

1 An approach to the verification of high-
2 resolution ocean models using spatial methods

3 Ric Crocker¹, Jan Maksymczuk², Marion Mittermaier¹, Marina Tonani², Christine Pequignet²

4 ¹Verification, Impacts and Post-Processing, Weather Science, Met Office, Exeter, EX1 3PB, United Kingdom

5 ²Ocean Forecasting Research & Development, Weather Science, Met Office, Exeter, EX1 3PB, United Kingdom

6

7

8 *Corresponding author:* ric.crocker@metoffice.gov.uk

9 **Abstract**

10 The Met Office currently runs two operational ocean forecasting configurations for the North
11 West European Shelf, an eddy-permitting model with a resolution of 7 km (AMM7), and an eddy-
12 resolving model at 1.5 km (AMM15).

13 Whilst qualitative assessments have demonstrated the benefits brought by the increased
14 resolution of AMM15, particularly in the ability to resolve finer-scale features, it has been difficult
15 to show this quantitatively, especially in forecast mode. Application of typical assessment metrics
16 such as the root mean square error have been inconclusive, as the high-resolution model tends
17 to be penalised more severely, referred to as the double-penalty effect. This effect occurs in
18 point-to point comparisons whereby features correctly forecast but misplaced with respect to
19 the observations are penalised twice; once for not occurring at the observed location, and
20 secondly for occurring at the forecast location, where they have not been observed.

21 An exploratory assessment of sea surface temperature (SST) has been made at in-situ
22 observation locations using a single-observation-neighbourhood-forecast (SO-NF) spatial
23 verification method known as the High-Resolution Assessment (HiRA) framework. The primary
24 focus of the assessment was to capture important aspects of methodology to consider when
25 applying the HiRA framework. Forecast grid points within neighbourhoods centred on the
26 observing location are considered as pseudo ensemble members, so that typical ensemble and
27 probabilistic forecast verification metrics such as the Continuous Ranked Probability Score (CRPS)
28 can be utilised. It is found that through the application of HiRA it is possible to identify
29 improvements in the higher resolution model which were not apparent using typical grid scale
30 assessments.

31 This work suggests that future comparative assessments of ocean models with different
32 resolutions would benefit from using HiRA as part of the evaluation process, as it gives a more
33 equitable and appropriate reflection of model performance at higher resolutions.

34 **Keywords**

35 verification, ocean forecasts, SST, spatial methods, neighbourhood, validation, double-penalty

36

37 1. Introduction

38 When developing and improving forecast models an important aspect is to assess whether model
39 changes have truly improved the forecast. Assessment can be a mixture of subjective approaches,
40 such as visualising forecasts and assessing whether the broad structure of a field is appropriate,
41 or objective methods, comparing the difference between the forecast and an observed or
42 analysed value of 'truth' for the model domain.

43 Different types of intercomparison can be applied to identify different underlying behaviours:

- 44 • between different forecasting systems over an overlapping region to check for model
45 consistency between the two;
- 46 • between two versions of the same model to test the value of model upgrades prior to
47 operational implementation;
- 48 • parent-son intercomparison, evaluating the impact of downscaling or nesting of models;
- 49 • a forecast comparison against reanalysis of the same model, inferring the effect of
50 resolution and forcing, especially in coastal areas.

51 There are a number of works which have used these types of assessment to delve into the
52 characteristics of forecast models (e.g. Aznar et al., 2015, Mason et al., 2019, Juza et al., 2015)
53 and produce coordinated validation approaches (Hernandez et al., 2015).

54 To aid the production of quality model assessment, services exist which regularly produce multi-
55 model assessments to deliver to the ocean community (e.g. Lorente et al., 2019a)

56 One of the issues faced when assessing high-resolution models against lower resolution models
57 over the same domain is that often the coarser model appears to perform at least equivalently
58 or better when using typical verification metrics such as root-mean-squared-error (RMSE) or
59 mean error, which is a measure of the bias. Whereas a higher resolution model has the ability
60 and requirement to forecast greater variation, detail and extremes, a coarser model cannot
61 resolve the detail and will, by its nature, produce smoother features with less variation resulting
62 in smaller errors. This can lead to the situation that despite the higher resolution model looking
63 more realistic it may verify worse (e.g. Mass et al., 2002, Tonani et al., 2019).

64 This is particularly the case when assessing forecast models categorically. If the location of a
65 feature in the model is incorrect then two penalties will be accrued, one for not forecasting the
66 feature where it should have been and one for forecasting the same feature where it did not
67 occur (the double penalty effect, e.g. Rossa et al., 2008). This effect is more prevalent in higher-
68 resolution models due to their ability to, at least, partially resolve smaller-scale features of
69 interest. If the lower resolution model could not resolve the feature, and therefore did not
70 forecast it, that model would only be penalised once. Therefore, despite giving potentially better
71 guidance the higher resolution model will verify worse.

72 Yet, the underlying need to quantitatively show the value of high-resolution led to the
73 development of so-called “spatial” verification methods which aimed to account for the fact the
74 forecast produced realistic features that were not necessarily at the right place or at quite the
75 right time (e.g. Ebert, 2008 or Gilleland, 2009). These methods have been in routine use within
76 the atmospheric model community for a number of years with some long-term assessments and
77 model comparisons (e.g. Mittermaier *et al.* 2013 for precipitation).

78 Spatial methods allow forecast models to be assessed with respect to several different types of
79 focus. Initially these methods were classified into four groups. Some methods look at the ability
80 to forecast specific features (e.g. Davis et al., 2006), some look at how well the model performs
81 at different scales (scale-separation, e.g. Casati et al., 2004). Others look at field deformation
82 (how much a field would have to be transformed to match a ‘truth’ field (e.g. Keil and Craig,
83 2007). Finally, there is neighbourhood verification, many of which are equivalent to low band-
84 pass filters. In these methods forecasts are assessed at multiple spatial or temporal scales to see
85 how model skill changes as the scale is varied.

86 Dorninger et al. (2018) provides an updated classification of spatial methods, suggesting a fifth
87 class of methods, known as distance metrics, which sit between field deformation and feature-
88 based methods. These methods evaluate the distances between features, but instead of just
89 calculating the difference in object centroids (which is typical), the distances between all grid
90 point pairs are calculated, which makes distance metrics more similar to field deformation
91 approaches. Furthermore, there is no prior identification of features. This makes distance metrics

92 a distinct group that warrants being treated as such in terms of classification. Not all methods
93 are easy to classify. An example of this is the Integrated Ice Edge Error (IIEE) developed for
94 assessing the sea ice extent (Goessling et al., 2016).

95 This paper exploits the use of one such spatial technique for the verification of sea surface
96 temperature (SST), in order to determine the levels of forecast accuracy and skill across a range
97 of model resolutions. The High-Resolution Assessment framework (Mittermaier, 2014,
98 Mittermaier and Csima, 2017) is applied to the Met Office Atlantic Margin Model running at 7 km
99 (O’Dea et al., 2012, O’Dea et al., 2017, King et al., 2018) (AMM7), and 1.5 km (Graham et al.,
100 2018, Tonani et al., 2019) (AMM15) resolutions for the European North West Shelf (NWS). The
101 aim is to deliver an improved understanding beyond the use of basic biases and RMS errors for
102 assessing higher resolution ocean models, which would then better inform users on the quality
103 of regional forecast products. Atmospheric science has been using high-resolution convective-
104 scale models for over a decade, and so have experience in assessing forecast skill on these scales,
105 so it is appropriate to trial these methods on eddy-resolving ocean model data. As part of the
106 demonstration, the paper also looks at how the method should be applied to different ocean
107 areas, where variation at different scales occurs due to underlying driving processes.

108

109 The paper was influenced by discussions on how to quantify the added value from investments
110 in higher resolution modelling given the issues around the double-penalty effect discussed above,
111 which is currently an active area of research within the ocean community (Lorente et al., 2019b,
112 Hernández et al., 2018, Mourre et al., 2019).

113 Section 2 describes the model and observations used in this study along with the method applied.
114 Section 3 presents the results, and section 4 discusses the lessons learnt while using HiRA on
115 ocean forecasts and sets the path for future work by detailing the potential and limitations of the
116 method.

117

118 2. Data and Methods

119 **2.1 Forecasts**

120 The forecast data used in this study are from the two products available in the Copernicus Marine
121 Environment Monitoring Service (CMEMS, see e.g. Le Traon et al., 2019, for a summary of the
122 service) for the North West European Shelf area:

- 123 • NORTHWESTSHELF_ANALYSIS_FORECAST_PHYS_004_001_b (AMM7)
- 124 • NORTHWESTSHELF_ANALYSIS_FORECAST_PHY_004_013 (AMM15)

125 The major difference between these two products is the horizontal resolution, ~7 km for AMM7
126 and 1.5 km for AMM15. Both systems are based on a forecasting ocean assimilation model with
127 tides. The ocean model is NEMO (Nucleus for European Modelling of the Ocean, Madec, 2016),
128 using the 3DVar NEMOVAR system to assimilate observations (Mogensen et al., 2012). These are
129 surface temperature in-situ and satellite measurements, vertical profiles of temperature and
130 salinity, and along track satellite sea level anomaly data. The models are forced by lateral
131 boundary conditions from the UK Met Office North Atlantic Ocean forecast model and by the
132 CMEMS Baltic forecast product BALTICSEA_ANALYSIS_FORECAST_PHY_003_006. The
133 atmospheric forcing is given by the operational European Centre for Medium-Range Weather
134 Forecasts (ECMWF) Numerical Weather Prediction model for AMM15, and by the operational UK
135 Met Office Global Atmospheric model for AMM7.

136

	Resolution	Atmospheric forcing	Geographical model domain	
AMM7	~7 km	MetUM 10 km	40°N - 65°N	20°W -13°E
AMM15	~1.5 km	ECMWF IFS ~14 km	~45°N - 63°N	~20°W - 13°E

137

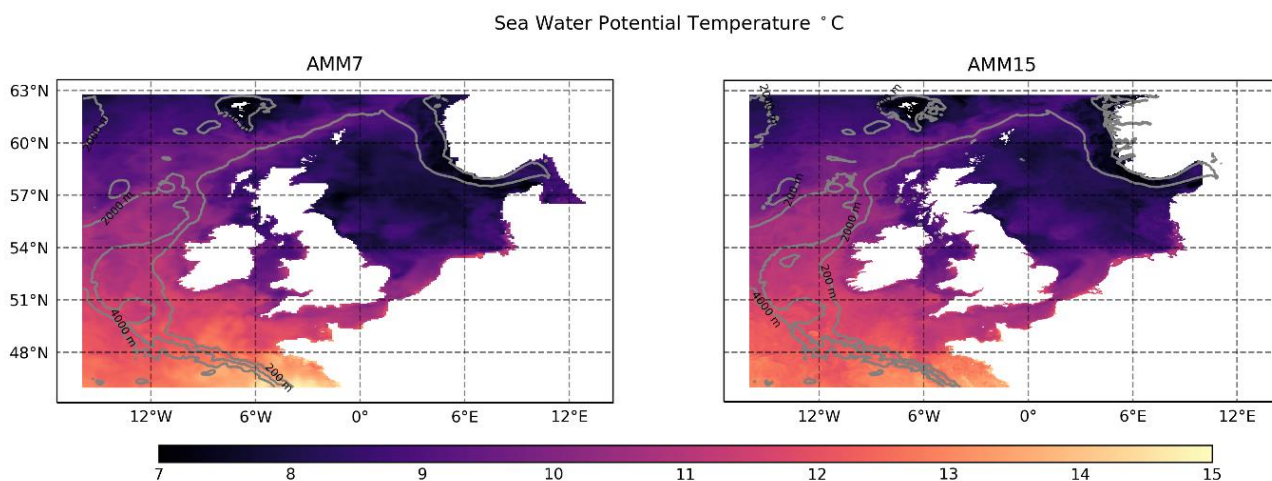
138 Table 1: Summary of the main differences between NORTHWESTSHELF_ANALYSIS_FORECAST_PHYS_004_001_b (AMM7) and
139 NORTHWESTSHELF_ANALYSIS_FORECAST_PHYS_004_013 (AMM15)

140 The AMM15 and AMM7 systems run once a day and provide forecasts for temperature, salinity,
141 horizontal currents, sea level, mixed layer depth, and bottom temperature. Hourly instantaneous
142 values and daily 25-hour, de-tided, averages are provided for the full water column.

143 AMM7 has a regular latitude-longitude grid, whilst AMM15 is computed on a rotated grid and re-
144 gridded to have both models delivered to the (CMEMS) data catalogue

145 (<http://marine.copernicus.eu/services-portfolio/access-to-products/>) on a regular grid. A fuller
146 description of the respective configurations of the two models can be found in Tonani et al.,
147 (2019).

148
149 For the purposes of this assessment the 5-day daily mean sea surface potential temperature (SST)
150 forecasts (with lead times of 12, 36, 60, 84, 108 hours) were utilised for the period from January
151 to September 2019. Forecasts were compared for the co-located areas of AMM7 and AMM15.
152 Figure 1 shows the AMM7 and AMM15 co-located domain along with the land-sea mask for each
153 of the models. AMM15 has a more detailed coastline and SST field than AMM7 due to its higher
154 resolution. When comparing two models with different resolutions it is important to know
155 whether increased detail actually translates into better forecast skill. Additionally, the differences
156 in coastline representation can have an impact on any HiRA results obtained, as will be discussed
157 in a later section.



158
159 *Figure 1 - AMM7 and AMM15 co-located areas. Note the difference in the land-sea boundaries due to the different resolutions,*
160 *notably around the Scandinavian coast. Contours show the model bathymetry at 200, 2000 and 4000 m.*

161
162 It should be noted that this study is an assessment of the application of spatial methods to ocean
163 forecast data, and as such, is not meant as a full and formal assessment and evaluation of the

164 forecast skill of the AMM7 and AMM15 ocean configurations. To this end, a number of
165 considerations have had to be taken into account in order to reduce the complexity of this initial
166 study. Specifically, it was decided at an early stage to use daily mean SST temperatures, as
167 opposed to hourly instantaneous SST, as this avoided any influence of the diurnal cycle and tides
168 on any conclusions made. AMM15 and AMM7 daily means are calculated as means over 25 hours
169 to remove both the diurnal cycle and the tides. The tidal signal is removed because the period of
170 the major tidal constituent, the semidiurnal lunar component M2, is 12 hr and 25 min (Howarth
171 and Pugh, 1983). Daily means are also one of the variables that are available from the majority
172 of the products within the CMEMS catalogue, including reanalysis, so the application of the
173 spatial methods could be relevant in other use cases beyond those considered here. In addition,
174 there are differences in both the source and frequency of the air-sea interface forcing used in
175 both the AMM7 and AMM15 configurations which could influence the results. Most notably, the
176 AMM7 uses hourly surface pressure and 10 m winds from the Met Office Unified Model (UM),
177 whereas the AMM15 uses 3-hourly data from ECMWF.

178 2.2 Observations

179 SST observations used in the verification were downloaded from the CMEMS catalogue from the
180 product

181

- 182 • INSITU_NWS_NRT_OBSERVATIONS_013_036

183

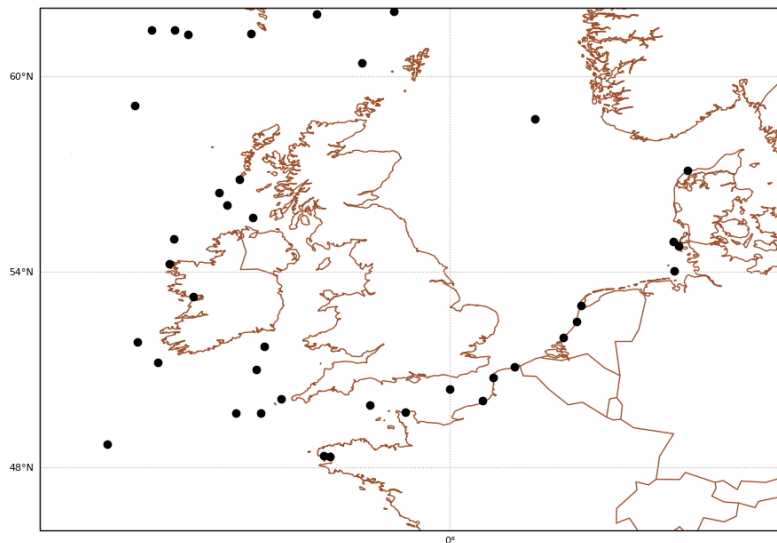
184 This dataset consists of in-situ observations only, including daily drifters, mooring, ferry-box and
185 Conductivity Temperature Depth (CTD) observations. This results in a varying number of
186 observations being available throughout the verification period, with uneven spatial coverage
187 over the verification domain. Figure 2 shows a snapshot of the typical observational coverage, in
188 this case for 1200 UTC 6th June 2019. This coverage is important when assessing the results,
189 notably when thinking about the size and type of area over which an observation is meant to be
190 representative of, and how close to the coastline each observation is.

191

192 This study was set up to detect issues that should be considered by users when applying HiRA
193 within a routine ocean verification set-up, using a broad assessment containing as much data as
194 was available in order to understand the impact of using HiRA for ocean forecasts. Several
195 assumptions were made in this study.

196
197 For example, there is a temporal mismatch between the forecasts and observations used. The
198 forecasts (which were available at the time of this study) are daily means of the SSTs from 00 UTC
199 to 00 UTC, whilst the observations are instantaneous and usually available hourly. For the
200 purposes of this assessment, we have focused on SSTs closest to the mid-point of the forecast
201 period for each day (nominally 12 UTC). Observation times had to be within 90 minutes of this
202 time, with any other times from the same observation site being rejected. A particular reason for
203 picking a single observation time rather than daily averages was so that moving observations,
204 such as drifting buoys, could be incorporated into the assessment. Creating daily mean
205 observations from moving observations would involve averaging reports from different forecast
206 grid- boxes, and hence contaminate the signal that HiRA is trying to evaluate.

207



208

209 *Figure 2 - Observation locations within the domain for 1200 UTC on 6th June 2019.*

210 Future applications would probably contain a stricter set-up, e.g. only using fixed daily mean
211 observations, or verifying instantaneous (hourly) forecasts so as to provide a sub-daily
212 assessment of the variable in question.

213

214 3. High Resolution Assessment (HiRA)

215 The HiRA framework (Mittermaier, 2014) was designed to overcome the difficulties encountered
216 in assessing the skill of high-resolution models when evaluating against point observations.
217 Traditional verification metrics such as RMSE and mean error rely on a precise matching in space
218 and time, by (typically) extracting the nearest model grid point to an observing location. The
219 method is an example of a single-observation-neighbourhood-forecast (SO-NF) approach, with
220 no smoothing. All the forecast grid points within a neighbourhood centred on an observing
221 location are treated as a pseudo ensemble, which is evaluated using well known ensemble and
222 probabilistic forecast metrics. Scores are computed for a range of (increasing) neighbourhood
223 sizes to understand the scale-error relationship. This approach assumes that the observation is
224 representative of not only its precise location but also has characteristics of the surrounding area
225 as well. WMO manual No 8 (2017) suggests that, in the atmosphere, observations can be
226 considered to be representative of an area within a 100 km radius of a land station, but this is
227 often very optimistic. The manual states further: “For small-scale or local applications the
228 considered area may have dimensions of 10 km or less.” A similar principle applies to the ocean,
229 i.e. observations can represent an area around the nominal observation location, though the
230 representative scales are likely to be very different from in the atmosphere. The representative
231 scale for an observation will also depend on local characteristics of the area, for example whether
232 the observation is on the shelf, or in open ocean or likely to be impacted by river discharge.

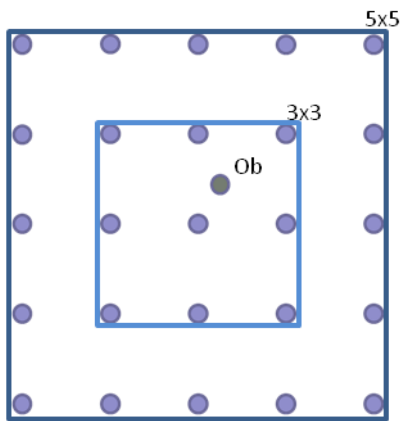
233 There will be a limit to the useful forecast neighbourhood size which can be used when comparing
234 to a point observation. This maximum neighbourhood size will depend on the representative
235 scale of the variable under consideration. Put differently, once the neighbourhoods become too
236 big there will be forecast values in the pseudo ensemble which will not be representative of the
237 observation (and the local climatology) and any skill calculated will be essentially random.

238 Combining results for multiple observations with very different representative scales (for
239 example a mixture of deep ocean and coastal observations) could contaminate results, due to
240 the forecast neighbourhood only being representative of a subset of the observations. The effect
241 of this is explored later in this paper.

242

243 HiRA can be based on a range of statistics, data thresholds and neighbourhood sizes in order to
244 assess a forecast model. When comparing deterministic models of different resolutions, the
245 approach is to equalise on the physical area of the neighbourhoods (i.e. having the same
246 “footprint”). By choosing sequences of neighbourhoods that provide (at least) approximate
247 equivalent neighbourhoods (in terms of area), two or more models can be fairly compared.

248 HiRA works as follows. For each observation, several neighbourhood sizes are constructed,
249 representing the length in forecast grid points of a square domain around the observation points,
250 centred on the grid point closest to the observation (Fig. 3). There is no interpolation applied to
251 the forecast data to bring it to the observation point, all the data values are used unaltered.



252

253 *Figure 3 - Example of forecast grid point selections for different HiRA neighbourhoods for a single observation point. A 3x3 domain*
254 *returns 9 points that represent the nearest forecast grid points in a square around the observation. A 5x5 domain encompasses*
255 *more points.*

256

257 Once neighbourhoods have been constructed, the data can be assessed using a range of well-
258 known ensemble or probabilistic scores. The choice of statistic usually depends on the

259 characteristics of the parameter being assessed. Parameters with significant thresholds can be
260 assessed using the Brier score (Brier, 1950) or the Ranked Probability Score (RPS) (Epstein, 1969),
261 i.e. assessing the ability of the forecast to correctly locate a forecast in the correct threshold
262 band. For continuous variables such as SST, the data has been assessed using the continuous
263 ranked probability score (CRPS) (Brown, 1974, Hersbach, 2000).

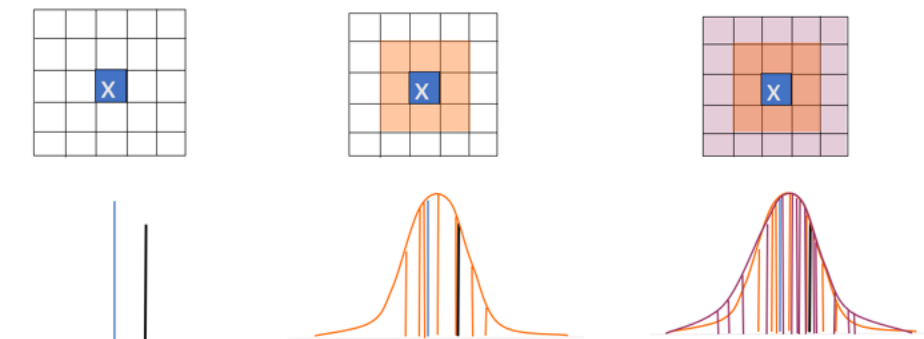
264 The CRPS is a continuous extension of the RPS. Whereas the RPS is effectively an average of a
265 user-defined set of Brier scores over a finite number of thresholds, the CRPS extends this by
266 considering an integral over all possible thresholds. It lends itself well to ensemble forecasts of
267 continuous variables such as temperature and has the useful property that the score reduces to
268 the mean absolute error (MAE) for a single grid point deterministic model comparison. This
269 means that if required, both deterministic and probabilistic forecasts can be compared using the
270 same score.

$$271 \quad CRPS = \int_{-\infty}^{\infty} [P_{fcst}(x) - P_{obs}(x)]^2 dx \quad (1)$$

272

273 Equation (1) defines the CRPS, where for a parameter x , $P_{fcst}(x)$ is the cumulative distribution of
274 the neighbourhood forecast and $P_{obs}(x)$ is the cumulative distribution of the observed value,
275 represented by a Heaviside function (see Hersbach, 2000). The CRPS is an error-based score
276 where a perfect forecast has a value of zero. It measures the difference between two cumulative
277 distributions, a forecast distribution formed by ranking the (in this case quasi) -ensemble
278 members represented by the forecast values in the neighbourhood, and a step function
279 describing the observed state. To use an ensemble, HiRA makes the assumption that all grid
280 points within a neighbourhood are equi-probable outcomes at the observing location. Therefore,
281 aside from the observation representativeness limit, as the neighbourhood sizes increase, this
282 assumption of equi-probability will break down as well, and scores become random. Care must
283 therefore be taken to decide whether a particular neighbourhood size is appropriately
284 representative. This decision will be based on the length scales appropriate for a variable as well
285 as the resolution of the forecast model being assessed. Figure 4 shows a schematic of how

286 different neighbourhood sizes contribute towards constructing forecast probability density
287 functions around a single observation.



288
289 *Figure 4 – Example of how different forecast neighbourhood sizes would contribute to generation of a probability density function*
290 *around an observation (denoted by x). The larger the neighbourhood, the better described the pdf, though potentially at the*
291 *expense of larger spread. Where a forecast point is invalid within the forecast neighbourhood then that site is rejected from the*
292 *calculations for that neighbourhood size.*

293
294 AMM7 and AMM15 resolve different length scale of motion, due to their horizontal resolution.
295 This should be taken into account when assessing the results of different neighbourhood sizes.
296 Both models can resolve the large barotropic scale (~200 km) and the shorter baroclinic scale off
297 the shelf, in deep water. On the continental shelf, only the resolution of ~1.5 km of AMM15,
298 permits motions at the smallest baroclinic scale since the first baroclinic Rossby radius is of order
299 of 4 km (O’Dea et al., 2012). AMM15 represents a step change in representing the eddy dynamics
300 variability on the continental shelf. This difference has an impact also on the data assimilation
301 scheme, where two horizontal correlation length scales (Mirouze et al., 2016) are used to
302 represent large and small scales of ocean variability. The long length scale is 100 km while the
303 short correlation length scale aims to account for internal ocean processes variability,
304 characterized by the Rossby radius of deformation. Computational requirements restrict the
305 short length scale to be at least 3 model grid points, 4.5 km and 21 km respectively for AMM15
306 and AMM7 (Tonani et al., 2019). Although AMM15 resolves smaller scale processes, comparing

307 AMM7 and AMM15 in neighbourhood sizes between the AMM7 resolution and multiples of this
308 resolution will address processes that should be accounted for in both models.

309

310 As the methodology is based on ensemble and probabilistic metrics it is naturally extensible to
311 ensemble forecasts (see Mittermaier and Csima, 2017), which are currently being developed in
312 research-mode by the ocean community, allowing for inter-comparison between deterministic
313 and probabilistic forecast models in an equitable and consistent way.

314

315 4. Model Evaluation Tools (MET)

316 Verification was performed using the Point-Stat tool, which is part of the Model Evaluation Tools
317 (MET) verification package, that was developed by the National Center for Atmospheric Research
318 (NCAR), and which can be configured to generate CRPS results using the HiRA framework. MET is
319 free to download from GitHub at <https://github.com/NCAR/MET>.

320

321 5. Equivalent neighbourhoods and equalisation

322 When comparing neighbourhoods between models, the preference is to look for similar-sized
323 areas around an observation and then transforming this to the closest odd-numbered, square
324 neighbourhood, which will be called the 'equivalent neighbourhood'. In the case of the two
325 models used, the most appropriate neighbourhood size can change depending on the structure
326 of the grid so the user needs to take into consideration what is an accurate match between the
327 models being compared.

328

329 The two model configurations used in this assessment are provided on standard latitude-
330 longitude grids via the CMEMS catalogue. The AMM7 and AMM15 configurations are stated to
331 have resolutions approximating 7 km and 1.5 km respectively. Thus, equivalent neighbourhoods
332 should simply be a case of matching neighbourhoods with similar spatial distances. In fact, the

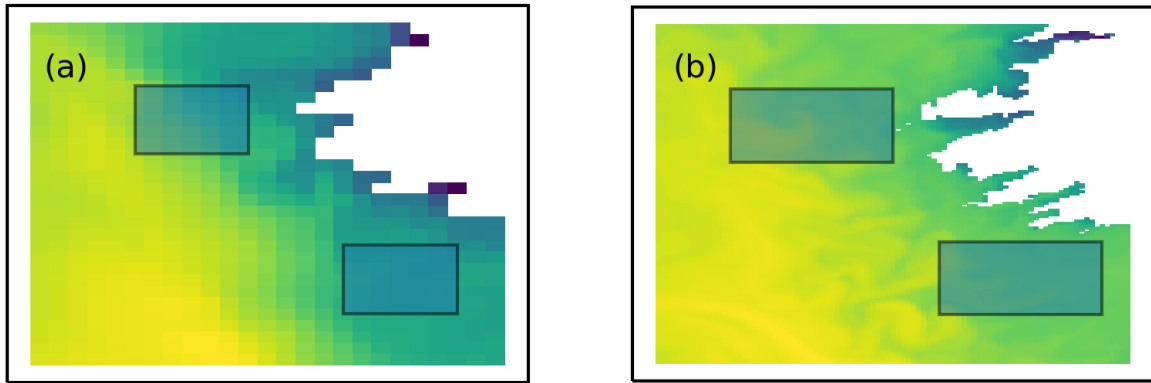
333 AMM15 is originally run on a rotated latitude-longitude grid where the resolution is closely
334 approximated by 1.5 km and subsequently provided to the CMEMS catalogue on the standard
335 latitude-longitude grid. Once the grid has been transformed to a regular latitude-longitude grid
336 the 1.5 km nominal spatial resolution is not as accurate. This is particularly important when
337 neighbourhood sizes become larger, since any error in the approximation of the resolution will
338 become multiplied as the number of points being used increases.

339

340 Additionally, the two model configurations do not have the same aspect ratio of grid points.
341 AMM7 has a longitudinal resolution of $\sim 0.11^\circ$ and a latitudinal resolution of $\sim 0.066^\circ$ (a ratio of
342 3:5) whilst the AMM15 grid has a resolution of $\sim 0.03^\circ$ and $\sim 0.0135^\circ$ respectively (a ratio of 5:11).
343 HiRA neighbourhoods typically contain the same number of grid-points in the zonal and
344 meridional directions which will lead to discrepancies in the area selected when comparing
345 models with different grid aspect ratios, depending on whether the comparison is based on
346 neighbourhoods with a similar longitudinal or similar latitudinal size. This difference will scale as
347 the neighbourhood size increases as shown in Fig. 4 and Table 2. The onus is therefore on the
348 user to understand any difference in grid structure, and therefore within the HiRA
349 neighbourhoods, between models being compared and to allow for this when comparing
350 equivalent neighbourhoods.

351

352



353

354 *Figure 5 - Similar neighbourhood sizes for a 49 km neighbourhood using the approximate resolutions (7 km and 1.5 km) with a)*
 355 *AMM7 with a 7x7 neighbourhood, b) AMM15 with a 33x33 neighbourhood. Whilst the neighbourhoods are similar sizes in the*
 356 *latitudinal direction, the AMM15 neighbourhood is sampling a much larger area due to different scales in the longitudinal*
 357 *direction. This means that a comparison with a 25x25 AMM15 neighbourhood is more appropriate.*

358 *Table 2 - Details of equivalent neighbourhoods used when comparing AMM7 and AMM15.*

Name	AMM7				AMM15			
	Total Points	Shape	Size (E-W)		Total Points	Shape	Size (E-W)	
			Actual (°)	Nominal (km)			Actual (°)	Nominal (km)
NB1	1	1x1	0.11	7	25	5x5	0.15	7.5
NB2	9	3x3	0.33	21	121	11x11	0.33	16.5
NB3	25	5x5	0.55	35	361	19x19	0.57	28.5
NB4	49	7x7	0.77	49	625	25x25	0.76	37.5
NB5	81	9x9	0.99	63	1089	33x33	0.99	49.5

359

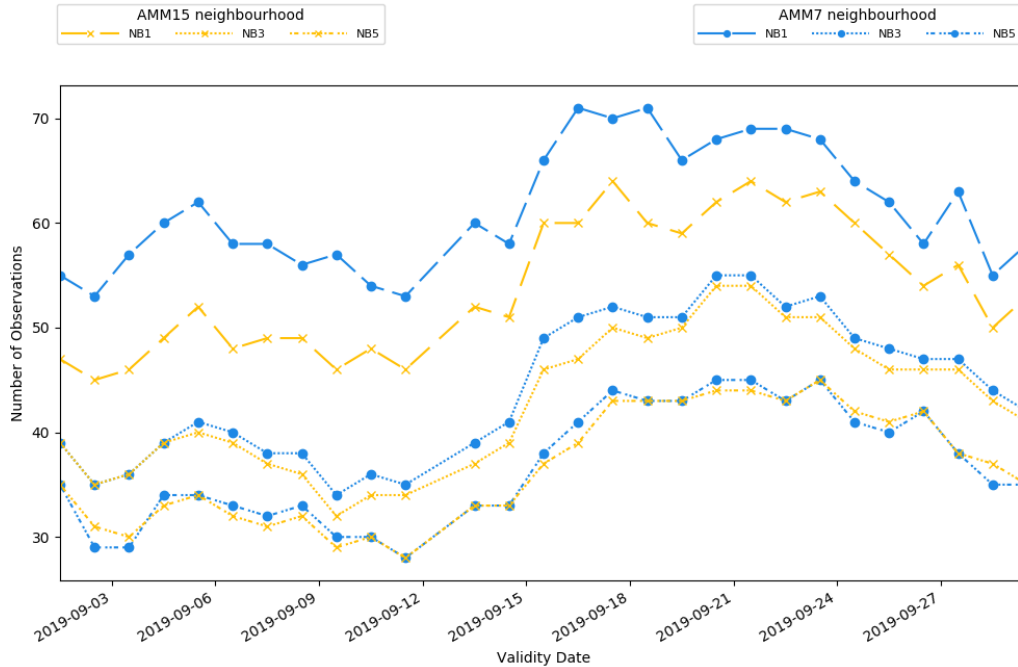
360 For this study we have matched neighbourhoods between model configurations based on their
 361 longitudinal size. The equivalent neighbourhoods used to show similar areas within the two
 362 configurations are indicated in Table 2 along with the bar style and naming convention used
 363 throughout.

364

365 For ocean applications there are other aspects of the processing to be aware of when using
 366 neighbourhood methods. This is mainly related to the presence of coastlines and how their
 367 representation changes resolution (as defined by the land-sea mask) and the treatment of

368 observations within HiRA neighbourhoods. Figure 5 illustrates the contrasting land-sea
369 boundaries due to the different resolutions of the two configurations. When calculating HiRA
370 neighbourhood values, all forecast values in the specific neighbourhood around an observation
371 must be present for a score to be calculated. If any forecast points within a neighbourhood
372 contain missing data then that observation at that neighbourhood size is rejected. This is to
373 ensure that the resolution of the “ensemble”, which is defined or determined by the number of
374 members, remains the same. For typical atmospheric fields such as screen temperature this is
375 not an issue, but with parameters that have physical boundaries (coastlines), such as SST, there
376 will be discontinuities in the forecast field that depend on the location of the land-sea boundary.
377 For coastal observations, this means that as the neighbourhood size increases, it is more likely
378 that an observation will be rejected from the comparison due to missing data. Even at the grid
379 scale, the nearest model grid point to an observation may not be a sea point. In addition, different
380 land-sea borders between models mean that potentially some observations will be rejected from
381 one model comparison but will be retained in the other because of missing forecast points within
382 their respective neighbourhoods. Care should be taken when implementing HiRA to check the
383 observations available to each model configuration when assessing the results and make a
384 judgement as to whether the differences are important.

385 There are potential ways to ensure equalisation, for example only using observations that are
386 available in both configurations for a location and neighborhoods, or only observations away
387 from the coast. For the purposes of this study, which aims to show the utility of the method, it
388 was judged important to use as many observations as possible, so as to capture any potential
389 pitfalls in the application of the framework, which would be relevant to any future application of
390 it.



391

392 *Figure 6- Number of observation sites within NB1, NB3 and NB5 for AMM15 and AMM7. Numbers are those used during*
 393 *September 2019 but represent typical total observations during a month. Matching line styles represent equivalent*
 394 *neighbourhoods.*

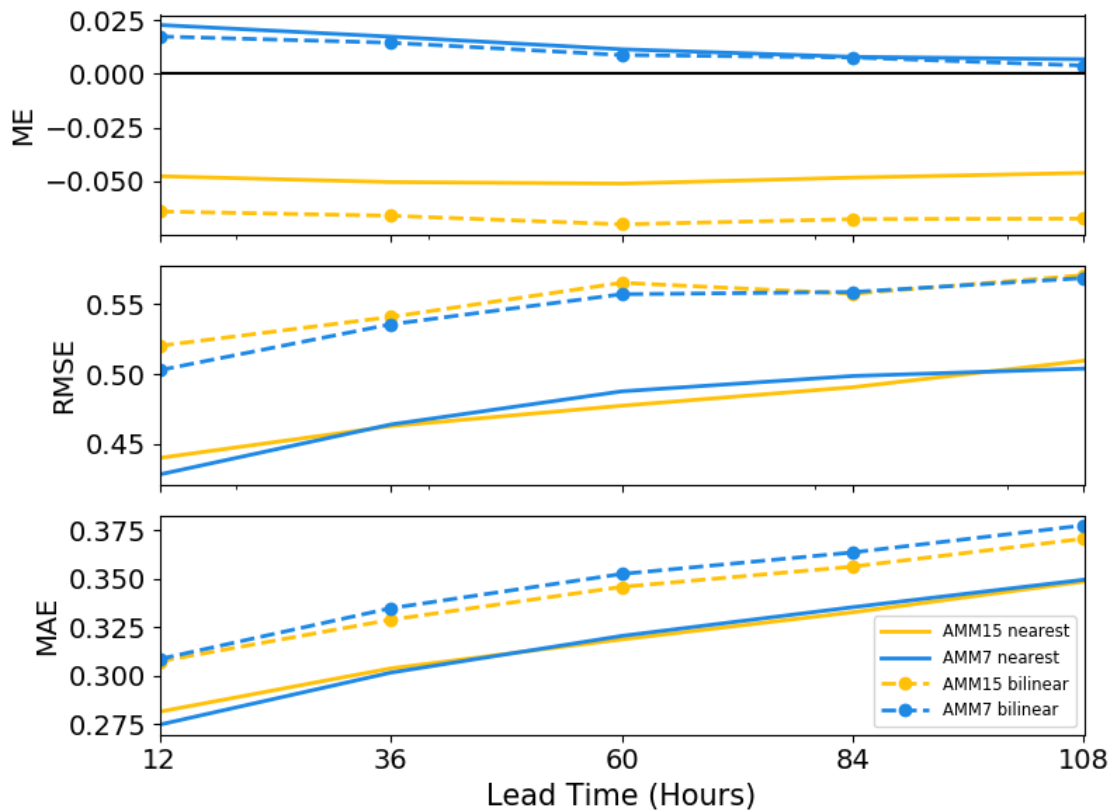
395

396 Figure 6 shows the number of observations available to each neighbourhood for each day during
 397 September 2019. For each model configuration it shows how these observations vary within the
 398 HiRA framework. There are several reasons for the differences shown in the plot. There is the
 399 difference mentioned previously whereby a model neighbourhood includes a land point, and
 400 therefore is rejected from the calculations because the number of quasi-ensemble members is
 401 no longer the same. This is more likely for coastal observations and depends on the particularities
 402 of the model land-sea mask near each observation. This rejection is more likely for the high-
 403 resolution AMM15 when looking at equivalent areas, in part due to the larger number of grid
 404 boxes being used; however, there are also instances of observations being rejected from the
 405 coarser resolution AMM7 and not the higher-resolution AMM15 due to nuances of the land-sea
 406 mask.

407 It is apparent that for equivalent neighbourhoods there are typically more observations available
 408 for the coarser model configuration and that this difference is largest for the smallest equivalent
 409 neighbourhood size but becoming less obvious at larger neighbourhoods. It could therefore be
 410 worth considering that the large benefit in AMM15 when looking at the first equivalent
 411 neighbourhood is potentially influenced by the difference in observations. As the neighbourhood
 412 sizes increase, the number of observations reduces due to the higher likelihood of a land point
 413 being part of a larger neighbourhood. It is also noted that there is a general daily variability in the
 414 number of observations present, based on differences in the observations reporting on any
 415 particular day within the co-located domain.

416

417 6. Results



418

419 *Figure 7 - Verification results using a typical statistics approach for January – September 2019. Mean error (top), root mean square*
 420 *error (middle) and mean absolute error (bottom) results are shown for the two model configurations. Two methods of matching*

421 *forecast to observations points have been used; a nearest neighbor approach (solid) representing the single grid point results from*
422 *HiRA, and a bilinear interpolation approach (dashed) more typically used in operational ocean verification.*

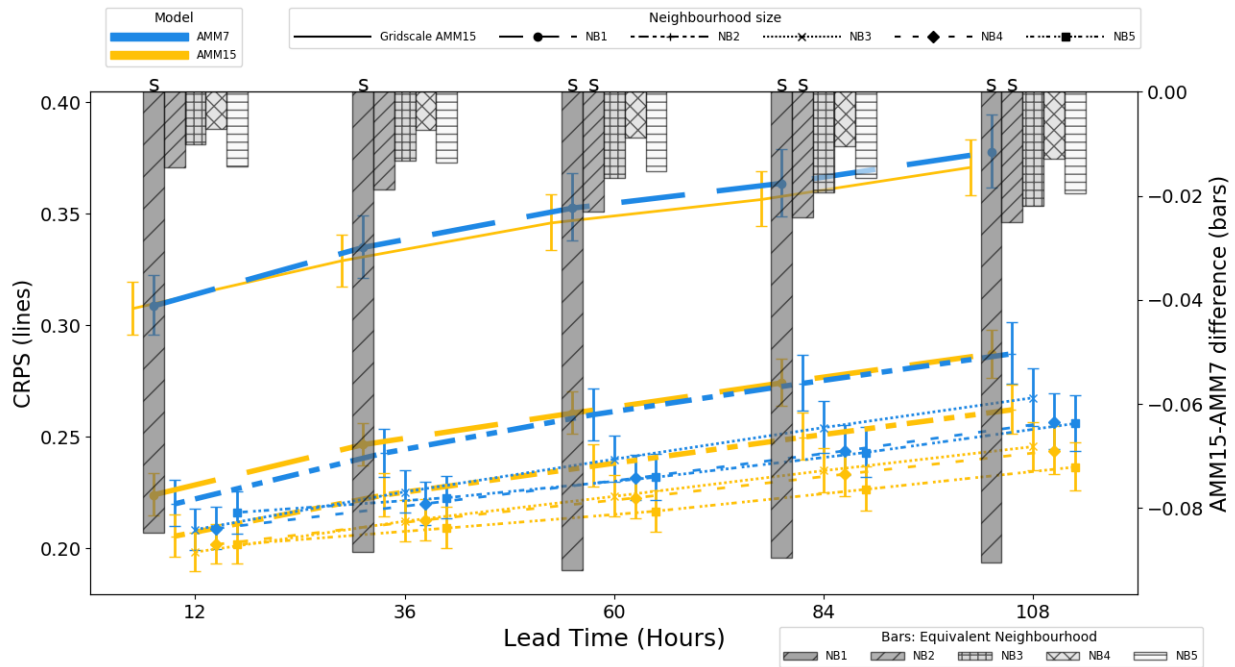
423 Figure 7 shows the aggregated results from the study period defined in Section 2 by applying
424 typical verification statistics. Results have been averaged across the entire period from January
425 to September and output relative to the forecast validity time. Two methods of matching forecast
426 grid points to observation locations have been used. Bilinear interpolation is typically the
427 approach used in traditional verification of SST, as it is a smoothly varying field. A nearest
428 neighbour approach has also been shown, as this is the method that would be used for HiRA
429 when applying it at the grid scale.

430 It is noted that the two methods of matching forecasts to observation locations give quite
431 different results. For the mean error, the impact of moving from a single grid point approach to
432 a bilinear interpolation method appears to be minor for the AMM7 model, but is more severe for
433 the AMM15, resulting in a larger error across all lead times. For the RMSE the picture is more
434 mixed, generally suggesting that the AMM7 forecasts are better when using a bilinear
435 interpolation method but giving no clear overall steer when the nearest grid point is used.
436 However, the impact of taking a bilinear approach results in much higher gross errors across all
437 lead times when compared to the nearest grid point approach.

438 The MAE has been suggested as a more appropriate metric than the RMSE for ocean fields using
439 (as is the case here) near real time observation data (Brassington, 2017). In Fig. 6 it can be seen
440 that the nearest grid point approach for both AMM7 and AMM15 gives almost exactly the same
441 results, except for the shortest of lead times. For the bilinear interpolation method, AMM15 has
442 a smaller error than AMM7 as lead time increases, behavior which is not apparent when RMSE is
443 applied.

444 Based on the interpolated RMSE results in Fig. 6 it would be hard to conclude that there was a
445 significant benefit to using high-resolution ocean models for forecasting SSTs. This is where the
446 HiRA framework can be applied. It can be used to provide more information, which can better
447 inform any conclusions on model error.

448



450
 451 *Figure 8- Summary of CRPS (left axis, lines) and CRPS difference (right axis, bars) for the period January 2019 to September 2019*
 452 *for AMM7 and AMM15 models at different neighbourhood sizes. Error bars represent 95 % confidence intervals generated using*
 453 *a bootstrap with replacement method for 10000 samples. An 'S' above the bar denotes that 95 % error bars for the two models*
 454 *do not overlap.*

455 Figure 8 shows the results for AMM7 and AMM15 for the period January - September 2019 using
 456 the HiRA framework with the CRPS. The lines on the plot show the CRPS for the two model
 457 configurations for different neighbourhood sizes, each plotted against lead-time. Similar line
 458 styles are used to represent equivalent neighbourhood sizes. Confidence intervals have been
 459 generated by applying a bootstrap with replacement method, using 10000 samples, to the
 460 domain-averaged CRPS (e.g. Efron and Tibshirani, 1993). The error bars represent the 95 %
 461 confidence level. The results for the single grid-point show the MAE and are the same as would
 462 be obtained using a traditional (precise) matching. In the case of CRPS, where a lower score is
 463 better, we see that AMM15 is better than AMM7, though not significantly so, except at shorter
 464 lead-times where there is little difference.

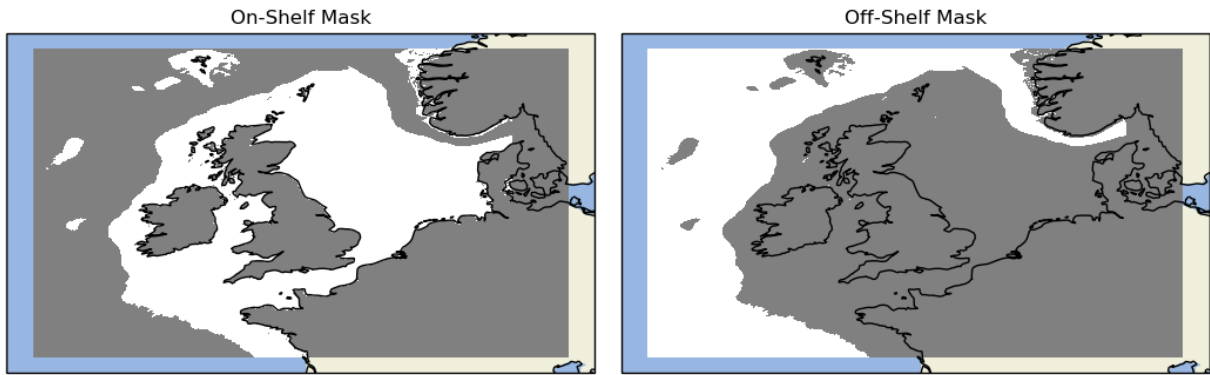
465 The differences at equivalent neighbourhood sizes are displayed as a bar plot on the same figure,
 466 with scores referenced with respect to the right-hand axis. Line markers and error bars have been
 467 offset to aid visualization, such that results for equivalent neighbourhoods are displayed in the

468 same vertical column as the difference indicated by the barplot. The details of the equivalent
469 neighbourhood sizes are presented in Table 2. Since a lower CRPS score is better, a positively
470 orientated (upwards) bar implies AMM7 is better, whilst a negatively orientated (downwards)
471 bar means AMM15 is better.

472 As indicated in Table 2, NB1 compares the single grid-point results of AMM7 with a 25-member
473 pseudo-ensemble constructed from a 5x5 AMM15 neighbourhood. Given the different
474 resolutions of the two configurations, these two neighbourhoods represent similar physical areas
475 from each model domain, with AMM7 only represented by a single forecast value for each
476 observation, but AMM15 represented by 25 values cover the same area, and as such potentially
477 better able to represent small-scale variability within that area.

478 At this equivalent scale the AMM15 results are markedly better than AMM7, with lower errors,
479 suggesting that overall the AMM15 neighbourhood better represents the variation around the
480 observation than the coarser single grid point of AMM7. At the next set of equivalent
481 neighbourhoods (NB2), the gap between the two configurations has closed, but AMM15 is still
482 consistently better than AMM7 as lead time increases. Above this scale the neighbourhood
483 values tend towards similarity, and then start to diverge again suggesting that the representative
484 scale of the neighbourhoods has been reached and that errors are essentially random.

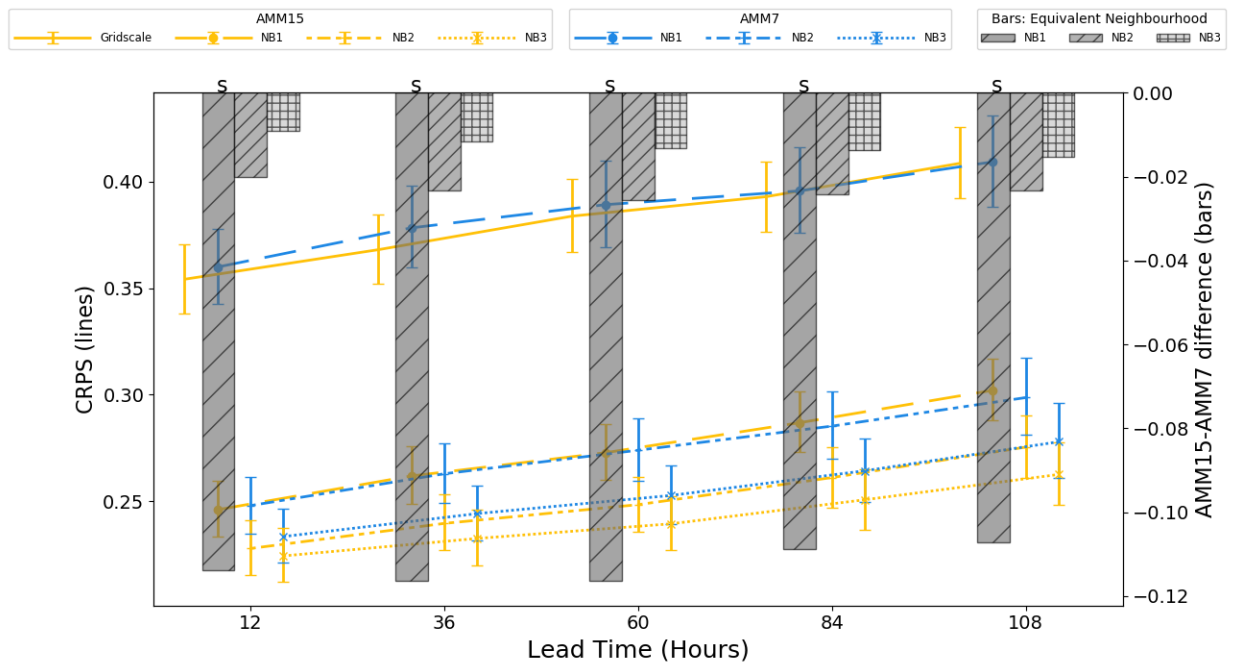
485 Whilst the overall HiRA neighbourhood results for the co-located domains appear to show a
486 benefit to using a higher resolution model forecast, it could be that these results are influenced
487 by the spatial distribution of observations within the domain and the characteristics of the
488 forecasts at those locations. In order to investigate whether this was important behaviour, the
489 results were separated into two domains, one representing the continental shelf part of the
490 domain (where the bathymetry < 200 m), and the other representing the deeper, off-shelf, ocean
491 component (Fig. 8). HiRA results were compared for observations only within each masked
492 domain.



493

494 *Figure 9 - On-shelf and off-shelf masking regions within the co-located AMM7 and AMM15 domain (data within the grey areas is*
 495 *masked).*

496



497

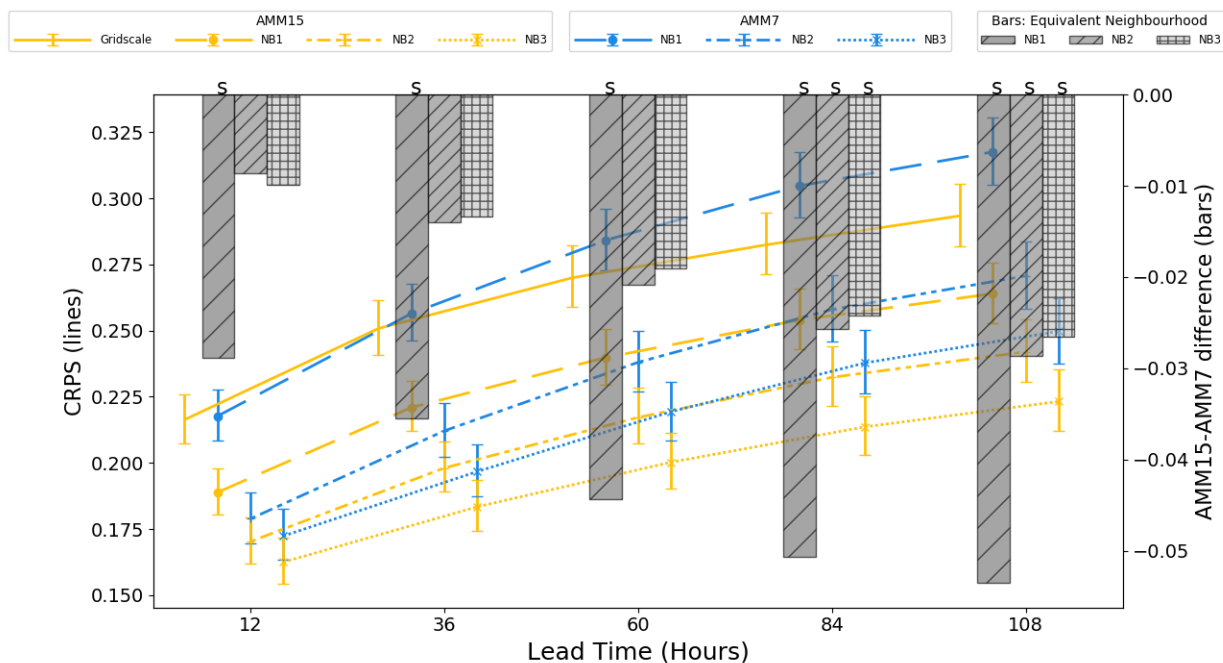
498 *Figure 10- Summary of on-shelf CRPS (left axis, lines) and CRPS difference (right axis, bars) for the period January 2019 to*
 499 *September 2019 for AMM7 and AMM15 models at different neighbourhood sizes. Error bars represent 95 % confidence values*
 500 *obtained from 10000 samples using bootstrap with replacement. An 'S' above the bar denotes that 95 % error bars for the two*
 501 *models do not overlap.*

502 On-shelf results (Fig. 10) show that at the grid scale the results for both AMM7 and AMM15 are
 503 worse for this sub-domain. This could be explained by both the complexity of processes (tides,
 504 friction, river mixing, topographical effects, etc.), and the small dynamical scales associated with
 505 shallow waters on the shelf (Holt et al., 2017).

506
 507 The on-shelf spatial variability in SST across a neighbourhood is likely to be higher than for an
 508 equivalent deep ocean neighbourhood due to small-scale changes in bathymetry, and for some
 509 observations, the impact of coastal effects. Both AMM7 and AMM15 show improvement in CRPS
 510 with increased neighbourhood size until the CRPS plateaus in the range 0.225 to 0.25, with
 511 AMM15 generally better than AMM7 for equivalent neighbourhood sizes. Scores get worse
 512 (errors increase) for both model configurations as the forecast lead time increases.

513

514



515
 516 *Figure 11 – Summary of off-shelf CRPS (left axis, lines) and CRPS difference (right axis, bars) for the period January 2019 to*
 517 *September 2019 for AMM7 and AMM15 models at different neighbourhood sizes. Error bars represent 95 % confidence values*
 518 *obtained from 10000 samples using bootstrap with replacement. An ‘S’ above the bar denotes that 95 % error bars for the two*
 519 *models do not overlap.*

520

521 For off-shelf results (Fig. 11), the CRPS is much better (smaller error), at both the grid scale and
522 for HiRA neighbourhoods, suggesting that both configurations are better at forecasting these
523 deep ocean SSTs (or that it is easier to do so). There is still an improvement in CRPS when going
524 from the grid scale (single grid box) to neighbourhoods, but the value of that change is much
525 smaller than for the on-shelf sub-domain. When comparing equivalent neighbourhoods, the
526 AMM15 still gives consistently better results (smaller errors) and appears to improve over AMM7
527 as lead time increases in contrast to the on-shelf results.

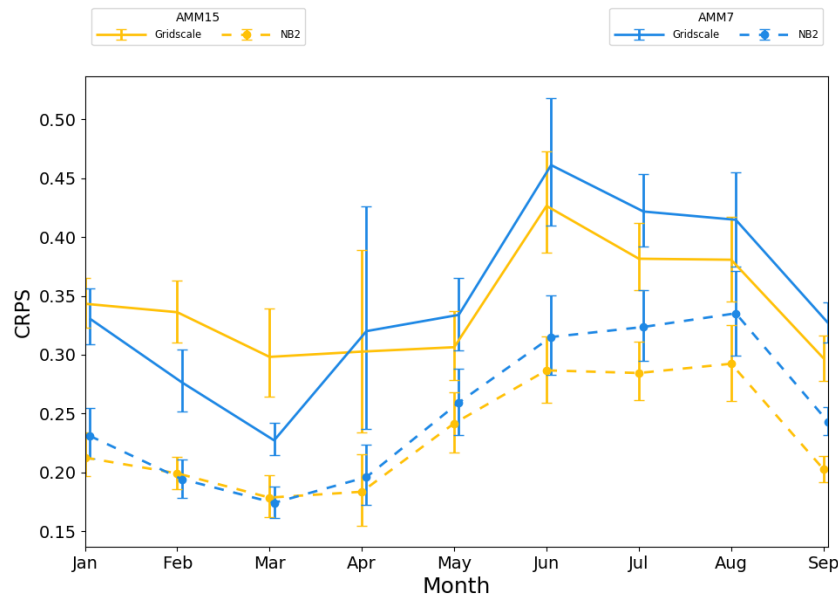
528 It is likely that the neighbourhood at which we lose representativity will be larger for the deeper
529 ocean than the shelf area because of the larger scale of dynamical processes in deep water. When
530 choosing an optimum neighbourhood to use for assessment, care should be taken to check
531 whether there are different representativity levels in the data (such as here for on-shelf and off-
532 shelf) and pragmatically choose the smaller of those equivalent neighbourhoods when looking at
533 data combining the different representativity levels.

534 Overall, for the period January-September 2019, the AMM15 demonstrates a lower (better) CRPS
535 than AMM7 when looking at the HiRA neighbourhoods. However, this also appears to be true at
536 the grid scale over the assessment period. One of the aspects that HiRA is trying to provide
537 additional information about is whether higher resolution models can demonstrate improvement
538 over coarser models against a perception that the coarser models score better in standard
539 verification forecast assessments. Assessed over the whole period, this initial premise does not
540 appear to hold true, therefore a deeper look at the data is required to assess whether this signal
541 is consistent within shorter time periods, or whether there are underlying periods contributing
542 significant and contrasting results to the whole-period aggregate.

543 Figure 12 shows a monthly breakdown of the grid scale and the NB2 HiRA neighbourhood scores
544 at T+60. This shows the underlying monthly variability not immediately apparent in the whole-
545 period plots. Notably for the January to March period, AMM7 outperforms AMM15 at the grid
546 scale. With the introduction of HiRA neighbourhoods, AMM7 still performs better for February
547 and March but the difference between the models is significantly reduced. For these monthly

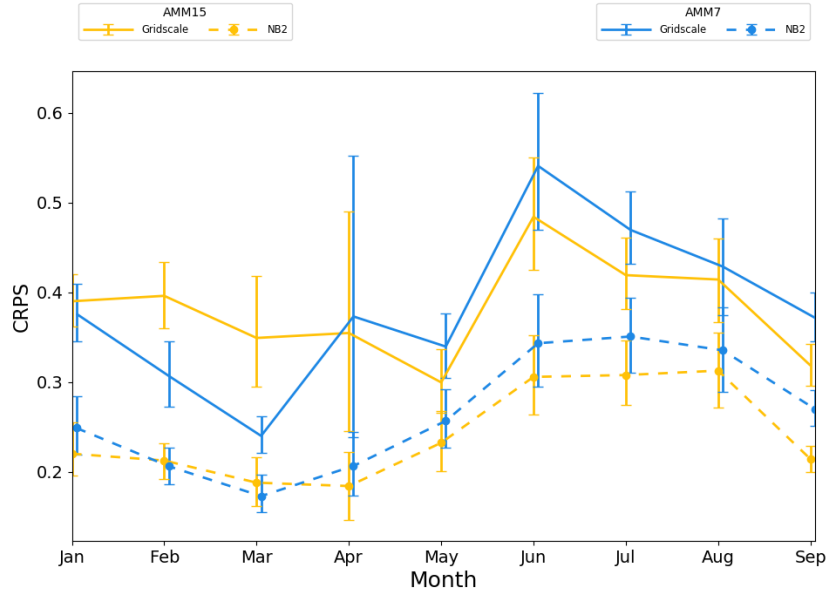
548 timeseries the error bars increase in size relative to the summary plots (e.g. Fig 8) due to the
 549 reduction in data available. The sample size will have an impact on the error bars as the smaller
 550 the sample, the less representative of the true population the data is likely to be. April in
 551 particular contained several days of missing forecast data, leading to a reduction in sample size
 552 and corresponding increase in error bar size, whilst during May there was a period with reduced
 553 numbers of observations.

554



555

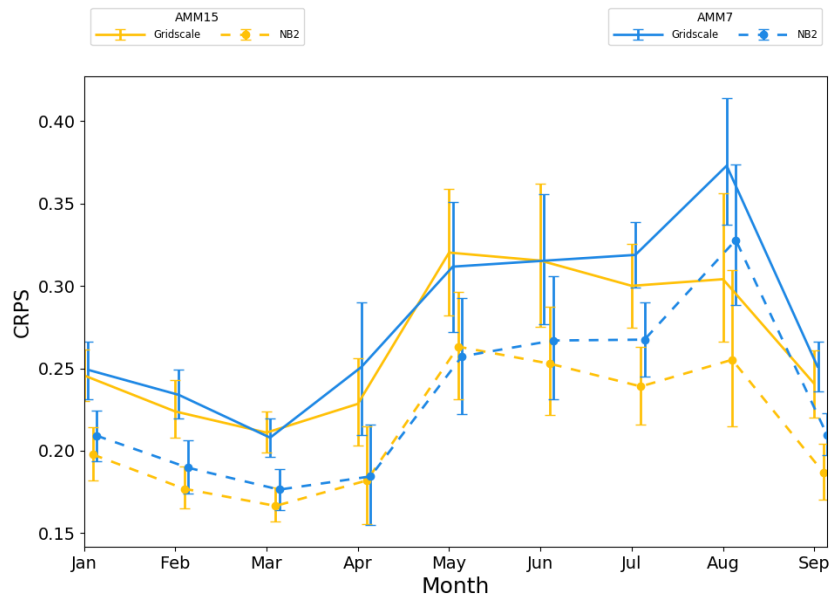
556 *Figure 12 – Monthly time series of whole-domain CRPS scores for grid scale (solid line) and NB2 neighbourhood (dashes) for T+60*
 557 *forecasts. Error bars represent 95 % confidence values obtained from 10000 samples using bootstrap with replacement. Error bars*
 558 *have been staggered in the x-direction to aid clarity.*



559

560 *Figure 13 - On-shelf monthly time series of CRPS. Error bars represent 95 % confidence values obtained from 10000 samples using*
 561 *bootstrap with replacement.*

562



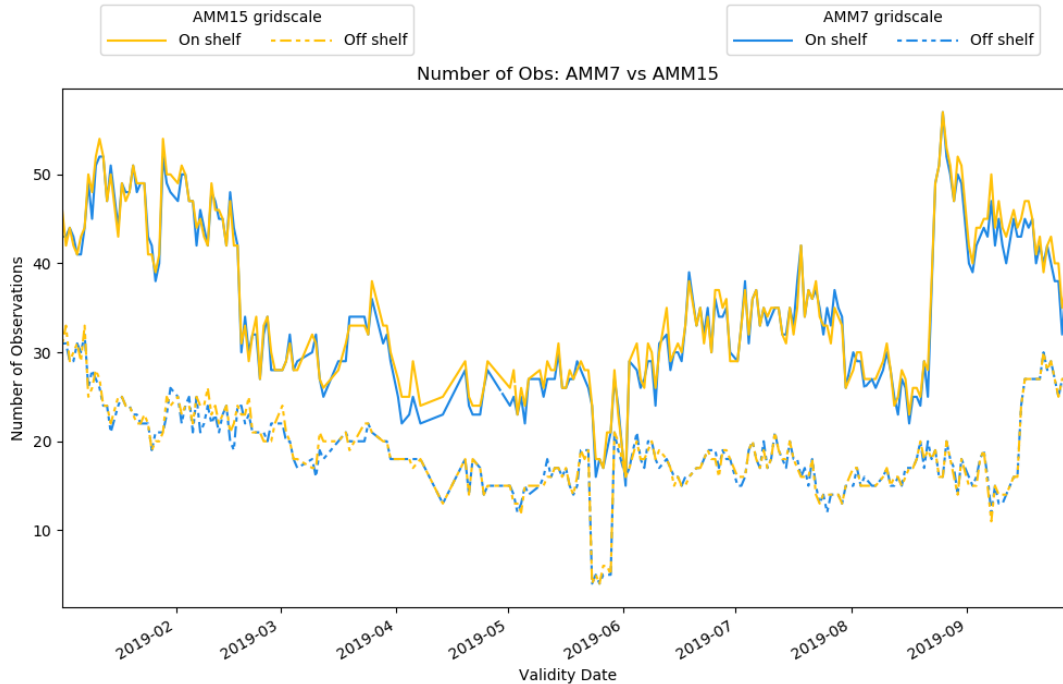
563

564 *Figure 14 - Off-shelf monthly time series of CRPS. Error bars represent 95 % confidence values obtained from 10000 samples using*
 565 *bootstrap with replacement.*

566

567 The same pattern is present for the on-shelf sub-domain (Fig. 13), where what appears to be a
568 significant benefit for the AMM7 during February and March is less clear-cut at the NB2
569 neighbourhood. For the off-shelf sub-domain (Fig. 14), differences between the two
570 configurations at the grid scale are mainly apparent during the summer months. At the NB2 scale,
571 the AMM15 potentially demonstrates more benefit than AMM7 except for April and May, where
572 the two show similar results. There is a balance to be struck in this conclusion as the differences
573 between the two models are rarely greater than the 95 % error bars. This in itself does not mean
574 that the results are not significant. However, care should be taken when interpreting such a result
575 as a statistical conclusion rather than broad guidance as to model performance. Attempts to
576 reduce the error bar size, such as increasing the number of observations, or number of times
577 within the period would aid this interpretation.

578 One noticeable aspect of the time series plots is that the whole-domain plot is heavily influenced
579 by the on-shelf results. This is due to the difference in observation numbers as shown in Fig. 15,
580 with the on-shelf domain having more observations overall, sometimes significantly more, for
581 example during January or mid-late August. For the overall domain, the on-shelf observations
582 will contribute more to the overall score and hence the underlying off-shelf signal will tend to be
583 masked. This is an indication of why verification is more useful when done over smaller, more
584 homogeneous sub-regions, rather than verifying everything together, with the caveat that
585 sample sizes are large enough, since underlying signals can be swamped by dominant error types.



586

587 *Figure 15 - Number of grid scale observations for the on and off-shelf domains.*

588

589 7. Discussion and Conclusions

590 In this study, the HiRA framework has been applied to SST forecasts from two ocean models with
 591 different resolutions. This enables a different view of the forecast errors than obtained using
 592 traditional (precise) grid scale matching against ocean observations. Particularly it enables us to
 593 demonstrate the additional value of high-resolution model. When considered more
 594 appropriately high-resolution models (with the ability to forecast small-scale detail) have lower
 595 errors when compared to the smoother forecasts provided by a coarser-resolution model.

596 The HiRA framework was intended to address the question ‘Does moving to higher resolution
 597 add value?’ This study has identified and highlighted aspects that need to be considered when
 598 setting up such an assessment. Prior to this study, routine verification statistics typically showed
 599 that coarser resolution models had equivalent or more skill than higher resolution models (e.g.
 600 Mass et al., 2002, Tonani et al., 2019). During the period January to September 2019, grid scale

601 verification within this assessment showed that the coarser-resolution AMM7 often
602 demonstrated lower errors than the AMM15.

603 HiRA neighbourhoods were applied and the data then assessed using the CRPS, showing a large
604 reduction (improvement) in errors for AMM15 when going from a grid scale, point-based
605 verification assessment to a neighbourhood, ensemble approach. When applying an equivalent-
606 sized neighbourhood to both configurations, AMM15 typically demonstrated lower (better)
607 scores. These scores were in turn broken down into off-shelf and on-shelf sub-domains and
608 showed that the different physical processes in these areas affected the results. Forecast
609 verification studies tailored for the coastal/shelf areas are needed to properly understand the
610 forecast skills in areas with high complexity and fast evolving dynamics.

611 When constructing HiRA neighbourhoods the spatial scales that are appropriate for the
612 parameter must be considered carefully. This often means running at several neighbourhood
613 sizes and determining where the scores no longer seem physically representative. When
614 comparing models, care should be taken to construct neighbourhood sizes that are similarly sized
615 spatially, the details of the neighbourhood sizes will depend on the structure and resolution of
616 the model grid.

617 Treatment of observations is also important in any verification set-up. For this study, the fact that
618 there are different numbers of observations present at each neighbourhood scale (as
619 observations are rejected due to land contamination) means that there is never an optimally
620 equalized data set (i.e. the same observations for all models and for all neighbourhood sizes). It
621 also means that comparison of the different neighbourhood results from a single model is ill
622 advised, in this case, as the observations numbers can be very different, and therefore the model
623 forecast is being sampled at different locations. Despite this, observation numbers should be
624 similar when looking at matched spatially sized neighbourhoods from different models if results
625 are to be compared. One of the main constraints identified through this work is both the sparsity
626 and geographical distribution of observations throughout the North West Shelf domain, with
627 several viable locations rejected during the HiRA processing due to their proximity to coastlines.

628 The purest assessment, in terms of observations, would involve a fixed set of observations,
629 equalized across both model configurations and all neighbourhoods at every time. This would
630 remove the variation in observation numbers seen as neighbourhood sizes increase as well as
631 those seen between the two models and give a clean comparison between two models.

632 Care should be taken when applying strict equalization rules as this could result in only a small
633 number of observations being used. The total number of observations used should be large
634 enough to ensure that the sample is large enough to produce robust results and satisfy rules for
635 statistical significance. Equalisation rules could also unfairly affect the spatial sampling of the
636 verification domain. For example, in this study coastal observations would be affected more than
637 deep ocean observations if neighbourhood equalization were applied, due to the proximity of
638 the coast.

639 To a lesser extent, the variation in observation numbers on a day-to-day timescale also has an
640 impact on any results and could mean that incorrect importance is attributed to certain results,
641 which are simply due to fluctuations in observation numbers.

642 The fact that the errors can be reduced through the use of neighbourhoods shows that the ocean
643 and the atmosphere have similarities in the way the forecasts behave as a function of resolution.
644 This study did not consider the concept of skill, which incorporates the performance of the
645 forecast relative to a pre-defined benchmark. For the ocean the choice of reference needs to be
646 considered. This could be the subject of further work.

647 To our knowledge, this work is the first attempt to use neighbourhood techniques to assess ocean
648 models. The promising results showing reductions in errors of the finer resolution configuration
649 warrant further work. We see a number of directions the current study could be extended.

650 The study was conducted on daily output which should be appropriate to address eddy mesoscale
651 variability, but observations are distributed at hourly resolution, and so the next logical step
652 would be to assess the hourly forecasts against the hourly observation and see how this impacted
653 the results. This will increase the sample size, if all hourly observations were considered together.
654 However, it is impossible to speculate on whether considering hourly forecasts would lead to
655 more noisy statistics, counteracting the larger sample size.

656 This assessment only looked at SST for this initial examination. Consideration of other ocean
657 variables would also be of interest, including looking at derived diagnostics such as mixed layer
658 depth, but the sparsity of observations available for some variables may limit the case studies
659 available. HiRA as a framework is not remaining static. Enhancements to introduce non-regular
660 flow-dependent neighbourhoods are planned and may be of benefit to ocean applications in the
661 future. Finally, an advantage of using the HiRA framework is that results obtained from
662 deterministic ocean models could also be compared against results from ensemble models when
663 these become available for ocean applications.

664

665 8. References

666 Aznar, R., Sotillo, M., Cailleau, S., Lorente, P., Levier, B., Amo-Baladrón, A., Reffray, G. and Alvarez Fanjul,
667 E.: Strengths and weaknesses of the CMEMS forecasted and reanalyzed solutions for the Iberia-Biscay-
668 Ireland (IBI) waters. *J. Marine. Syst.*, 159, <https://doi.org/10.1016/j.jmarsys.2016.02.007>, 2016.

669 Brassington, G.: Forecast Errors, Goodness, and Verification in Ocean Forecasting, *J. Marine Res.*, 75,
670 403-433, <https://doi.org/10.1357/002224017821836851>, 2017.

671 Brier, G. W.: Verification of Forecasts Expressed in Terms of Probability, *Mon. Wea. Rev.*, 78, 1-3,
672 [http://dx.doi.org/10.1175/1520-0493\(1950\)078<0001:VOFEIT>2.0.CO;2](http://dx.doi.org/10.1175/1520-0493(1950)078<0001:VOFEIT>2.0.CO;2), 1950.

673 Brown, T. A.: Admissible scoring systems for continuous distributions, Santa Monica, CA: RAND
674 Corporation, available at <https://www.rand.org/pubs/papers/P5235.html>, 1974.

675 Casati, B., Ross, G. and Stephenson, D. B.: A new intensity-scale approach for the verification of spatial
676 precipitation forecasts, *Met. Apps.*, 11, 141-154, <https://doi.org/10.1017/S1350482704001239>, 2004.

677 Davis, C., Brown, B. and Bullock, R.: Object-Based Verification of Precipitation Forecasts. Part I:
678 Methodology and Application to Mesoscale Rain Areas, *Mon. Wea. Rev.*, **134**, 1772–1784,
679 <https://doi.org/10.1175/MWR3145.1>, 2006.

680 Dorninger, M., Gilleland, E., Casati, B., Mittermaier, M. P., Ebert, E. E., Brown, B. G., and Wilson, L. J.: The
681 Setup of the MesoVICT Project, *Bull. Amer. Meteor. Soc.*, 99, 1887–1906,
682 <https://doi.org/10.1175/BAMS-D-17-0164.1>, 2008.

683 Ebert, E. E.: Fuzzy verification of high-resolution gridded forecasts: a review and proposed framework,
684 Met. Apps, 15, 51-64, <https://doi.org/10.1002/met.25>, 2008.

685 Efron, B. and Tibshirani, R.: Bootstrap methods for standard errors, confidence intervals, and other
686 measures of statistical accuracy, Statistical Science, 1, 54-77, 1986.

687 Epstein, E. S.: A Scoring System for Probability Forecasts of Ranked Categories, J. Appl. Meteor., 8, 985–
688 987, 1969.

689 Gilleland, E., Ahijevych, D., Brown, B. G., Casati, B., and Ebert, E. E.: Intercomparison of Spatial Forecast
690 Verification Methods, Wea. Forecasting, 24, 1416–1430, <https://doi.org/10.1175/2009WAF2222269.1>,
691 2009.

692 Goessling, H. F., Tietsche, S., Day, J. J., Hawkins, E., and Jung, T. : Predictability of the Arctic sea ice edge,
693 Geophys. Res. Lett., 43, 1642– 1650, doi:[10.1002/2015GL067232](https://doi.org/10.1002/2015GL067232), 2016.

694 Graham, J. A., O'Dea, E., Holt, J., Polton, J., Hewitt, H. T., Furner, R., Guihou, K., Brereton, A., Arnold, A.,
695 Wakelin, S., Castillo Sanchez, J. M., and Mayorga Adame, C. G.: AMM15: a new high-resolution NEMO
696 configuration for operational simulation of the European north-west shelf, Geosci. Model Dev., 11, 681–
697 696, <https://doi.org/10.5194/gmd-11-681-2018>, 2018.

698 Hernandez, F., Blockley, E., Brassington, G. B., Davidson, F., Divakaran, P., Drévillon, M., Ishizaki, S.,
699 Garcia-Sotillo, M., Hogan, P. J., Lagemaat, P., Levier, B., Martin, M., Mehra, A., Mooers, C., Ferry, N.,
700 Ryan, A., Regnier, C., Sellar, A., Smith, G. C., Sofianos, S., Spindler, T., Volpe, G., Wilkin, J., Zaron, E. D.,
701 and Zhang, A.: Recent progress in performance evaluations and near real-time assessment of
702 operational ocean products, J. Oper. Oceanogr., 8, 221–238,
703 <https://doi.org/10.1080/1755876X.2015.1050282>, 2015.

704 Hernandez, F., Smith, G., Baetens, K., Cossarini, G., Garcia-Hermosa, I., Drevillon, M., Maksymczuk, J.,
705 Melet, A., Regnier, C., and von Schuckmann, K.: Measuring Performances, Skill and Accuracy in
706 Operational Oceanography: New Challenges and Approaches. In "New Frontiers in Operational
707 Oceanography", E. Chassignet, A. Pascual, J. Tintore, and J. Verron, Eds., GODAE OceanView, 759-796,
708 doi:10.17125/gov2018.ch29, 2018.

709 Hersbach, H.: Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction
710 Systems, Wea. Forecasting, 15, 559–570, [https://doi.org/10.1175/1520-
711 0434\(2000\)015<0559:DOTCRP>2.0.CO;2](https://doi.org/10.1175/1520-0434(2000)015<0559:DOTCRP>2.0.CO;2), 2000.

712 Holt, J., Hyder, P., Ashworth, M., Harle, J., Hewitt, H. T., Liu, H., New, A. L., Pickles, S., Porter, A., Popova,
713 E. and Allen, J.: Prospects for improving the representation of coastal and shelf seas in global ocean
714 models, *Geosci. Model Dev.*, 10, 499-523, 2017.

715 Howarth, M. and Pugh, D.: Chapter 4 Observations of Tides Over the Continental Shelf of North-West
716 Europe, Elsevier Oceanography Series, 35, 135–188, [https://doi.org/10.1016/S04229894\(08\)70502-6](https://doi.org/10.1016/S04229894(08)70502-6),
717 1983

718 Juza, M., Mourre, B., Lellouche, J. M., Tonani M., and Tintoré, J.: From basin to sub-basin scale
719 assessment and intercomparison of numerical simulations in the western Mediterranean Sea
720 *J. Mar. Syst.*, 149, 36-49, <https://doi.org/10.1016/j.jmarsys.2015.04.010>, 2015.

721 Keil, C. and Craig, G. C.: A Displacement-Based Error Measure Applied in a Regional Ensemble
722 Forecasting System, *Mon. Wea. Rev.*, 135, 3248–3259, <https://doi.org/10.1175/MWR3457.1>, 2007.

723 King, R., While, J., Martin, M. J., Lea, D. J., Lemieux-Dudon, B, Waters, J., O’Dea, E.: Improving the
724 initialisation of the Met Office operational shelf-seas model. *Ocean Model*, 130, 1-14, 2018.

725 Le Traon PY, Reppucci A, Alvarez Fanjul E, Aouf L, Behrens A, Belmonte M, Bentamy A, Bertino L, Brando
726 VE, Kreiner MB, Benkiran M, Carval T, Ciliberti SA, Claustre H, Clementi E, Coppini G, Cossarini G, De
727 Alfonso Alonso-Muñoyerro M, Delamarche A, Dibarboure G, Dinessen F, Drevillon M, Drillet Y, Faugere
728 Y, Fernández V, Fleming A, Garcia-Hermosa MI, Sotillo MG, Garric G, Gasparin F, Giordan C, Gehlen M,
729 Gregoire ML, Guinehut S, Hamon M, Harris C, Hernandez F, Hinkler JB, Hoyer J, Karvonen J, Kay S, King R,
730 Lavergne T, Lemieux-Dudon B, Lima L, Mao C, Martin MJ, Masina S, Melet A, Buongiorno Nardelli B,
731 Nolan G, Pascual A, Pistoia J, Palazov A, Piolle JF, Pujol MI, Pequignet AC, Peneva E, Pérez Gómez B, Petit
732 de la Villeon L, Pinardi N, Pisano A, Pouliquen S, Reid R, Remy E, Santoleri R, Siddorn J, She J, Staneva J,
733 Stoffelen A, Tonani M, Vandenbulcke L, von Schuckmann K, Volpe G, Wettre C and Zacharioudaki A:
734 From Observation to Information and Users: The Copernicus Marine Service Perspective. *Front. Mar. Sci.*
735 6:234. doi: 10.3389/fmars.2019.00234, 2019.

736 Lorente, P., Sotillo, M., Amo-Baladrón, A., Aznar, R., Levier, B., Aouf, L., Dabrowski, T., Pascual, Á.,
737 Reffray, G., Dalphinnet, A., Toledano Lozano, C., Rainaud, R., and Alvarez Fanjul, E. : The NARVAL Software
738 Toolbox in Support of Ocean Models Skill Assessment at Regional and Coastal Scales.
739 http://doi.org/10.1007/978-3-030-22747-0_25 2019a.

740 Lorente, P., García-Sotillo, M., Amo-Baladrón, A., Aznar, R., Levier, B., Sánchez-Garrido, J. C.,
741 Sammartino, S., de Pascual-Collar, Á., Reffray, G., Toledano, C., and Álvarez-Fanjul, E.: Skill assessment of
742 global, regional, and coastal circulation forecast models: evaluating the benefits of dynamical
743 downscaling in IBI (Iberia–Biscay–Ireland) surface waters, *Ocean Sci.*, 15, 967–996,
744 <https://doi.org/10.5194/os-15-967-2019>, 2019b.

745 Madec, G. and the NEMO team: NEMO ocean engine. Note du Pôle de modélisation, Institut Pierre-
746 Simon Laplace (IPSL), France, No 27 ISSN No 1288-1619, 2016.

747 Mason, E., Ruiz, S., Bourdalle-Badie, R., Reffray, G., García-Sotillo, M., and Pascual, A.: New insight into
748 3-D mesoscale eddy properties from CMEMS operational models in the western Mediterranean, *Ocean*
749 *Sci.*, 15, 1111–1131, <https://doi.org/10.5194/os-15-1111-2019>, 2019.

750 Mass, C. F., Ovens, D., Westrick, K., and Colle, B. A.: DOES INCREASING HORIZONTAL RESOLUTION
751 PRODUCE MORE SKILLFUL FORECASTS?, *Bull. Amer. Meteor. Soc.*, 83, 407–430,
752 [https://doi.org/10.1175/1520-0477\(2002\)083<0407:DIHRPM>2.3.CO;2](https://doi.org/10.1175/1520-0477(2002)083<0407:DIHRPM>2.3.CO;2), 2002.

753 Mirouze, I., Blockley, E. W., Lea, D. J., Martin, M. J. and Bell, M. J.: A multiple length scale correlation
754 operator for ocean data assimilation, *Tellus A*, 68, 29744, <https://doi.org/10.3402/tellusa.v68.29744>,
755 2016.

756 Mittermaier, M., Roberts, N., and Thompson, S. A.: A long-term assessment of precipitation forecast skill
757 using the Fractions Skill Score, *Met. Apps*, 20, 176-186, <https://doi.org/10.1002/met.296>, 2013.

758 Mittermaier, M. P.: A Strategy for Verifying Near-Convection-Resolving Model Forecasts at Observing
759 Sites, *Wea. Forecasting*, 29, 185–204, <https://doi.org/10.1175/WAF-D-12-00075.1>, 2014.

760 Mittermaier, M. P., and Csima, G.: Ensemble versus Deterministic Performance at the Kilometer Scale,
761 *Wea. Forecasting*, 32, 1697–1709, <https://doi.org/10.1175/WAF-D-16-0164.1>, 2017.

762 Mogensen, K, Balmaseda, M. A., Weaver, A.: The NEMOVAR ocean data assimilation system as
763 implemented in the ECMWF ocean analysis for System 4. *European Centre for Medium-Range Weather*
764 *Forecasts*, 2012.

765 Mourre B., E. Aguiar, M. Juza, J. Hernandez-Lasheras, E. Reyes, E. Heslop, R. Escudier, E. Cutolo, S. Ruiz,
766 E. Mason, A. Pascual and J. Tintoré: Assessment of high-resolution regional ocean prediction systems
767 using multi-platform observations: illustrations in the Western Mediterranean Sea. In “New Frontiers in
768 Operational Oceanography”, E. Chassignet, A. Pascual, J. Tintoré and J. Verron, Eds, GODAE Ocean View,
769 663-694, doi: 10.17125/gov2018.ch24, 2018.

770 O'Dea, E. J., Arnold, A. K., Edwards, K. P., Furner, R., Hyder, P., Martin, M. J., Siddorn, J. R., Storkey, D.,
771 While, J., Holt, J. T., and Liu, H.: An operational ocean forecast system incorporating NEMO and SST data
772 assimilation for the tidally driven European North-West shelf, *J. Oper. Oceanogr.*, 5, 3–17,
773 <https://doi.org/10.1080/1755876X.2012.11020128>, 2012.

774 O'Dea, E., Furner, R., Wakelin, S., Siddorn, J., While, J., Sykes, P., King, R., Holt, J., and Hewitt, H.: The
775 CO5 configuration of the 7 km Atlantic Margin Model: large-scale biases and sensitivity to forcing,
776 physics options and vertical resolution, *Geosci. Model Dev.*, 10, 2947–2969,
777 <https://doi.org/10.5194/gmd-10-2947-2017>, 2017.

778 Rossa A., Nurmi P., Ebert E.: Overview of methods for the verification of quantitative precipitation
779 forecasts, in: *Precipitation: Advances in Measurement, Estimation and Prediction*, edited by:
780 Michaelides, S., Springer, Berlin, Heidelberg, 419–452, https://doi.org/10.1007/978-3-540-77655-0_16,
781 2008.

782 Tonani, M., Sykes, P., King, R. R., McConnell, N., Péquignat, A.-C., O'Dea, E., Graham, J. A., Polton, J., and
783 Siddorn, J.: The impact of a new high-resolution ocean model on the Met Office North-West European
784 Shelf forecasting system, *Ocean Sci.*, 15, 1133–1158, <https://doi.org/10.5194/os-15-1133-2019>, 2019.

785
786 World Meteorological Organisation: Guide to Meteorological Instruments and Methods of Observation
787 (WMO-No. 8, the CIMO Guide) –available at
788 https://library.wmo.int/opac/doc_num.php?explnum_id=4147, 2017.

789 9. Author contributions

790 All authors contributed to the introduction, data and methods, and conclusions. RC, JM and MM
791 contributed to the scientific evaluation and analysis of the results. RC and JM designed and ran

792 the model assessments. CP supported the assessments through the provision and reformatting
793 of the data used. MT provided detail on the model configurations used.

794

795 **10. Competing interests**

796 The authors declare that they have no conflict of interest.

797

798 **11. Acknowledgements**

799 This study has been conducted using E.U. Copernicus Marine Service Information.

800 This work has been carried out as part of the Copernicus Marine Environment Monitoring Service
801 (CMEMS) HiVE project. CMEMS is implemented by Mercator Ocean International in the
802 framework of a delegation agreement with the European Union.

803 Model Evaluation Tools (MET) was developed at the National Center for Atmospheric Research
804 (NCAR) through grants from the National Science Foundation (NSF), the National Oceanic and
805 Atmospheric Administration (NOAA), The United States Air Force (USAF), and the United States
806 Department of Energy (DOE). NCAR is sponsored by the United States National Science
807 Foundation.