within the ensemble members. For more discussion, see Jolliffe and Stephenson (2003).

Rank histograms for all wave measurement sites for forecast range 24 and 48 h are shown in Figure S10 and S11 for forecast range 24 h ~~resp.~~ and 48 h respectively. We note that all histograms show the U-shape, indicating an unrealistically small variance within the ensembles. For most sites the U-shape is symmetric, except for Vahemadal, where the U-shape is strongly asymmetrical. This corresponds well with the bias mentioned in Section 6.1.

7 Discussion

Our main finding in the previous section is that for most wave measurement sites included in this study, the LOWENSMEAN and the LOWENS forecast classes in many cases have a better performance superior to the LOW and HIGH forecast classes. Only for one site results are different; namely that the HIGH forecast class has the superior performance. These conclusions hold, whether based on overall RMSE, CRPS or the Brier score.

In the discussion below, it should be mentioned that improving wave forecasts is not the only driving factor in reducing the grid size of the wave model. Coupling the wave model with atmosphere or ocean circulation models may give a better description of vertical fluxes of heat and momentum (Cavaleri et al., 2012). For instance, Alari et al. (2016) documented a significant improvement of modelled sea-surface temperatures SSTs by the NEMO circulation model in the Baltic Sea when a two-way coupling to the wave model WAM was introduced. Introducing such coupling may demand a high horizontal resolution, in atmosphere, wave and ocean models, in order to describe the fluxes. Doing such couplings may demand a high horizontal resolution to describe the fluxes most satisfactorily.

satisfactorily. Note also that the methodology applied in this study is a site-specific verification and inter-comparison of the different forecast families. This is a valid approach, since most applications of uses of the wave forecasts are site-specific. However, it must be remembered, that the approach has a risk of under-

estimating the overall performance due to double-counting errors in both space and time. We have made no attempt to assess the magnitude of this potential effect.

7.1 Comparison with other operational forecast systems

Multi-year verification results from two operational deterministic wave forecast systems [that covers the region in focus](#) have been published, and can be compared to results from the present study. Both these systems are based on the third generation WAM; the system described in (Tuomi et al., 2008) has about 22 km horizontal resolution, while the system described in (Tuomi et al., 2017) has 1 naut. mile horizontal resolution.

For certain sites, the RMSE of the 6 hour forecasts of SWH are available for at least one of the aforementioned forecast systems in addition to the DMI-WAM forecasts; thus comparison of the systems is possible. All sites have a water depth of more than 46 m and therefore represent offshore conditions.

Table 5 Comparison of RMSE for SWH of 6h forecast runs for selected sites. FIMR values are from (Tuomi et al., 2008) and FMI values are from (Tuomi et al., 2017)

	FIMR	FMI	DMI LOW	DMI LOWENSMEAN	DMI HIGH
Horizontal resolution WAM	~ 22 km	1 naut. mile	10 km	10 km	5 km
Horizontal resolution NWP	~ 22 km	2.5 km	3 km	5 km	3 km
Arkona WR	-	0.28	0.26	0.24	0.26
Bothnian Sea	-	0.28	0.25	0.23	0.25
Finngrundet WR	-	0.27	0.24	0.22	0.23
Helsinki Buoy	0.25	0.26	-	-	-
Northern Baltic	0.31	0.26	0.24	0.23	0.24

Fore completeness, wWe remind the reader that the cases compared in Table 5Table 5Table 6 have different wind forcing and probably also different version of WAM. Therefore the figures cannot be directly compared and differences seen cannot with certainty be attributed to differences in horizontal resolution.

From [Table 5Table 5Table 6](#) one can see that for the sites considered, the LOWENSMEAN has the lowest RMSE. This supports the finding of this study that for offshore conditions, there is no reason to improve the resolution further than that of the LOW configuration. In addition, the results emphasize the value of describing the uncertainties of in the atmospheric forcing by introducing ensembles, as this leads to a lower RMSE of the forecasts. This is also in line with our findings in the previous section.

433 Test runs of a few months duration of deterministic and ensemble wave forecasts of SWH for the Baltic Sea
434 (Behrens, 2015) also show ~~s~~ slight improvement of ensemble mean forecasts, compared to deterministic
435 forecasts, and thus support ~~s~~ our findings.

436 ~~For completeness, we remind the reader that the cases compared in Table 6 have different wind forcing~~
437 ~~and probably also different version of WAM. Therefore the differences seen cannot with certainty be~~
438 ~~attributed to differences in horizontal resolution.~~

439 7.2 Limitations of the study

440 7.2.1 Length of verification period

441 Operational centers typically renew their computer installations every 5-6 years with about an order of
442 magnitude increase in performance. At DMI, a new installation was introduced ~~prime-early~~ 2016, allowing
443 the HIGH and LOWENS configurations to replace the LOW configuration. Presently (~~medio-mid~~-2018) the
444 system is mid-term upgraded and this makes it appropriate to do the inter-~~com~~parison now as a guidance
445 for any changes in the operational setup.

446 ~~Thus~~For this reason, the operational forecasts performed on the present system, supplemented by
447 delayed-mode forecasts has determined the three-year verification period used in our study. A longer
448 verification period could evidently have reduced the sampling uncertainty in the analyses and thereby
449 sharpened the conclusions. On the other hand, the three-year verification is not short compared to ~~other~~
450 ~~the study ies, e.g. by~~ Bunney and Saulter (2015) or the CMEMS verification report by Tuomi et al.(2017)

451 7.2.2 Choice of observational base

452 The present verification is based on observations ~~in-at~~ near-hourly resolution from a number of sites in the
453 Baltic Sea. Therefore, in the major ~~part~~sity of the Baltic Sea, verification is not possible, which ~~puts a~~ limits
454 ~~on how strong the firmness of our~~ conclusions ~~can be made~~.

455 SWH derived from satellite-borne altimeters (Kudryavtseva and Soomere, 2016) offers an alternative,
456 which could be pursued in a future study. These data ~~has~~ have a fair spatial ~~data~~ coverage but at the cost
457 of a temporal resolution of one day or less. This means that maximum wave heights connected to severe
458 storms may easily be missed. Nevertheless, these data has proven useful for verification in the Baltic Sea by
459 (Tuomi et al., 2011)

460 7.3 Effect of sea ice coverage

461 The main effect of sea ice on formation of waves is to limit the fetch. Furthermore, when a developed wave
462 field approach ~~es~~ an ice-covered area, the wind and the waves decouple, so that the waves act more like
463 swell, propagating through ice-covered areas while losing energy by breaking up the ice cover. The WAM
464 model does not account for such interactions, and sea ice, when dense enough, act ~~s~~ as a solid shield that
465 effectively remove ~~s~~ all local wave energy in the model. It is implicitly assumed that dense ice will also be
466 thick enough for this to be approximately correct. In the Baltic Sea, that may not always be the case, and
467 therefore sea ice occurrence may represent a systematic error source in the present study. Another effect
468 of sea ice in the Baltic is that the wave observing systems are withdrawn ~~;~~ when ice is expected. This may
469 cause a systematic bias in the verification analysis ~~;~~ if strong winds during winter are left out.

Based on Copernicus sea ice charts produced by the Finnish Meteorological ~~institute~~Institute the ice conditions for the Baltic have been evaluated. The Finnish ice charts are produced on a grid of approximately 1 km^2 with a temporal resolution of approximately one day in the ice season. Data is available from 2010 onwards. The average ice conditions for February for all years and the three years in focus can be found in Figure S12. All three years 2015-2017, and in particular 2015, have a smaller ice cover relative to the period 2010-2018.

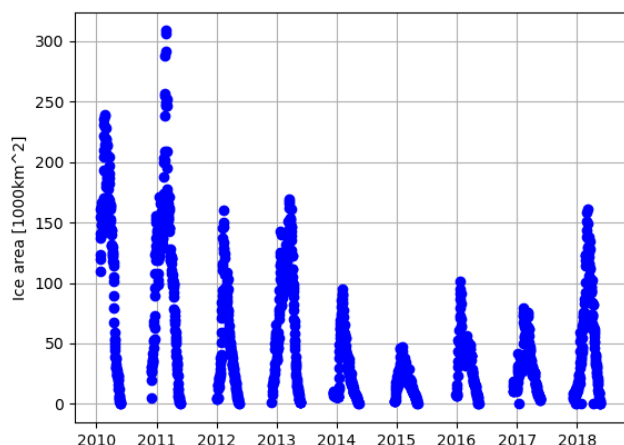


Figure 13 Integrated sea ice area of the Baltic Sea based on Finnish ice charts

Another way to illustrate this is considering the Baltic Sea integrated sea ice area, depicted in Figure 13, which shows that the years 2015-2017 have the lowest sea ice area over the whole period 2010-2018. Therefore, we may anticipate that systematic errors arise from the occurrence of sea ice are relatively small.

8 Conclusion

For most sites, we find that the HIGH forecast class does not perform ~~superior to~~better than the LOW forecast class in forecasting SWH. These sites are all positioned well away from coasts in deep water and are thus freely exposed from all directions. This suggests that the resolution of the bathymetry and the spectral resolution are adequate. For these offshore sites, introducing ensembles increases the performance of the forecasts, whether as in the LOWENSMEAN deterministic forecasts ~~or in~~and the LOWENS probabilistic forecasts. A similar conclusion generally holds for the binary forecast of exceeding a threshold.

For one site, Vahemedal just outside Tallin, the HIGH forecast class performs better than the other classes. The bathymetry near Vahemedal is complex and relatively shallow, thus the bathymetry affects the wave field and an improved description will therefore improve the ~~modeled~~modelled wave field. - Further verification with near-coast stations may reveal whether this conclusion ~~holds in general~~is general for coastal areas.

For high wave heights, there are significant systematic biases for most sites shared among all three forecast configurations. These are most probably to be ascribed to model deficiencies and act to mask any

differences in performance between the different forecast classes. Also the RMSE becomes large for large observed SWH. This is expected since small timing errors in the predicted wave time series will have larger impacts on the model-observation match-up when the SWH is large. The present study therefore suggests that for offshore conditions, there are no indications ~~that a~~ further increase of the resolution of the WAM model will result in enhanced forecast performance. In addition, the results show that introducing ensembles increases the performances. This is both true for deterministic forecast in the form of ensemble mean and for probabilistic forecast.

For nearshore conditions conclusions are based on only one site, but results from this indicates that increasing the resolution gives better forecasts, while introducing ensembles does not. This can be due to both enhanced spatial resolution, allowing a better representation of shadow and shallow water effects, and/or spectral resolution.

The results of the present study thus underpins that a wave model setup with an equidistant grid cannot deliver optimal wave forecasts for both coastal and offshore conditions. This is particularly true for the Baltic Sea, where very small spatial scales are found in the archipelago near the coasts of Sweden and Finland (Björkqvist et al., 2017b). Besides implementing a 0.1 naut. miles model, these authors improved forecasts by introducing semi-empirical modifications to the wave model. ~~The issue is described in Cavalieri et al. (2018), (2018) also write about this and discuss where other approaches are discussed.~~ These include one-way nesting, used in the present study (see Section 2), multi-cell grids (Bunney and Saulter, 2015), and triangular unstructured grids (e.g. Zijlema, 2010). These techniques may be worth testing for the Baltic Sea.

Finally, we note the under-spread in the ensemble forecasts demonstrated in Section 6.4. This points to a potential for improving the combined weather-wave system.

Data availability. Model data is available from the authors upon request, whereas wave observations can be found on the CMEMS server.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. This work was carried out under the EfficienSea2 project and supported by European Union's Horizon 2020 Research and Innovation Programme under Grant Agreement No. 636329. Observational data were kindly provided by Copernicus Marine Environmental Monitoring System (CMEMS). We thank two anonymous referees for valuable suggestions- and dr. Ruth Mottram for help with the English language.

References

Alari, V., Staneva, J., Breivik, Ø., Bidlot, J.-R., Mogensen, K. and Janssen, P.: Surface wave effects on water temperature in the Baltic Sea: simulations with the coupled NEMO-WAM model, Ocean Dyn., 66(8), 917–930, 2016.

533 Alves, J.-H. G., Wittmann, P., Sestak, M., Schauer, J., Stripling, S., Bernier, N. B., McLean, J., Chao, Y.,
534 Chawla, A., Tolman, H. and others: The NCEP–FNMOCC combined wave ensemble product: Expanding
535 benefits of interagency probabilistic forecasts to the oceanic environment, *Bull. Am. Meteorol. Soc.*, 94(12),
536 1893–1905, 2013.

537 Amante, C. and Eakins, B. W.: ETOPO1 1 ARC-MINUTE GLOBAL RELIEF MODEL: PROCEDURES, DATA
538 SOURCES AND ANALYSIS., 2009.

539 Battjes, J. A. and Janssen, J. P. F. M.: Energy Loss and Set-Up Due to Breaking of Random Waves, in *Coastal
540 Engineering 1978.*, 1978.

541 Behrens, A.: Development of an ensemble prediction system for ocean surface waves in a coastal area,
542 *Ocean Dyn.*, 65(4), 469–486, doi:10.1007/s10236-015-0825-y, 2015.

543 Björkqvist, J.-V., Tuomi, L., Tollman, N., Kangas, A., Pettersson, H., Marjamaa, R., Jokinen, H. and Fortelius,
544 C.: Brief communication: Characteristic properties of extreme wave events observed in the northern Baltic
545 Proper, Baltic Sea, *Nat. Hazards Earth Syst. Sci.*, 17(9), 1653–1658, doi:10.5194/nhess-17-1653-2017, 2017a.

546 Björkqvist, J.-V., Tuomi, L., Fortelius, C., Pettersson, H., Tikka, K. and Kahma, K. K.: Improved estimates of
547 nearshore wave conditions in the Gulf of Finland, *J. Mar. Syst.*, 171, 43–53,
548 doi:10.1016/j.jmarsys.2016.07.005, 2017b.

549 Bunney, C. and Saulter, A.: An ensemble forecast system for prediction of Atlantic–UK wind waves, *Ocean
550 Model.*, 96, 103–116, doi:10.1016/j.ocemod.2015.07.005, 2015.

551 Cao, D., Tolman, H. L., Chen, H. S., Chawla, A. and Gerald, V. M.: Performance of the ocean wave ensemble
552 forecast system at NCEP, in *The 11th International Workshop on Wave Hindcasting & Forecasting and 2nd
553 Coastal Hazards Symposium*. [online] Available from:
554 <http://nopp.ncep.noaa.gov/mmab/papers/tn279/mmab279.pdf> (Accessed 22 September 2017), 2009.

555 Carrasco, A. and Saetra, Ø.: A limited-area wave ensemble prediction system for the Nordic Seas and the
556 North Sea, *Norwegian Meteorological Institute.*, 2008.

557 Cavaleri, L., Fox-Kemper, B. and Hemer, M.: Wind Waves in the Coupled Climate System, *Bull. Am.
558 Meteorol. Soc.*, 93(11), 1651–1661, doi:10.1175/BAMS-D-11-00170.1, 2012.

559 Cavaleri, L., Abdalla, S., Benetazzo, A., Bertotti, L., Bidlot, J.-R., Breivik, Ø., Carniel, S., Jensen, R. E., Portilla-
560 Yandun, J., Rogers, W. E., Roland, A., Sanchez-Arcilla, A., Smith, J. M., Staneva, J., Toledo, Y., van Vledder, G.
561 P. and van der Westhuysen, A. J.: Wave modelling in coastal and inner seas, *Prog. Oceanogr.*,
562 doi:10.1016/j.pocean.2018.03.010, 2018.

563 Günther, H., Hasselmann, S. and Janssen, P. A. E. M.: The WAM Model cycle 4, *World Data Center for
564 Climate (WDCC) at DKRZ.*, 1992.

565 Hasselmann, K., Barnett, T., Bouws, E., Carlson, H., Cartwright, D., Enke, K., Ewing, J., Gienapp, H.,
566 Hasselmann, D., Kruseman, P., Meerburg, A., Müller, P., Olbers, D., Richter, K., Sell, W. and Walden, H.:
567 Measurements of wind-wave growth and swell decay during the {Joint North Sea Wave Project}, *Deut
568 Hydrogr Z.*, 8(12), 1–95, 1973.

569 Jolliffe, I. T. and Stephenson, D. B.: *Forecast verification: a practitioner’s guide in atmospheric science*, John
570 Wiley & Sons., 2003.

571 Kudryavtseva, N. A. and Soomere, T.: Validation of the multi-mission altimeter wave height data for the
572 Baltic Sea region, *Est. J. Earth Sci.*, 65(3), 161, doi:10.3176/earth.2016.13, 2016.

573 Saetra, Ø. and Bidlot, J.-R.: Assessment of the ECMWF Ensemble Prediction Sytem for Waves and Marine
574 Winds, European Centre for Medium-Range Weather Forecasts., 2002.

575 Schaffer, J., Timmermann, R., Arndt, J. E., Kristensen, S. S., Mayer, C., Morlighem, M. and Steinhage, D.: A
576 global, high-resolution data set of ice sheet topography, cavity geometry, and ocean bathymetry, *Earth*
577 *Syst. Sci. Data*, 8(2), 543–557, doi:10.5194/essd-8-543-2016, 2016.

578 She, J., Allen, I., Buch, E., Crise, A., Johannessen, J. A., Le Traon, P.-Y., Lips, U., Nolan, G., Pinardi, N.,
579 Reißmann, J. H., Siddorn, J., Stanev, E. and Wehde, H.: Developing European operational oceanography for
580 Blue Growth, climate change adaptation and mitigation, and ecosystem-based management, *Ocean Sci.*,
581 12(4), 953–976, doi:10.5194/os-12-953-2016, 2016.

582 Tuomi, L., Kangas, A., Leinonen, J. and Boman, H.: The Accuracy of FIMR Wave Forecasts in 2002-2005.,
583 2008.

584 Tuomi, L., Kahma, K. K. and Pettersson, H.: Wave hindcast statistics in the seasonally ice-covered Baltic sea.,
585 *Boreal Environ. Res.*, 16, 2011.

586 Tuomi, L., Vähä-Piikkiö, O. and Alari, V.: Baltic Sea Wave Analysis and Forecasting Product
587 BALTICSEA_ANALYSIS_FORECAST_WAV_003_010, 2017.

588 Zijlema, M.: Computation of wind-wave spectra in coastal waters with SWAN on unstructured grids, *Coast.*
589 *Eng.*, 57(3), 267–277, doi:10.1016/j.coastaleng.2009.10.011, 2010.

590